



# cuBLAS Library

## User Guide



---

# Chapter 1. Introduction

The cuBLAS library is an implementation of BLAS (Basic Linear Algebra Subprograms) on top of the NVIDIA® CUDA® runtime. It allows the user to access the computational resources of NVIDIA Graphics Processing Unit (GPU).

The cuBLAS Library exposes three sets of API:

- ▶ The [cuBLAS API](#), which is simply called cuBLAS API in this document (starting with CUDA 6.0),
- ▶ The [cuBLASXt API](#) (starting with CUDA 6.0), and
- ▶ The [cuBLASLt API](#) (starting with CUDA 10.1)

To use the cuBLAS API, the application must allocate the required matrices and vectors in the GPU memory space, fill them with data, call the sequence of desired cuBLAS functions, and then upload the results from the GPU memory space back to the host. The cuBLAS API also provides helper functions for writing and retrieving data from the GPU.

To use the cuBLASXt API, the application may have the data on the Host or any of the devices involved in the computation, and the Library will take care of dispatching the operation to, and transferring the data to, one or multiple GPUs present in the system, depending on the user request.

The cuBLASLt is a lightweight library dedicated to GEneral Matrix-to-matrix Multiply (GEMM) operations with a new flexible API. This library adds flexibility in matrix data layouts, input types, compute types, and also in choosing the algorithmic implementations and heuristics through parameter programmability. After a set of options for the intended GEMM operation are identified by the user, these options can be used repeatedly for different inputs. This is analogous to how cuFFT and FFTW first create a plan and reuse for same size and type FFTs with different input data.

## 1.1. Data Layout

For maximum compatibility with existing Fortran environments, the cuBLAS library uses column-major storage, and 1-based indexing. Since C and C++ use row-major storage, applications written in these languages can not use the native array semantics for two-dimensional arrays. Instead, macros or inline functions should be defined to implement matrices on top of one-dimensional arrays. For Fortran code ported to C in mechanical fashion, one may chose to retain 1-based indexing to avoid the need to transform loops. In this

case, the array index of a matrix element in row “i” and column “j” can be computed via the following macro

```
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1)
```

Here, `ld` refers to the leading dimension of the matrix, which in the case of column-major storage is the number of rows of the allocated matrix (even if only a submatrix of it is being used). For natively written C and C++ code, one would most likely choose 0-based indexing, in which case the array index of a matrix element in row “i” and column “j” can be computed via the following macro

```
#define IDX2C(i,j,ld) ((j)*(ld))+i)
```

## 1.2. New and Legacy cuBLAS API

Starting with version 4.0, the cuBLAS Library provides a new API, in addition to the existing legacy API. This section discusses why a new API is provided, the advantages of using it, and the differences with the existing legacy API.



**WARNING: The legacy cuBLAS API is deprecated and will be removed in a future release.**

The new cuBLAS library API can be used by including the header file “`cublas_v2.h`”. It has the following features that the legacy cuBLAS API does not have:

- ▶ The `handle` to the cuBLAS library context is initialized using the function and is explicitly passed to every subsequent library function call. This allows the user to have more control over the library setup when using multiple host threads and multiple GPUs. This also allows the cuBLAS APIs to be reentrant.
- ▶ The scalars  $\alpha$  and  $\beta$  can be passed by reference on the host or the device, instead of only being allowed to be passed by value on the host. This change allows library functions to execute asynchronously using streams even when  $\alpha$  and  $\beta$  are generated by a previous kernel.
- ▶ When a library routine returns a scalar result, it can be returned by reference on the host or the device, instead of only being allowed to be returned by value only on the host. This change allows library routines to be called asynchronously when the scalar result is generated and returned by reference on the device resulting in maximum parallelism.
- ▶ The error status `cublasStatus_t` is returned by all cuBLAS library function calls. This change facilitates debugging and simplifies software development. Note that `cublasStatus` was renamed `cublasStatus_t` to be more consistent with other types in the cuBLAS library.
- ▶ The `cublasAlloc()` and `cublasFree()` functions have been deprecated. This change removes these unnecessary wrappers around `cudaMalloc()` and `cudaFree()`, respectively.
- ▶ The function `cublasSetKernelStream()` was renamed `cublasSetStream()` to be more consistent with the other CUDA libraries.

The legacy cuBLAS API, explained in more detail in the Appendix A, can be used by including the header file “cublas.h”. Since the legacy API is identical to the previously released cuBLAS library API, existing applications will work out of the box and automatically use this legacy API without any source code changes.

The current and the legacy cuBLAS APIs cannot be used simultaneously in a single translation unit: including both “cublas.h” and “cublas\_v2.h” header files will lead to compilation errors due to incompatible symbol redeclarations.

In general, new applications should not use the legacy cuBLAS API, and existing applications should convert to using the new API if it requires sophisticated and optimal stream parallelism, or if it calls cuBLAS routines concurrently from multiple threads.

For the rest of the document, the new cuBLAS Library API will simply be referred to as the cuBLAS Library API.

As mentioned earlier the interfaces to the legacy and the cuBLAS library APIs are the header file “cublas.h” and “cublas\_v2.h”, respectively. In addition, applications using the cuBLAS library need to link against:

- ▶ The DSO `cublas.so` for Linux,
- ▶ The DLL `cublas.dll` for Windows, or
- ▶ The dynamic library `cublas.dylib` for Mac OS X.



**Note:** The same dynamic library implements both the new and legacy cuBLAS APIs.

## 1.3. Example Code

The following code examples show an application written in C using the cuBLAS library API with two indexing styles. Example 1 shows 1-based indexing and Example 2 shows 0-based indexing.

```
//Example 1. Application Using C and cuBLAS: 1-based indexing
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"
#define M 6
#define N 5
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1))

static __inline__ void modify (cublasHandle_t handle, float *m, int ldm, int n, int
p, int q, float alpha, float beta){
    cublasSscal (handle, n-q+1, &alpha, &m[IDX2F(p,q,ldm)], ldm);
    cublasSscal (handle, ldm-p+1, &beta, &m[IDX2F(p,q,ldm)], 1);
}

int main (void){
    cudaError_t cudaStat;
    cublasStatus_t stat;
    cublasHandle_t handle;
```

```

int i, j;
float* devPtrA;
float* a = 0;
a = (float *)malloc (M * N * sizeof (*a));
if (!a) {
    printf ("host memory allocation failed");
    return EXIT_FAILURE;
}
for (j = 1; j <= N; j++) {
    for (i = 1; i <= M; i++) {
        a[IDX2F(i,j,M)] = (float)((i-1) * N + j);
    }
}
cudaStat = cudaMalloc ((void*)&devPtrA, M*N*sizeof(*a));
if (cudaStat != cudaSuccess) {
    printf ("device memory allocation failed");
    return EXIT_FAILURE;
}
stat = cublasCreate(&handle);
if (stat != CUBLAS_STATUS_SUCCESS) {
    printf ("CUBLAS initialization failed\n");
    return EXIT_FAILURE;
}
stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
if (stat != CUBLAS_STATUS_SUCCESS) {
    printf ("data download failed");
    cudaFree (devPtrA);
    cublasDestroy(handle);
    return EXIT_FAILURE;
}
modify (handle, devPtrA, M, N, 2, 3, 16.0f, 12.0f);
stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
if (stat != CUBLAS_STATUS_SUCCESS) {
    printf ("data upload failed");
    cudaFree (devPtrA);
    cublasDestroy(handle);
    return EXIT_FAILURE;
}
cudaFree (devPtrA);
cublasDestroy(handle);
for (j = 1; j <= N; j++) {
    for (i = 1; i <= M; i++) {
        printf ("%7.0f", a[IDX2F(i,j,M)]);
    }
    printf ("\n");
}
free(a);
return EXIT_SUCCESS;
}

```

-----  
 //Example 2. Application Using C and cuBLAS: 0-based indexing

//-----

```

#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"
#define M 6
#define N 5
#define IDX2C(i,j,ld) (((j)*(ld))+(i))

static __inline__ void modify (cublasHandle_t handle, float *m, int ldm, int n, int
p, int q, float alpha, float beta){
    cublasSscal (handle, n-q, &alpha, &m[IDX2C(p,q,ldm)], ldm);
    cublasSscal (handle, ldm-p, &beta, &m[IDX2C(p,q,ldm)], 1);
}

```

```

}

int main (void){
    cudaError_t cudaStat;
    cublasStatus_t stat;
    cublasHandle_t handle;
    int i, j;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            a[IDX2C(i,j,M)] = (float)(i * N + j + 1);
        }
    }
    cudaStat = cudaMalloc ((void**)&devPtrA, M*N*sizeof(*a));
    if (cudaStat != cudaSuccess) {
        printf ("device memory allocation failed");
        return EXIT_FAILURE;
    }
    stat = cublasCreate(&handle);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("CUBLAS initialization failed\n");
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    modify (handle, devPtrA, M, N, 1, 2, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    cudaFree (devPtrA);
    cublasDestroy(handle);
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            printf ("%7.0f", a[IDX2C(i,j,M)]);
        }
        printf ("\n");
    }
    free(a);
    return EXIT_SUCCESS;
}

```

---

## Chapter 2. Using the cuBLAS API

### 2.1. General Description

This section describes how to use the cuBLAS library API.

#### 2.1.1. Error status

All cuBLAS library function calls return the error status `cublasStatus_t`.

#### 2.1.2. cuBLAS context

The application must initialize the `handle` to the cuBLAS library context by calling the `cublasCreate()` function. Then, the `handle` is explicitly passed to every subsequent library function call. Once the application finishes using the library, it must call the function `cublasDestroy()` to release the resources associated with the cuBLAS library context.

This approach allows the user to explicitly control the library setup when using multiple host threads and multiple GPUs. For example, the application can use `cudaSetDevice()` to associate different devices with different host threads and in each of those host threads it can initialize a unique `handle` to the cuBLAS library context, which will use the particular device associated with that host thread. Then, the cuBLAS library function calls made with different `handle` will automatically dispatch the computation to different devices.

The device associated with a particular cuBLAS context is assumed to remain unchanged between the corresponding `cublasCreate()` and `cublasDestroy()` calls. In order for the cuBLAS library to use a different device in the same host thread, the application must set the new device to be used by calling `cudaSetDevice()` and then create another cuBLAS context, which will be associated with the new device, by calling `cublasCreate()`.

A cuBLAS library context is tightly coupled with the CUDA context that is current at the time of the `cublasCreate()` call. An application that uses multiple CUDA contexts is required to create a cuBLAS context per CUDA context and make sure the former never outlives the latter.

#### 2.1.3. Thread Safety

The library is thread safe and its functions can be called from multiple host threads, even with the same `handle`. When multiple threads share the same `handle`, extreme care needs to be taken when the `handle` configuration is changed because that change will affect potentially



subsequent cuBLAS calls in all threads. It is even more true for the destruction of the handle. So it is not recommended that multiple thread share the same cuBLAS handle.

## 2.1.4. Results Reproducibility

By design, all cuBLAS API routines from a given toolkit version generate the same bit-wise results at every run when executed on GPUs with the same architecture and the same number of SMs. However, bit-wise reproducibility is not guaranteed across toolkit versions because the implementation might differ due to some implementation changes.

This guarantee holds when a single CUDA stream is active only. If multiple concurrent streams are active, the library may optimize total performance by picking different internal implementations.

Note: The non-deterministic behavior of multi-stream execution is due to library optimizations in selecting internal workspace for the routines running in parallel streams. To avoid this effect, users can either:

- ▶ provide a separate workspace for each used stream using the `cublasSetWorkspace()` function, or
- ▶ have one cuBLAS handle per stream, or
- ▶ use `cublasLtMatmul()` instead of the `*gemm*` family of functions and provide user-owned workspace, or
- ▶ set the debug environment variable `CUBLAS_WORKSPACE_CONFIG` to `":16:8"` (may limit overall performance) or `":4096:8"` (will increase library footprint in GPU memory by approximately 24MiB).

Any of those settings will allow for deterministic behavior even with multiple concurrent streams sharing a single cuBLAS handle.

This behavior is expected to change in a future release.

For some routines such as `cublas<t>symv` and `cublas<t>hemv`, an alternate significantly faster routine can be chosen using the routine `cublasSetAtomicMode()`. In that case, the results are not guaranteed to be bit-wise reproducible because atomics are used for the computation.

## 2.1.5. Scalar Parameters

There are two categories of the functions that use scalar parameters:

- ▶ Functions that take `alpha` and `beta` parameters by reference on the host or the device as scaling factors, such as `gemm`.
- ▶ Functions that return a scalar result on the host or the device such as `amax()`, `amin()`, `asum()`, `rotg()`, `rotmg()`, `dot()`, and `nrm2()`.

For the functions of the first category, when the pointer mode is set to `CUBLAS_POINTER_MODE_HOST`, the scalar parameters `alpha` and `beta` can be on the stack or allocated on the heap, and should not be placed in managed memory. Underneath, the CUDA kernels related to those functions will be launched with the value of `alpha` and

beta. Therefore if they were allocated on the heap, they can be freed just after the return of the call even though the kernel launch is asynchronous. When the pointer mode is set to `CUBLAS_POINTER_MODE_DEVICE`, alpha and beta must be accessible on the device, and their values should not be modified until the kernel is done. Note that since `cudaFree()` does an implicit `cudaDeviceSynchronize()`, `cudaFree()` can still be called on alpha and beta just after the call but it would defeat the purpose of using this pointer mode in that case.

For the functions of the second category, when the pointer mode is set to `CUBLAS_POINTER_MODE_HOST`, these functions block the CPU, until the GPU has completed its computation and the results have been copied back to the host. When the pointer mode is set to `CUBLAS_POINTER_MODE_DEVICE`, these functions return immediately. In this case, similar to matrix and vector results, the scalar result is ready only when execution of the routine on the GPU has completed. This requires proper synchronization in order to read the result from the host.

In either case, the pointer mode `CUBLAS_POINTER_MODE_DEVICE` allows the library functions to execute completely asynchronously from the Host even when alpha and/or beta are generated by a previous kernel. For example, this situation can arise when iterative methods for solution of linear systems and eigenvalue problems are implemented using the cuBLAS library.

## 2.1.6. Parallelism with Streams

If the application uses the results computed by multiple independent tasks, CUDA streams can be used to overlap the computation performed in these tasks.

The application can conceptually associate each stream with each task. In order to achieve the overlap of computation between the tasks, the user should create CUDA streams using the function `cudaStreamCreate()` and set the stream to be used by each individual cuBLAS library routine by calling `cublasSetStream()` just before calling the actual cuBLAS routine. Note that `cublasSetStream()` resets the user-provided workspace to the default workspace pool; refer to [cublasSetWorkspace\(\)](#). Then, the computation performed in separate streams would be overlapped automatically when possible on the GPU. This approach is especially useful when the computation performed by a single task is relatively small and is not enough to fill the GPU with work.

We recommend using the new cuBLAS API with scalar parameters and results passed by reference in the device memory to achieve maximum overlap of the computation when using streams.

A particular application of streams, batching of multiple small kernels, is described in the following section.

## 2.1.7. Batching Kernels

In this section, we explain how to use streams to batch the execution of small kernels. For instance, suppose that we have an application where we need to make many small independent matrix-matrix multiplications with dense matrices.

It is clear that even with millions of small independent matrices we will not be able to achieve the same *GFLOPS* rate as with a one large matrix. For example, a single  $n \times n$  large matrix-

matrix multiplication performs  $n^3$  operations for  $n^2$  input size, while  $1024 \frac{n}{32} \times \frac{n}{32}$  small matrix-matrix multiplications perform  $1024 \left(\frac{n}{32}\right)^3 = \frac{n^3}{32}$  operations for the same input size. However, it is also clear that we can achieve a significantly better performance with many small independent matrices compared with a single small matrix.

The architecture family of GPUs allows us to execute multiple kernels simultaneously. Hence, in order to batch the execution of independent kernels, we can run each of them in a separate stream. In particular, in the above example we could create 1024 CUDA™ streams using the function `cudaStreamCreate()`, then preface each call to `cublas<T>gemm()` with a call to `cublasSetStream()` with a different stream for each of the matrix-matrix multiplications (note that `cublasSetStream()` resets user-provided workspace to the default workspace pool, see `cublasSetWorkspace()`). This will ensure that when possible the different computations will be executed concurrently. Although the user can create many streams, in practice it is not possible to have more than 32 concurrent kernels executing at the same time.

## 2.1.8. Cache Configuration

On some devices, L1 cache and shared memory use the same hardware resources. The cache configuration can be set directly with the CUDA Runtime function `cudaDeviceSetCacheConfig`. The cache configuration can also be set specifically for some functions using the routine `cudaFuncSetCacheConfig`. Refer to the [CUDA Runtime API](#) documentation for details about the cache configuration settings.

Because switching from one configuration to another can affect kernel's concurrency, the cuBLAS Library does not set any cache configuration preference and relies on the current setting. However, some cuBLAS routines, especially Level-3 routines, rely heavily on shared memory. Thus the cache preference setting might affect adversely their performance.

## 2.1.9. Static Library Support

The cuBLAS Library is also delivered in a static form as `libcublas_static.a` on Linux. The static cuBLAS library and all other static math libraries depend on a common thread abstraction layer library called `libcublibos.a`.

For example, on Linux, to compile a small application using cuBLAS, against the dynamic library, the following command can be used:

```
nvcc myCublasApp.c -lcublas -o myCublasApp
```

Whereas to compile against the static cuBLAS library, the following command must be used:

```
nvcc myCublasApp.c -lcublas_static -lcublibos -o myCublasApp
```

It is also possible to use the native Host C++ compiler. Depending on the Host operating system, some additional libraries like `pthread` or `d1` might be needed on the linking line. The following command on Linux is suggested :

```
g++ myCublasApp.c -lcublas_static -lculibos -lcudart_static -lpthread -ldl -I
<cuda-toolkit-path>/include -L <cuda-toolkit-path>/lib64 -o myCublasApp
```

Note that in the latter case, the library `cuda` is not needed. The CUDA Runtime will try to open explicitly the `cuda` library if needed. In the case of a system which does not have the CUDA driver installed, this allows the application to gracefully manage this issue and potentially run if a CPU-only path is available.

Starting with release 11.2, using the typed functions instead of the extension functions (`cublas**Ex()`) helps in reducing the binary size when linking to static cuBLAS Library.

## 2.1.10. GEMM Algorithms Numerical Behavior

Some GEMM algorithms split the computation along the dimension `K` to increase the GPU occupancy, especially when the dimension `K` is large compared to dimensions `M` and `N`. When this type of algorithm is chosen by the cuBLAS heuristics or explicitly by the user, the results of each split is summed deterministically into the resulting matrix to get the final result.

For the routines `cublas<t>gemmEx` and `cublasGemmEx`, when the compute type is greater than the output type, the sum of the split chunks can potentially lead to some intermediate overflows thus producing a final resulting matrix with some overflows. Those overflows might not have occurred if all the dot products had been accumulated in the compute type before being converted at the end in the output type. This computation side-effect can be easily exposed when the computeType is `CUDA_R_32F` and `Atype`, `Btype` and `Ctype` are in `CUDA_R_16F`. This behavior can be controlled using the compute precision mode `CUBLAS_MATH_DISALLOW_REDUCED_PRECISION_REDUCTION` with [`cublasSetMathMode\(\)`](#)

## 2.1.11. Tensor Core Usage

Tensor cores were first introduced with Volta GPUs (compute capability  $\geq$  sm\_70) and significantly accelerate matrix multiplications. Starting with cuBLAS version 11.0.0, the library will automatically make use of Tensor Core capabilities wherever possible, unless they are explicitly disabled by selecting pedantic compute modes in cuBLAS (refer to [`cublasSetMathMode\(\)`](#), [`cublasMath\_t`](#)).

It should be noted that the library will pick a Tensor Core enabled implementation wherever it determines that it would provide the best performance.

Starting with cuBLAS version 11.0.0 there are no longer any restrictions on matrix dimensions and memory alignments to use Tensor Cores. However, the best performance when using Tensor Cores can be achieved when the matrix dimensions and pointers meet certain memory alignment requirements. Specifically, all of the following conditions must be satisfied to get the most performance out of Tensor Cores:

- ▶ `m % 8 == 0`
- ▶ `k % 8 == 0`
- ▶ `op_B == CUBLAS_OP_N || n%8 == 0`
- ▶ `intptr_t(A) % 16 == 0`
- ▶ `intptr_t(B) % 16 == 0`

- ▶ `intptr_t(C) % 16 == 0`
- ▶ `intptr_t(A+lda) % 16 == 0`
- ▶ `intptr_t(B+ldb) % 16 == 0`
- ▶ `intptr_t(C+ldc) % 16 == 0`

## 2.1.12. CUDA Graphs Support

cuBLAS routines can be captured in CUDA Graph stream capture without restrictions in most situations.

The exception are routines that output results into host buffers (e.g. `cublas<T>dot` while pointer mode `CUBLAS_POINTER_MODE_HOST` is configured), as it enforces synchronization.

For input coefficients (such as `alpha`, `beta`) behavior depends on the pointer mode setting:

- ▶ In the case of `CUBLAS(LT)_POINTER_MODE_HOST`, coefficient values are captured in the graph.
- ▶ In the case of pointer modes with device pointers, coefficient value is accessed using the device pointer at the time of graph execution.



**Note:** Every time cuBLAS routines are captured in a new CUDA Graph, cuBLAS will allocate workspace memory on the device. This memory is only freed when the cuBLAS handle used during capture is deleted. To avoid this, use `cublasSetWorkspace()` function to provide user-owned workspace memory.

## 2.2. cuBLAS Datatypes Reference

### 2.2.1. `cublasHandle_t`

The `cublasHandle_t` type is a pointer type to an opaque structure holding the cuBLAS library context. The cuBLAS library context must be initialized using `cublasCreate()` and the returned handle must be passed to all subsequent library function calls. The context should be destroyed at the end using `cublasDestroy()`.

### 2.2.2. `cublasStatus_t`

The type is used for function status returns. All cuBLAS library functions return their status, which can have the following values.

Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The cuBLAS library was not initialized. This is usually caused by the lack of a prior <code>cublasCreate()</code> call, an error in the CUDA

Value	Meaning
	<p>Runtime API called by the cuBLAS routine, or an error in the hardware setup.</p> <p>To correct: call <code>cublasCreate()</code> before the function call; and check that the hardware, an appropriate version of the driver, and the cuBLAS library are correctly installed.</p>
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	<p>Resource allocation failed inside the cuBLAS library. This is usually caused by a <code>cudaMalloc()</code> failure.</p> <p>To correct: prior to the function call, deallocate previously allocated memory as much as possible.</p>
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<p>An unsupported value or parameter was passed to the function (a negative vector size, for example).</p> <p>To correct: ensure that all the parameters being passed have valid values.</p>
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	<p>The function requires a feature absent from the device architecture; usually caused by compute capability lower than 5.0.</p> <p>To correct: compile and run the application on a device with appropriate compute capability.</p>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	<p>An access to GPU memory space failed, which is usually caused by a failure to bind a texture.</p> <p>To correct: before the function call, unbind any previously bound textures.</p>
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	<p>The GPU program failed to execute. This is often caused by a launch failure of the kernel on the GPU, which can be caused by multiple reasons.</p> <p>To correct: check that the hardware, an appropriate version of the driver, and the cuBLAS library are correctly installed.</p>
<code>CUBLAS_STATUS_INTERNAL_ERROR</code>	<p>An internal cuBLAS operation failed. This error is usually caused by a <code>cudaMemcpyAsync()</code> failure.</p> <p>To correct: check that the hardware, an appropriate version of the driver, and the cuBLAS library are correctly installed. Also, check that the memory passed as a parameter to the routine is not being deallocated before the routine's completion.</p>
<code>CUBLAS_STATUS_NOT_SUPPORTED</code>	The functionality requested is not supported.
<code>CUBLAS_STATUS_LICENSE_ERROR</code>	The functionality requested requires some license and an error was detected when trying to check

Value	Meaning
	the current licensing. This error can happen if the license is not present or is expired or if the environment variable <code>NVIDIA_LICENSE_FILE</code> is not set properly.

### 2.2.3. `cublasOperation_t`

The `cublasOperation_t` type indicates which operation needs to be performed with the dense matrix. Its values correspond to Fortran characters `'N'` or `'n'` (non-transpose), `'T'` or `'t'` (transpose) and `'C'` or `'c'` (conjugate transpose) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_OP_N</code>	The non-transpose operation is selected.
<code>CUBLAS_OP_T</code>	The transpose operation is selected.
<code>CUBLAS_OP_C</code>	The conjugate transpose operation is selected.

### 2.2.4. `cublasFillMode_t`

The type indicates which part (lower or upper) of the dense matrix was filled and consequently should be used by the function. Its values correspond to Fortran characters `'L'` or `'l'` (lower) and `'U'` or `'u'` (upper) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_FILL_MODE_LOWER</code>	The lower part of the matrix is filled.
<code>CUBLAS_FILL_MODE_UPPER</code>	The upper part of the matrix is filled.
<code>CUBLAS_FILL_MODE_FULL</code>	The full matrix is filled.

### 2.2.5. `cublasDiagType_t`

The type indicates whether the main diagonal of the dense matrix is unity and consequently should not be touched or modified by the function. Its values correspond to Fortran characters `'N'` or `'n'` (non-unit) and `'U'` or `'u'` (unit) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_DIAG_NON_UNIT</code>	The matrix diagonal has non-unit elements.
<code>CUBLAS_DIAG_UNIT</code>	The matrix diagonal has unit elements.

## 2.2.6. `cublasSideMode_t`

The type indicates whether the dense matrix is on the left or right side in the matrix equation solved by a particular function. Its values correspond to Fortran characters `'L'` or `'1'` (left) and `'R'` or `'r'` (right) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_SIDE_LEFT</code>	The matrix is on the left side in the equation.
<code>CUBLAS_SIDE_RIGHT</code>	The matrix is on the right side in the equation.

## 2.2.7. `cublasPointerMode_t`

The `cublasPointerMode_t` type indicates whether the scalar values are passed by reference on the host or device. It is important to point out that if several scalar values are present in the function call, all of them must conform to the same single pointer mode. The pointer mode can be set and retrieved using `cublasSetPointerMode()` and `cublasGetPointerMode()` routines, respectively.

Value	Meaning
<code>CUBLAS_POINTER_MODE_HOST</code>	The scalars are passed by reference on the host.
<code>CUBLAS_POINTER_MODE_DEVICE</code>	The scalars are passed by reference on the device.

## 2.2.8. `cublasAtomsMode_t`

The type indicates whether cuBLAS routines which have an alternate implementation using atomics can be used. The atomics mode can be set and queried using `cublasSetAtomsMode()` and `cublasGetAtomsMode()` and routines, respectively.

Value	Meaning
<code>CUBLAS_ATOMICS_NOT_ALLOWED</code>	The usage of atomics is not allowed.
<code>CUBLAS_ATOMICS_ALLOWED</code>	The usage of atomics is allowed.

## 2.2.9. `cublasGemmAlgo_t`

`cublasGemmAlgo_t` type is an enumerator to specify the algorithm for matrix-matrix multiplication on GPU architectures up to `sm_75`. On `sm_80` and newer GPU architectures, this enumerator has no effect. cuBLAS has the following algorithm options:

Value	Meaning
<code>CUBLAS_GEMM_DEFAULT</code>	Apply Heuristics to select the GEMM algorithm.
<code>CUBLAS_GEMM_ALGO0</code> to <code>CUBLAS_GEMM_ALGO23</code>	Explicitly choose an Algorithm [0,23]. Note: Has no effect on NVIDIA Ampere Architecture GPUs and newer.



Value	Meaning
CUBLAS_GEMM_DEFAULT_TENSOR_OP[DEPRECATED]	This mode is deprecated and will be removed in a future release. Apply Heuristics to select the GEMM algorithm, while allowing use of reduced precision CUBLAS_COMPUTE_32F_FAST_16F kernels (for backward compatibility).
CUBLAS_GEMM_ALGO0_TENSOR_OP to CUBLAS_GEMM_ALGO15_TENSOR_OP[DEPRECATED]	Those values are deprecated and will be removed in a future release. Explicitly choose a Tensor core GEMM Algorithm [0,15]. Allows use of reduced precision CUBLAS_COMPUTE_32F_FAST_16F kernels (for backward compatibility). Note: Has no effect on NVIDIA Ampere Architecture GPUs and newer.

## 2.2.10. cublasMath\_t

cublasMath\_t enumerate type is used in cublasSetMathMode() to choose compute precision modes as defined in the following table. Since this setting does not directly control the use of Tensor Cores, the mode CUBLAS\_TENSOR\_OP\_MATH is being deprecated, and will be removed in a future release.

Value	Meaning
CUBLAS_DEFAULT_MATH	This is the default and highest-performance mode that uses compute and intermediate storage precisions with at least the same number of mantissa and exponent bits as requested. Tensor Cores will be used whenever possible.
CUBLAS_PEDANTIC_MATH	This mode uses the prescribed precision and standardized arithmetic for all phases of calculations and is primarily intended for numerical robustness studies, testing, and debugging. This mode might not be as performant as the other modes.
CUBLAS_TF32_TENSOR_OP_MATH	Enable acceleration of single-precision routines using TF32 tensor cores.
CUBLAS_MATH_DISALLOW_REduced_PRECISION_REDUCED	Force many reductions during matrix multiplications to use the accumulator type (that is, compute type) and not the output type in case of mixed precision routines where output type precision is less than the compute type precision. This is a flag that can be set (using a bitwise or operation) alongside any of the other values.
CUBLAS_TENSOR_OP_MATH [DEPRECATED]	This mode is deprecated and will be removed in a future release. Allows the library to use Tensor Core operations whenever possible. For single-precision GEMM routines cuBLAS will use the CUBLAS_COMPUTE_32F_FAST_16F compute type.

## 2.2.11. `cublasComputeType_t`

`cublasComputeType_t` enumerate type is used in `cublasGemmEx` and `cublasLtMatmul` (including all batched and strided batched variants) to choose compute precision modes as defined below.

Value	Meaning
<code>CUBLAS_COMPUTE_16F</code>	This is the default and highest-performance mode for 16-bit half precision floating point and all compute and intermediate storage precisions with at least 16-bit half precision. Tensor Cores will be used whenever possible.
<code>CUBLAS_COMPUTE_16F_PEDANTIC</code>	This mode uses 16-bit half precision floating point standardized arithmetic for all phases of calculations and is primarily intended for numerical robustness studies, testing, and debugging. This mode might not be as performant as the other modes since it disables use of tensor cores.
<code>CUBLAS_COMPUTE_32F</code>	This is the default 32-bit single precision floating point and uses compute and intermediate storage precisions of at least 32-bits.
<code>CUBLAS_COMPUTE_32F_PEDANTIC</code>	Uses 32-bit single-precision floating point arithmetic for all phases of calculations and also disables algorithmic optimizations such as Gaussian complexity reduction (3M).
<code>CUBLAS_COMPUTE_32F_FAST_16F</code>	Allows the library to use Tensor Cores with automatic down-conversion and 16-bit half-precision compute for 32-bit input and output matrices.
<code>CUBLAS_COMPUTE_32F_FAST_16BF</code>	Allows the library to use Tensor Cores with automatic down-conversion and bfloat16 compute for 32-bit input and output matrices. See <a href="#">Alternate Floating Point</a> section for more details on bfloat16.
<code>CUBLAS_COMPUTE_32F_FAST_TF32</code>	Allows the library to use Tensor cores with TF32 compute for 32-bit input and output matrices. See <a href="#">Alternate Floating Point</a> section for more details on TF32 compute.
<code>CUBLAS_COMPUTE_64F</code>	This is the default 64-bit double precision floating point and uses compute and intermediate storage precisions of at least 64-bits.
<code>CUBLAS_COMPUTE_64F_PEDANTIC</code>	Uses 64-bit double precision floating point arithmetic for all phases of calculations and

Value	Meaning
	also disables algorithmic optimizations such as Gaussian complexity reduction (3M).
CUBLAS_COMPUTE_32I	This is the default 32-bit integer mode and uses compute and intermediate storage precisions of at least 32-bits.
CUBLAS_COMPUTE_32I_PEDANTIC	Uses 32-bit integer arithmetic for all phases of calculations.



**Note:** Setting the environment variable `NVIDIA_TF32_OVERRIDE = 0` will override any defaults or programmatic configuration of NVIDIA libraries, and consequently, cuBLAS will not accelerate FP32 computations with TF32 tensor cores.

## 2.3. CUDA Datatypes Reference

This section describes types shared by multiple CUDA Libraries and defined in the header file `library_types.h`.

### 2.3.1. `cudaDataType_t`

The `cudaDataType_t` type is an enumerator to specify the data precision. It is used when the data reference does not carry the type itself (such as `void *`).

For example, it is used in the routine `cublasSgemvEx`.

Value	Meaning
CUDA_R_16F	The data type is a 16-bit real half precision floating-point.
CUDA_C_16F	The data type is a 32-bit structure comprised of two half precision floating-points representing a complex number.
CUDA_R_16BF	The data type is a 16-bit real bfloat16 floating-point.
CUDA_C_16BF	The data type is a 32-bit structure comprised of two bfloat16 floating-points representing a complex number.
CUDA_R_32F	The data type is a 32-bit real single precision floating-point.
CUDA_C_32F	The data type is a 64-bit structure comprised of two single precision floating-points representing a complex number.
CUDA_R_64F	The data type is a 64-bit real double precision floating-point.

Value	Meaning
CUDA_C_64F	The data type is a 128-bit structure comprised of two double precision floating-points representing a complex number.
CUDA_R_8I	The data type is an 8-bit real signed integer.
CUDA_C_8I	The data type is a 16-bit structure comprised of two 8-bit signed integers representing a complex number.
CUDA_R_8U	The data type is an 8-bit real unsigned integer.
CUDA_C_8U	The data type is a 16-bit structure comprised of two 8-bit unsigned integers representing a complex number.
CUDA_R_32I	The data type is a 32-bit real signed integer.
CUDA_C_32I	The data type is a 64-bit structure comprised of two 32-bit signed integers representing a complex number.
CUDA_R_8F_E4M3	The data type is an 8-bit real floating point in E4M3 format.
CUDA_R_8F_E5M2	The data type is an 8-bit real floating point in E5M2 format.

### 2.3.2. libraryPropertyType\_t

The `libraryPropertyType_t` is used as a parameter to specify which property is requested when using the routine `cublasGetProperty`.

Value	Meaning
MAJOR_VERSION	Enumerant to query the major version.
MINOR_VERSION	Enumerant to query the minor version.
PATCH_LEVEL	Number to identify the patch level.

## 2.4. cuBLAS Helper Function Reference

### 2.4.1. cublasCreate()

```
cublasStatus_t
cublasCreate(cublasHandle_t *handle)
```

This function initializes the cuBLAS library and creates a handle to an opaque structure holding the cuBLAS library context. It allocates hardware resources on the host and device

and must be called prior to making any other cuBLAS library calls. The cuBLAS library context is tied to the current CUDA device. To use the library on multiple devices, one cuBLAS handle needs to be created for each device. Furthermore, for a given device, multiple cuBLAS handles with different configurations can be created. Because `cublasCreate()` allocates some internal resources and the release of those resources by calling `cublasDestroy()` will implicitly call `cublasDeviceSynchronize()`, it is recommended to minimize the number of `cublasCreate()/cublasDestroy()` occurrences. For multi-threaded applications that use the same device from different threads, the recommended programming model is to create one cuBLAS handle per thread and use that cuBLAS handle for the entire life of the thread.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The initialization succeeded.
CUBLAS_STATUS_NOT_INITIALIZED	The CUDA™ Runtime initialization failed.
CUBLAS_STATUS_ALLOC_FAILED	The resources could not be allocated.
CUBLAS_STATUS_INVALID_VALUE	<code>handle == NULL</code>

## 2.4.2. `cublasDestroy()`

```
cublasStatus_t
cublasDestroy(cublasHandle_t handle)
```

This function releases hardware resources used by the cuBLAS library. This function is usually the last call with a particular handle to the cuBLAS library. Because `cublasCreate()` allocates some internal resources and the release of those resources by calling `cublasDestroy()` will implicitly call `cublasDeviceSynchronize()`, it is recommended to minimize the number of `cublasCreate()/cublasDestroy()` occurrences.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The shut down succeeded.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.3. `cublasGetVersion()`

```
cublasStatus_t
cublasGetVersion(cublasHandle_t handle, int *version)
```

This function returns the version number of the cuBLAS library.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The provided storage for library version number is not initialized (NULL).

## 2.4.4. `cublasGetProperty()`

```
cublasStatus_t
cublasGetProperty(libraryPropertyType type, int *value)
```

This function returns the value of the requested property in memory pointed to by `value`. Refer to `libraryPropertyType` for supported types.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	Invalid type value: <ul style="list-style-type: none"> <li>▶ If invalid type value or</li> <li>▶ <code>value == NULL</code></li> </ul>

## 2.4.5. `cublasGetStatusName()`

```
const char* cublasGetStatusName(cublasStatus_t status)
```

This function returns the string representation of a given status.

Return Value	Meaning
NULL-terminated string	The string representation of the <code>status</code> .

## 2.4.6. `cublasGetStatusString()`

```
const char* cublasGetStatusString(cublasStatus_t status)
```

This function returns the description string for a given status.

Return Value	Meaning
NULL-terminated string	The description of the <code>status</code> .

## 2.4.7. `cublasSetStream()`

```
cublasStatus_t
cublasSetStream(cublasHandle_t handle, cudaStream_t streamId)
```

This function sets the cuBLAS library stream, which will be used to execute all subsequent calls to the cuBLAS library functions. If the cuBLAS library stream is not set, all kernels use the *default* `NULL` stream. In particular, this routine can be used to change the stream between kernel launches and then to reset the cuBLAS library stream back to `NULL`. Additionally this function unconditionally resets the cuBLAS library workspace back to the default workspace pool (see `cublasSetWorkspace()`).

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The stream was set successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.8. cublasSetWorkspace()

```
cublasStatus_t
cublasSetWorkspace(cublasHandle_t handle, void *workspace, size_t
workspaceSizeInBytes)
```

This function sets the cuBLAS library workspace to a user-owned device buffer, which will be used to execute all subsequent calls to the cuBLAS library functions (on the currently set stream). If the cuBLAS library workspace is not set, all kernels will use the default workspace pool allocated during the cuBLAS context creation. In particular, this routine can be used to change the workspace between kernel launches. The workspace pointer has to be aligned to at least 256 bytes, otherwise CUBLAS\_STATUS\_INVALID\_VALUE error is returned. The `cublasSetStream()` function unconditionally resets the cuBLAS library workspace back to the default workspace pool. Too small `workspaceSizeInBytes` may cause some routines to fail with CUBLAS\_STATUS\_ALLOC\_FAILED error returned or cause large regressions in performance. Workspace size equal to or larger than 16KiB is enough to prevent CUBLAS\_STATUS\_ALLOC\_FAILED error, while a larger workspace can provide performance benefits for some routines. Recommended size of user-provided workspace is at least 4MiB (to match cuBLAS' default workspace pool).

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The stream was set successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The workspace pointer wasn't aligned to at least 256 bytes.

## 2.4.9. cublasGetStream()

```
cublasStatus_t
cublasGetStream(cublasHandle_t handle, cudaStream_t *streamId)
```

This function gets the cuBLAS library stream, which is being used to execute all calls to the cuBLAS library functions. If the cuBLAS library stream is not set, all kernels use the *default* NULL stream.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The stream was returned successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<code>streamId == NULL</code>

## 2.4.10. cublasGetPointerMode()

```
cublasStatus_t
cublasGetPointerMode(cublasHandle_t handle, cublasPointerMode_t *mode)
```

This function obtains the pointer mode used by the cuBLAS library. Please see the section on the `cublasPointerMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The pointer mode was obtained successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	mode == NULL

## 2.4.11. cublasSetPointerMode()

```
cublasStatus_t
cublasSetPointerMode(cublasHandle_t handle, cublasPointerMode_t mode)
```

This function sets the pointer mode used by the cuBLAS library. The *default* is for the values to be passed by reference on the host. Please see the section on the `cublasPointerMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The pointer mode was set successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	mode is not CUBLAS_POINTER_MODE_HOST or CUBLAS_POINTER_MODE_DEVICE

## 2.4.12. cublasSetVector()

```
cublasStatus_t
cublasSetVector(int n, int elemSize,
                const void *x, int incx, void *y, int incy)
```

This function copies `n` elements from a vector `x` in host memory space to a vector `y` in GPU memory space. Elements in both vectors are assumed to have a size of `elemSize` bytes. The storage spacing between consecutive elements is given by `incx` for the source vector `x` and by `incy` for the destination vector `y`.

Since column-major format for two-dimensional matrices is assumed, if a vector is part of a matrix, a vector increment equal to 1 accesses a (partial) column of that matrix. Similarly, using an increment equal to the leading dimension of the matrix results in accesses to a (partial) row of that matrix.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.



Return Value	Meaning
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>incx</code> , <code>incy</code> , <code>elemSize</code> ≤ 0.
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.13. `cublasGetVector()`

```
cublasStatus_t
cublasGetVector(int n, int elemSize,
               const void *x, int incx, void *y, int incy)
```

This function copies `n` elements from a vector `x` in GPU memory space to a vector `y` in host memory space. Elements in both vectors are assumed to have a size of `elemSize` bytes. The storage spacing between consecutive elements is given by `incx` for the source vector and `incy` for the destination vector `y`.

Since column-major format for two-dimensional matrices is assumed, if a vector is part of a matrix, a vector increment equal to 1 accesses a (partial) column of that matrix. Similarly, using an increment equal to the leading dimension of the matrix results in accesses to a (partial) row of that matrix.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>incx</code> , <code>incy</code> , <code>elemSize</code> ≤ 0.
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.14. `cublasSetMatrix()`

```
cublasStatus_t
cublasSetMatrix(int rows, int cols, int elemSize,
               const void *A, int lda, void *B, int ldb)
```

This function copies a tile of `rows` × `cols` elements from a matrix `A` in host memory space to a matrix `B` in GPU memory space. It is assumed that each element requires storage of `elemSize` bytes and that both matrices are stored in column-major format, with the leading dimension of the source matrix `A` and destination matrix `B` given in `lda` and `ldb`, respectively. The leading dimension indicates the number of rows of the allocated matrix, even if only a submatrix of it is being used.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>rows</code> , <code>cols</code> < 0 or <code>elemSize</code> , <code>lda</code> , <code>ldb</code> ≤ 0.
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.15. cublasGetMatrix()

```
cublasStatus_t
cublasGetMatrix(int rows, int cols, int elemSize,
               const void *A, int lda, void *B, int ldb)
```

This function copies a tile of `rows` x `cols` elements from a matrix `A` in GPU memory space to a matrix `B` in host memory space. It is assumed that each element requires storage of `elemSize` bytes and that both matrices are stored in column-major format, with the leading dimension of the source matrix `A` and destination matrix `B` given in `lda` and `ldb`, respectively. The leading dimension indicates the number of rows of the allocated matrix, even if only a submatrix of it is being used.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>rows</code> , <code>cols</code> <0 or <code>elemSize</code> , <code>lda</code> , <code>ldb</code> <=0.
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.16. cublasSetVectorAsync()

```
cublasStatus_t
cublasSetVectorAsync(int n, int elemSize, const void *hostPtr, int incx,
                   void *devicePtr, int incy, cudaStream_t stream)
```

This function has the same functionality as `cublasSetVector()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>incx</code> , <code>incy</code> , <code>elemSize</code> <=0.
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.17. cublasGetVectorAsync()

```
cublasStatus_t
cublasGetVectorAsync(int n, int elemSize, const void *devicePtr, int incx,
                   void *hostPtr, int incy, cudaStream_t stream)
```

This function has the same functionality as `cublasGetVector()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Return Value	Meaning
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>incx</code> , <code>incy</code> , <code>elemSize</code> $\leq 0$ .
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.18. `cublasSetMatrixAsync()`

```
cublasStatus_t
cublasSetMatrixAsync(int rows, int cols, int elemSize, const void *A,
                    int lda, void *B, int ldb, cudaStream_t stream)
```

This function has the same functionality as `cublasSetMatrix()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>rows</code> , <code>cols</code> $< 0$ or <code>elemSize</code> , <code>lda</code> , <code>ldb</code> $\leq 0$ .
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.19. `cublasGetMatrixAsync()`

```
cublasStatus_t
cublasGetMatrixAsync(int rows, int cols, int elemSize, const void *A,
                    int lda, void *B, int ldb, cudaStream_t stream)
```

This function has the same functionality as `cublasGetMatrix()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	The parameters <code>rows</code> , <code>cols</code> $< 0$ or <code>elemSize</code> , <code>lda</code> , <code>ldb</code> $\leq 0$ .
CUBLAS_STATUS_MAPPING_ERROR	There was an error accessing GPU memory.

## 2.4.20. `cublasSetAtomicMode()`

```
cublasStatus_t cublasSetAtomicMode(cublasHandle_t handle, cublasAtomicMode_t mode)
```

Some routines like `cublas<t>symv` and `cublas<t>hemv` have an alternate implementation that uses atomics to cumulate results. This implementation is generally significantly faster but can generate results that are not strictly identical from one run to the others. Mathematically, those different results are not significant but when debugging those differences can be prejudicial.

This function allows or disallows the usage of atomics in the cuBLAS library for all routines which have an alternate implementation. When not explicitly specified in the documentation of any cuBLAS routine, it means that this routine does not have an alternate implementation that use atomics. When atomics mode is disabled, each cuBLAS routine should produce the same results from one run to the other when called with identical parameters on the same hardware.

The default atomics mode of default initialized `cublasHandle_t` object is `CUBLAS_ATOMICS_NOT_ALLOWED`. Please see the section on the type for more details.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The atomics mode was set successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.

## 2.4.21. `cublasGetAtomicsMode()`

```
cublasStatus_t cublasGetAtomicsMode(cublasHandle_t handle, cublasAtomicsMode_t *mode)
```

This function queries the atomic mode of a specific cuBLAS context.

The default atomics mode of default initialized `cublasHandle_t` object is `CUBLAS_ATOMICS_NOT_ALLOWED`. Please see the section on the type for more details.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The atomics mode was queried successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	The argument <code>mode</code> is a NULL pointer.

## 2.4.22. `cublasSetMathMode()`

```
cublasStatus_t cublasSetMathMode(cublasHandle_t handle, cublasMath_t mode)
```

The `cublasSetMathMode` function enables you to choose the compute precision modes as defined by `cublasMath_t` (refer to [cublasMath\\_t](#)). Users are allowed to set the compute precision mode as a logical combination of them (except the deprecated `CUBLAS_TENSOR_OP_MATH`). For example, `cublasSetMathMode(handle, CUBLAS_DEFAULT_MATH | CUBLAS_MATH_DISALLOW_REDUCED_PRECISION_REDUCTION)`. Note that the default math mode is `CUBLAS_DEFAULT_MATH`.

For matrix and compute precisions allowed for `cublasGemmEx()` and `cublasLtMatmul()` APIs and their strided variants please refer to: [cublasGemmEx\(\)](#), [cublasGemmBatchedEx\(\)](#), [cublasGemmStridedBatchedEx\(\)](#) and [cublasLtMatmul\(\)](#).

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The math mode was set successfully.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	An invalid value for mode was specified.

Return Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.23. cublasGetMathMode()

```
cublasStatus_t cublasGetMathMode(cublasHandle_t handle, cublasMath_t *mode)
```

This function returns the math mode used by the library routines.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The math type was returned successfully.
CUBLAS_STATUS_INVALID_VALUE	If mode is NULL.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.24. cublasSetSmCountTarget()

```
cublasStatus_t cublasSetSmCountTarget(cublasHandle_t handle, int smCountTarget)
```

The `cublasSetSmCountTarget` function allows overriding the number of multiprocessors available to the library during kernels execution.

This option can be used to improve the library performance when cuBLAS routines are known to run concurrently with other work on different CUDA streams. E.g. a NVIDIA A100 GPU has 108 SM and there is a concurrent kernel running with grid size of 8, one can use `cublasSetSmCountTarget` with value 100 to override the library heuristics to optimize for running on 100 multiprocessors.

When set to 0 the library returns to its default behavior. The input value should not exceed the device's multiprocessor count, which can be obtained using `cudaDeviceGetAttribute`. Negative values are not accepted.

The user must ensure thread safety when modifying the library handle with this routine similar to when using `cublasSetStream`, and so on.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	SM count target was set successfully.
CUBLAS_STATUS_INVALID_VALUE	The value of <code>smCountTarget</code> outside of the allowed range.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.25. cublasGetSmCountTarget()

```
cublasStatus_t cublasGetSmCountTarget(cublasHandle_t handle, int *smCountTarget)
```

This function obtains the value previously programmed to the library handle.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	SM count target was set successfully.
CUBLAS_STATUS_INVALID_VALUE	smCountTarget is NULL.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

## 2.4.26. cublasLoggerConfigure()

```
cublasStatus_t cublasLoggerConfigure(
    int      logIsOn,
    int      logToStdOut,
    int      logToStdErr,
    const char* logFileName)
```

This function configures logging during runtime. Besides this type of configuration, it is possible to configure logging with special environment variables which will be checked by libcublas:

- CUBLAS\_LOGINFO\_DBG - Setup env. variable to "1" means turn on logging (by default, logging is off).
- CUBLAS\_LOGDEST\_DBG - Setup env. variable encodes how to log. "stdout", "stderr" means to output log messages to stdout or stderr, respectively. In the other case, it specifies "filename" of file.

### Parameters

#### logIsOn

*Input.* Turn on/off logging completely. By default is off, but is turned on by calling `cublasSetLoggerCallback` to user defined callback function.

#### logToStdOut

*Input.* Turn on/off logging to standard error I/O stream. By default is off.

#### logToStdErr

*Input.* Turn on/off logging to standard error I/O stream. By default is off.

#### logFileName

*Input.* Turn on/off logging to file in filesystem specified by its name. `cublasLoggerConfigure` copies content of `logFileName`. You should provide a null pointer if you're not interested in this type of logging.

### Returns

#### CUBLAS\_STATUS\_SUCCESS

Success.

## 2.4.27. cublasGetLoggerCallback()

```
cublasStatus_t cublasGetLoggerCallback(
    cublasLogCallback* userCallback)
```

This function retrieves function pointer to previously installed custom user defined callback function via `cublasSetLoggerCallback()` or zero otherwise.

### Parameters

#### **userCallback**

*Output.* Pointer to user defined callback function.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_INVALID_VALUE	userCallback is NULL.

## 2.4.28. cublasSetLoggerCallback()

```
cublasStatus_t cublasSetLoggerCallback(
    cublasLogCallback userCallback)
```

This function installs a custom user-defined callback function via cuBLAS C public API.

### Parameters

#### **userCallback**

*Input.* Pointer to user defined callback function.

### Returns

#### **CUBLAS\_STATUS\_SUCCESS**

Success.

## 2.5. cuBLAS Level-1 Function Reference

This section describes the Level-1 Basic Linear Algebra Subprograms (BLAS1) functions that perform scalar and vector-based operations. We will use abbreviations `<type>` for type and `<t>` for the corresponding short type to make a more concise and clear presentation of the implemented functions. Unless otherwise specified `<type>` and `<t>` have the following meanings:

<type>	<t>	Meaning
float	's' or 'S'	Real single-precision.
double	'd' or 'D'	Real double-precision.
cuComplex	'c' or 'C'	Complex single-precision.
cuDoubleComplex	'z' or 'Z'	Complex double-precision.

When the parameters and returned values of the function differ, which sometimes happens for complex input, the `<t>` can also have the following meanings 'Sc', 'Cs', 'Dz' and 'Zd'.

The abbreviation **Re(.)** and **Im(.)** will stand for the real and imaginary part of a number, respectively. Since the imaginary part of a real number does not exist, we will consider it to be zero, and can usually discard it from the equation where it is being used. Also, the  $\bar{\alpha}$  will denote the complex conjugate of  $\alpha$ .

In general throughout the documentation, the lowercase Greek symbols  $\alpha$  and  $\beta$  denote scalars, lowercase English letters in bold type **x** and **y** denote vectors and capital English letters **A**, **B** and **C** denote matrices.

## 2.5.1. `cublasI<t>amax()`

```
cublasStatus_t cublasIsamax(cublasHandle_t handle, int n,
                           const float *x, int incx, int *result)
cublasStatus_t cublasIdamax(cublasHandle_t handle, int n,
                           const double *x, int incx, int *result)
cublasStatus_t cublasIcamax(cublasHandle_t handle, int n,
                           const cuComplex *x, int incx, int *result)
cublasStatus_t cublasIzamax(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, int *result)
```

This function finds the (smallest) index of the element of the maximum magnitude.

Hence, the result is the first  $i$  such that  $|\text{Im}(x[j])| + |\text{Re}(x[j])|$  is maximum for  $i = 1, \dots, n$  and  $j = 1 + (i - 1) * \text{incx}$ . Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector x.
x	device	input	<type> vector with elements.
incx		input	Stride between consecutive elements of x.
result	host or device	output	The resulting index, which is 0 if $n, \text{incx} \leq 0$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	result is NULL.

For references. refer to:

[isamax](#), [idamax](#), [icamax](#), [izamax](#)



## 2.5.2. `cublasI<t>amin()`

```

cublasStatus_t cublasIsamin(cublasHandle_t handle, int n,
                           const float *x, int incx, int *result)
cublasStatus_t cublasIdamin(cublasHandle_t handle, int n,
                           const double *x, int incx, int *result)
cublasStatus_t cublasIcamin(cublasHandle_t handle, int n,
                           const cuComplex *x, int incx, int *result)
cublasStatus_t cublasIzamin(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, int *result)

```

This function finds the (smallest) index of the element of the minimum magnitude. Hence, the result is the first  $i$  such that  $|\text{Im}(x[j])| + |\text{Re}(x[j])|$  is minimum for  $i = 1, \dots, n$  and  $j = 1 + (i - 1) * \text{incx}$ . Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector x.
x	device	input	<type> vector with elements.
incx		input	Stride between consecutive elements of x.
result	host or device	output	The resulting index, which is 0 if $n, \text{incx} \leq 0$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	result is NULL.

For references please refer to:

[isamin](#)

## 2.5.3. `cublas<t>asum()`

```

cublasStatus_t cublasSasum(cublasHandle_t handle, int n,
                           const float *x, int incx, float *result)
cublasStatus_t cublasDasum(cublasHandle_t handle, int n,
                           const double *x, int incx, double *result)
cublasStatus_t cublasScasum(cublasHandle_t handle, int n,
                           const cuComplex *x, int incx, float *result)
cublasStatus_t cublasDzasum(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, double *result)

```

This function computes the sum of the absolute values of the elements of vector  $\mathbf{x}$ . Hence, the result is  $\sum_{i=1}^n |\text{Im}(\mathbf{x}[j])| + |\text{Re}(\mathbf{x}[j])|$  where  $j = 1 + (i - 1) * \text{incx}$ . Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector $\mathbf{x}$ .
$\mathbf{x}$	device	input	<type> vector with elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
result	host or device	output	The resulting index, which is 0.0 if $n, \text{incx} \leq 0$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	result is NULL.

For references please refer to:

[sasum](#), [dasum](#), [scasum](#), [dzasum](#)

## 2.5.4. cublas<t>axpy()

```

cublasStatus_t cublasSaxpy(cublasHandle_t handle, int n,
                           const float          *alpha,
                           const float          *x, int incx,
                           float               *y, int incy)
cublasStatus_t cublasDaxpy(cublasHandle_t handle, int n,
                           const double         *alpha,
                           const double         *x, int incx,
                           double              *y, int incy)
cublasStatus_t cublasCaxpy(cublasHandle_t handle, int n,
                           const cuComplex      *alpha,
                           const cuComplex      *x, int incx,
                           cuComplex            *y, int incy)
cublasStatus_t cublasZaxpy(cublasHandle_t handle, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           cuDoubleComplex      *y, int incy)

```

This function multiplies the vector  $\mathbf{x}$  by the scalar  $\alpha$  and adds it to the vector  $\mathbf{y}$  overwriting the latest vector with the result. Hence, the performed operation is  $\mathbf{y}[j] = \alpha \times \mathbf{x}[k] + \mathbf{y}[j]$  for  $i = 1, \dots, n$ ,  $k = 1 + (i - 1) * \text{incx}$  and  $j = 1 + (i - 1) * \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
alpha	host or device	input	<type> scalar used for multiplication.
n		input	Number of elements in the vector $x$ and $y$ .
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .
y	device	in/out	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $y$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[saxpy](#), [daxpy](#), [caxpy](#), [zaxpy](#)

## 2.5.5. `cublas<t>copy()`

```

cublasStatus_t cublasScopy(cublasHandle_t handle, int n,
                          const float *x, int incx,
                          float *y, int incy)
cublasStatus_t cublasDcopy(cublasHandle_t handle, int n,
                          const double *x, int incx,
                          double *y, int incy)
cublasStatus_t cublasCcopy(cublasHandle_t handle, int n,
                          const cuComplex *x, int incx,
                          cuComplex *y, int incy)
cublasStatus_t cublasZcopy(cublasHandle_t handle, int n,
                          const cuDoubleComplex *x, int incx,
                          cuDoubleComplex *y, int incy)

```

This function copies the vector  $x$  into the vector  $y$ . Hence, the performed operation is  $y[j] = x[k]$  for  $i = 1, \dots, n$ ,  $k = 1 + (i - 1) * \text{incx}$  and  $j = 1 + (i - 1) * \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector $x$ and $y$ .
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .

Param.	Memory	In/out	Meaning
y	device	output	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[scopy](#), [dcopy](#), [ccopy](#), [zcopy](#)

## 2.5.6. cublas<t>dot()

```

cublasStatus_t cublasSdot (cublasHandle_t handle, int n,
                          const float *x, int incx,
                          const float *y, int incy,
                          float *result)
cublasStatus_t cublasDdot (cublasHandle_t handle, int n,
                          const double *x, int incx,
                          const double *y, int incy,
                          double *result)
cublasStatus_t cublasCdotu(cublasHandle_t handle, int n,
                          const cuComplex *x, int incx,
                          const cuComplex *y, int incy,
                          cuComplex *result)
cublasStatus_t cublasCdotc(cublasHandle_t handle, int n,
                          const cuComplex *x, int incx,
                          const cuComplex *y, int incy,
                          cuComplex *result)
cublasStatus_t cublasZdotu(cublasHandle_t handle, int n,
                          const cuDoubleComplex *x, int incx,
                          const cuDoubleComplex *y, int incy,
                          cuDoubleComplex *result)
cublasStatus_t cublasZdotc(cublasHandle_t handle, int n,
                          const cuDoubleComplex *x, int incx,
                          const cuDoubleComplex *y, int incy,
                          cuDoubleComplex *result)

```

This function computes the dot product of vectors  $x$  and  $y$ . Hence, the result is  $\sum_{i=1}^n (x[k] \times y[j])$  where  $k = 1 + (i - 1) * \text{incx}$  and  $j = 1 + (i - 1) * \text{incy}$ . Notice that in the first equation the conjugate of the element of vector  $x$  should be used if the function name ends in character 'c' and that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vectors $x$ and $y$ .

Param.	Memory	In/out	Meaning
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.
result	host or device	output	The resulting dot product, which is 0.0 if n<=0.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sdot](#), [ddot](#), [cdotu](#), [cdotc](#), [zdotu](#), [zdotc](#)

## 2.5.7. cublas<t>nrm2()

```

cublasStatus_t cublasSnrm2(cublasHandle_t handle, int n,
                           const float *x, int incx, float *result)
cublasStatus_t cublasDnrm2(cublasHandle_t handle, int n,
                           const double *x, int incx, double *result)
cublasStatus_t cublasScnrm2(cublasHandle_t handle, int n,
                            const cuComplex *x, int incx, float *result)
cublasStatus_t cublasDznrm2(cublasHandle_t handle, int n,
                            const cuDoubleComplex *x, int incx, double *result)

```

This function computes the Euclidean norm of the vector x. The code uses a multiphase model of accumulation to avoid intermediate underflow and overflow, with the result being equivalent

to  $\sqrt{\sum_{i=1}^n (\mathbf{x}[j] \times \mathbf{x}[j])}$  where  $j = 1 + (i - 1) * \text{incx}$  in exact arithmetic. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector x.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
result	host or device	output	The resulting norm, which is 0.0 if n, incx<=0.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	result is NULL.

For references please refer to:

snrm2, snrm2, dnrn2, dnrn2, scnrm2, scnrm2, dznrm2

## 2.5.8. cublas<t>rot()

```

cublasStatus_t cublasSrot(cublasHandle_t handle, int n,
                        float *x, int incx,
                        float *y, int incy,
                        const float *c, const float *s)
cublasStatus_t cublasDrot(cublasHandle_t handle, int n,
                        double *x, int incx,
                        double *y, int incy,
                        const double *c, const double *s)
cublasStatus_t cublasCrot(cublasHandle_t handle, int n,
                        cuComplex *x, int incx,
                        cuComplex *y, int incy,
                        const float *c, const cuComplex *s)
cublasStatus_t cublasCsrot(cublasHandle_t handle, int n,
                        cuComplex *x, int incx,
                        cuComplex *y, int incy,
                        const float *c, const float *s)
cublasStatus_t cublasZrot(cublasHandle_t handle, int n,
                        cuDoubleComplex *x, int incx,
                        cuDoubleComplex *y, int incy,
                        const double *c, const cuDoubleComplex *s)
cublasStatus_t cublasZdrot(cublasHandle_t handle, int n,
                        cuDoubleComplex *x, int incx,
                        cuDoubleComplex *y, int incy,
                        const double *c, const double *s)

```

This function applies Givens rotation matrix (that is, rotation in the x,y plane counter-clockwise by angle defined by  $\cos(\alpha)=c$ ,  $\sin(\alpha)=s$ ):

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

to vectors  $x$  and  $y$ .

Hence, the result is  $x[k] = c \times x[k] + s \times y[j]$  and  $y[j] = -s \times x[k] + c \times y[j]$  where  $k = 1 + (i - 1) \times \text{incx}$  and  $j = 1 + (i - 1) \times \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.

Param.	Memory	In/out	Meaning
n		input	Number of elements in the vectors $x$ and $y$ .
x	device	in/out	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .
y	device	in/out	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $y$ .
c	host or device	input	Cosine element of the rotation matrix.
s	host or device	input	Sine element of the rotation matrix.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[srot](#), [drot](#), [crot](#), [csrot](#), [zrot](#), [zdrot](#)

## 2.5.9. cublas<t>rotg()

```

cublasStatus_t cublasSrotg(cublasHandle_t handle,
                          float *a, float *b,
                          float *c, float *s)
cublasStatus_t cublasDrotg(cublasHandle_t handle,
                          double *a, double *b,
                          double *c, double *s)
cublasStatus_t cublasCrotg(cublasHandle_t handle,
                          cuComplex *a, cuComplex *b,
                          float *c, cuComplex *s)
cublasStatus_t cublasZrotg(cublasHandle_t handle,
                          cuDoubleComplex *a, cuDoubleComplex *b,
                          double *c, cuDoubleComplex *s)

```

This function constructs the Givens rotation matrix

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

that zeros out the second entry of a  $2 \times 1$  vector  $(a, b)^T$ .

Then, for real numbers we can write

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$$

where  $c^2 + s^2 = 1$  and  $r = \sqrt{a^2 + b^2}$ . The parameters  $a$  and  $b$  are overwritten with  $r$  and  $z$ , respectively. The value of  $z$  is such that  $c$  and  $s$  may be recovered using the following rules:

$$(c, s) = \begin{cases} (\sqrt{1-z^2}, z) & \text{if } |z| < 1 \\ (0.0, 1.0) & \text{if } |z| = 1 \\ (1/z, \sqrt{1-z^2}) & \text{if } |z| > 1 \end{cases}$$

For complex numbers we can write

$$\begin{pmatrix} c & s \\ -\bar{s} & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$$

where  $c^2 + (s \times \bar{s}) = 1$  and  $r = \frac{a}{|a|} \times \|(a, b)^T\|_2$  with  $\|(a, b)^T\|_2 = \sqrt{|a|^2 + |b|^2}$  for  $a \neq 0$  and  $r = b$  for  $a = 0$ . Finally, the parameter  $a$  is overwritten with  $r$  on exit.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
a	host or device	in/out	<type> scalar that is overwritten with $r$ .
b	host or device	in/out	<type> scalar that is overwritten with $z$ .
c	host or device	output	Cosine element of the rotation matrix.
s	host or device	output	Sine element of the rotation matrix.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[srotg](#), [drotg](#), [crotg](#), [zrotg](#)

## 2.5.10. cublas<t>rotm()

```
cublasStatus_t cublasSrotm(cublasHandle_t handle, int n, float *x, int incx,
                          float *y, int incy, const float* param)
cublasStatus_t cublasDrotm(cublasHandle_t handle, int n, double *x, int incx,
                          double *y, int incy, const double* param)
```

This function applies the modified Givens transformation

$$H = \begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$$

to vectors  $x$  and  $y$ .



Hence, the result is  $\mathbf{x}[k] = h_{11} \times \mathbf{x}[k] + h_{12} \times \mathbf{y}[j]$  and  $\mathbf{y}[j] = h_{21} \times \mathbf{x}[k] + h_{22} \times \mathbf{y}[j]$  where  $k = 1 + (i - 1) \times \text{incx}$  and  $j = 1 + (i - 1) \times \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

The elements  $h_{11}$ ,  $h_{12}$ ,  $h_{21}$ , and  $h_{22}$  of matrix  $H$  are stored in `param[1]`, `param[2]`, `param[3]` and `param[4]`, respectively. The `flag=param[0]` defines the following predefined values for the matrix  $H$  entries

<code>flag=-1.0</code>	<code>flag= 0.0</code>	<code>flag= 1.0</code>	<code>flag=-2.0</code>
$\begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & h_{12} \\ h_{21} & 1.0 \end{pmatrix}$	$\begin{pmatrix} h_{11} & 1.0 \\ -1.0 & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{pmatrix}$

Notice that the values -1.0, 0.0 and 1.0 implied by the flag are not stored in `param`.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>n</code>		input	Number of elements in the vectors <code>x</code> and <code>y</code> .
<code>x</code>	device	in/out	<type> vector with <code>n</code> elements.
<code>incx</code>		input	Stride between consecutive elements of <code>x</code> .
<code>y</code>	device	in/out	<type> vector with <code>n</code> elements.
<code>incy</code>		input	Stride between consecutive elements of <code>y</code> .
<code>param</code>	host or device	input	<type> vector of 5 elements, where <code>param[0]</code> and <code>param[1-4]</code> contain the flag and matrix $H$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[srotm](#), [drotm](#)

## 2.5.11. `cublas<t>rotmg()`

```
cublasStatus_t cublasSrotmg(cublasHandle_t handle, float *d1, float *d2,
                           float *x1, const float *y1, float *param)
cublasStatus_t cublasDrotmg(cublasHandle_t handle, double *d1, double *d2,
                           double *x1, const double *y1, double *param)
```

This function constructs the modified Givens transformation

$$H = \begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$$

that zeros out the second entry of a  $2 \times 1$  vector  $(\sqrt{d1} * x1, \sqrt{d2} * y1)^T$ .

The `flag=param[0]` defines the following predefined values for the matrix  $H$  entries

<b>flag=-1.0</b>	<b>flag= 0.0</b>	<b>flag= 1.0</b>	<b>flag=-2.0</b>
$\begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & h_{12} \\ h_{21} & 1.0 \end{pmatrix}$	$\begin{pmatrix} h_{11} & 1.0 \\ -1.0 & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{pmatrix}$

Note that the values -1.0, 0.0 and 1.0 implied by the flag are not stored in `param`.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>d1</code>	host or device	in/out	<type> scalar that is overwritten on exit.
<code>d2</code>	host or device	in/out	<type> scalar that is overwritten on exit.
<code>x1</code>	host or device	in/out	<type> scalar that is overwritten on exit.
<code>y1</code>	host or device	input	<type> scalar.
<code>param</code>	host or device	output	<type> vector of 5 elements, where <code>param[0]</code> and <code>param[1-4]</code> contain the flag and matrix $H$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[srotmg](#), [drotmg](#)

## 2.5.12. `cublas<t>scal()`

```

cublasStatus_t cublasSscal(cublasHandle_t handle, int n,
                           const float *alpha,
                           float *x, int incx)
cublasStatus_t cublasDscal(cublasHandle_t handle, int n,
                           const double *alpha,
                           double *x, int incx)
cublasStatus_t cublasCscal(cublasHandle_t handle, int n,
                           const cuComplex *alpha,
                           cuComplex *x, int incx)
cublasStatus_t cublasCsscal(cublasHandle_t handle, int n,
                           const float *alpha,
                           cuComplex *x, int incx)
cublasStatus_t cublasZscal(cublasHandle_t handle, int n,
```

```

        const cuDoubleComplex *alpha,
        cuDoubleComplex *x, int incx)
cublasStatus_t cublasZdscal(cublasHandle_t handle, int n,
        const double *alpha,
        cuDoubleComplex *x, int incx)

```

This function scales the vector  $x$  by the scalar  $\alpha$  and overwrites it with the result. Hence, the performed operation is  $x[j] = \alpha \times x[j]$  for  $i = 1, \dots, n$  and  $j = 1 + (i - 1) \times \text{incx}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
alpha	host or device	input	<type> scalar used for multiplication.
n		input	Number of elements in the vector $x$ .
x	device	in/out	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sscal](#), [dscal](#), [cscal](#), [zscal](#), [cscal](#), [zscal](#), [zscal](#)

## 2.5.13. cublas<t>swap()

```

cublasStatus_t cublasSswap(cublasHandle_t handle, int n, float *x,
        int incx, float *y, int incy)
cublasStatus_t cublasDswap(cublasHandle_t handle, int n, double *x,
        int incx, double *y, int incy)
cublasStatus_t cublasCswap(cublasHandle_t handle, int n, cuComplex *x,
        int incx, cuComplex *y, int incy)
cublasStatus_t cublasZswap(cublasHandle_t handle, int n, cuDoubleComplex *x,
        int incx, cuDoubleComplex *y, int incy)

```

This function interchanges the elements of vector  $x$  and  $y$ . Hence, the performed operation is  $y[j] \leftrightarrow x[k]$  for  $i = 1, \dots, n$ ,  $k = 1 + (i - 1) \times \text{incx}$  and  $j = 1 + (i - 1) \times \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector $x$ and $y$ .
x	device	in/out	<type> vector with $n$ elements.

Param.	Memory	In/out	Meaning
incx		input	Stride between consecutive elements of x.
y	device	in/out	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sswap](#), [dswap](#), [cswap](#), [zswap](#)

## 2.6. cuBLAS Level-2 Function Reference

This section describes the Level-2 Basic Linear Algebra Subprograms (BLAS2) functions that perform matrix-vector operations.

### 2.6.1. cublas<t>gbmv()

```

cublasStatus_t cublasSgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx,
                           const float *beta,
                           float *y, int incy)
cublasStatus_t cublasDgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx,
                           const double *beta,
                           double *y, int incy)
cublasStatus_t cublasCgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const cuComplex *alpha,
                           const cuComplex *A, int lda,
                           const cuComplex *x, int incx,
                           const cuComplex *beta,
                           cuComplex *y, int incy)
cublasStatus_t cublasZgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the banded matrix-vector multiplication

$$\mathbf{y} = \alpha \text{op}(\mathbf{A})\mathbf{x} + \beta\mathbf{y}$$

where  $\mathbf{A}$  is a banded matrix with  $kl$  subdiagonals and  $ku$  superdiagonals,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars. Also, for matrix  $\mathbf{A}$

$$\text{op}(\mathbf{A}) = \begin{cases} \mathbf{A} & \text{if transa} == \text{CUBLAS\_OP\_N} \\ \mathbf{A}^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ \mathbf{A}^H & \text{if transa} == \text{CUBLAS\_OP\_H} \end{cases}$$

The banded matrix  $\mathbf{A}$  is stored column by column, with the main diagonal stored in row  $ku+1$  (starting in first position), the first superdiagonal stored in row  $ku$  (starting in second position), the first subdiagonal stored in row  $ku+2$  (starting in first position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(ku+1+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [\max(1, j-ku), \min(m, j+kl)]$ . Also, the elements in the array  $\mathbf{A}$  that do not conceptually correspond to the elements in the banded matrix (the top left  $ku \times ku$  and bottom right  $kl \times kl$  triangles) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
trans		input	Operation $\text{op}(\mathbf{A})$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\mathbf{A}$ .
n		input	Number of columns of matrix $\mathbf{A}$ .
kl		input	Number of subdiagonals of matrix $\mathbf{A}$ .
ku		input	Number of superdiagonals of matrix $\mathbf{A}$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times n$ with $lda \geq kl+ku+1$ .
lda		input	Leading dimension of two-dimensional array used to store matrix $\mathbf{A}$ .
x	device	input	<type> vector with $n$ elements if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $m$ elements otherwise.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ then $\mathbf{y}$ does not have to be a valid input.
y	device	in/out	<type> vector with $m$ elements if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $n$ elements otherwise.
incy		input	Stride between consecutive elements of $\mathbf{y}$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n, k_l, k_u &lt; 0</math> or</li> <li>▶ if <math>lda &lt; (k_l + k_u + 1)</math> or</li> <li>▶ if <math>incx, incy == 0</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N, CUBLAS\_OP\_T, CUBLAS\_OP\_C}</math> or</li> <li>▶ <math>alpha, beta == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgbmv](#), [dgbmv](#), [cgbmv](#), [zgbmv](#)

## 2.6.2. `cublas<t>gemv()`

```

cublasStatus_t cublasSgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx,
                           const float *beta,
                           float *y, int incy)
cublasStatus_t cublasDgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx,
                           const double *beta,
                           double *y, int incy)
cublasStatus_t cublasCgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const cuComplex *alpha,
                           const cuComplex *A, int lda,
                           const cuComplex *x, int incx,
                           const cuComplex *beta,
                           cuComplex *y, int incy)
cublasStatus_t cublasZgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the matrix-vector multiplication

$$\mathbf{y} = \alpha \text{op}(\mathbf{A})\mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $m \times n$  matrix stored in column-major format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars. Also, for matrix  $\mathbf{A}$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_H} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
trans		input	Operation op(A) that is non- or (conj.) transpose.
m		input	Number of rows of matrix A.
n		input	Number of columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension lda x n with lda >= max(1,m). Before entry, the leading m by n part of the array A must contain the matrix of coefficients. Unchanged on exit.
lda		input	Leading dimension of two-dimensional array used to store matrix A. lda must be at least max(1,m).
x	device	input	<type> vector at least (1+(n-1)*abs(incx)) elements if transa==CUBLAS_OP_N and at least (1+(m-1)*abs(incx)) elements otherwise.
incx		input	Stride between consecutive elements of x.
beta	host or device	input	<type> scalar used for multiplication, if beta==0 then y does not have to be a valid input.
y	device	in/out	<type> vector at least (1+(m-1)*abs(incy)) elements if transa==CUBLAS_OP_N and at least (1+(n-1)*abs(incy)) elements otherwise.
incy		input	Stride between consecutive elements of y.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters m,n<0 or incx,incy=0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgemv](#), [dgemv](#), [cgemv](#), [zgemv](#)

### 2.6.3. `cublas<t>ger()`

```

cublasStatus_t cublasSger(cublasHandle_t handle, int m, int n,
                          const float *alpha,
                          const float *x, int incx,
                          const float *y, int incy,
                          float *A, int lda)
cublasStatus_t cublasDger(cublasHandle_t handle, int m, int n,
                          const double *alpha,
                          const double *x, int incx,
                          const double *y, int incy,
                          double *A, int lda)
cublasStatus_t cublasCgeru(cublasHandle_t handle, int m, int n,
                          const cuComplex *alpha,
                          const cuComplex *x, int incx,
                          const cuComplex *y, int incy,
                          cuComplex *A, int lda)
cublasStatus_t cublasCgerc(cublasHandle_t handle, int m, int n,
                          const cuComplex *alpha,
                          const cuComplex *x, int incx,
                          const cuComplex *y, int incy,
                          cuComplex *A, int lda)
cublasStatus_t cublasZgeru(cublasHandle_t handle, int m, int n,
                          const cuDoubleComplex *alpha,
                          const cuDoubleComplex *x, int incx,
                          const cuDoubleComplex *y, int incy,
                          cuDoubleComplex *A, int lda)
cublasStatus_t cublasZgerc(cublasHandle_t handle, int m, int n,
                          const cuDoubleComplex *alpha,
                          const cuDoubleComplex *x, int incx,
                          const cuDoubleComplex *y, int incy,
                          cuDoubleComplex *A, int lda)

```

This function performs the rank-1 update

$$A = \begin{cases} \alpha \mathbf{xy}^T + A & \text{if ger(), geru() is called} \\ \alpha \mathbf{xy}^H + A & \text{if gerc() is called} \end{cases}$$

where  $A$  is a  $m \times n$  matrix stored in column-major format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  is a scalar.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
m		input	Number of rows of matrix A.
n		input	Number of columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with m elements.
incx		input	Stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.
A	device	in/out	<type> array of dimension lda x n with lda >= max(1, m).



Param.	Memory	In/out	Meaning
lda		input	Leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>► If <math>m &lt; 0</math> or <math>n &lt; 0</math></li> <li>► if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>► if <math>alpha == \text{NULL}</math> or</li> <li>► <math>lda &lt; \max\{1, m\}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sger](#), [dger](#), [cgeru](#), [cgerc](#), [zgeru](#), [zgerc](#)

## 2.6.4. `cublas<t>sbmv()`

```

cublasStatus_t cublasSsbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, int k, const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx,
                           const float *beta, float *y, int incy)
cublasStatus_t cublasDsbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, int k, const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx,
                           const double *beta, double *y, int incy)

```

This function performs the symmetric banded matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $A$  is a  $n \times n$  symmetric banded matrix with  $k$  subdiagonals and superdiagonals,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the symmetric banded matrix  $A$  is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [j, \min(m, j+k)]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the bottom right  $k \times k$  triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the symmetric banded matrix  $A$  is stored column by column, with the main diagonal of the matrix stored in row  $k+1$ , the first superdiagonal

in row  $k$  (starting at second position), the second superdiagonal in row  $k-1$  (starting at third position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+k+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [\max(1, j-k), j]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the top left  $k \times k$  triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $A$ .
k		input	Number of sub- and super-diagonals of matrix $A$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times n$ with $lda \geq k+1$ .
lda		input	Leading dimension of two-dimensional array used to store matrix $A$ .
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .
beta	host or device	input	<type> scalar used for multiplication, if $\beta == 0$ then $y$ does not have to be a valid input.
y	device	in/out	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $y$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math> or</li> <li>▶ <math>lda &lt; \{1 + k\}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssbmv](#), [dsbmv](#)

## 2.6.5. `cublas<t>spmv()`

```
cublasStatus_t cublasSspmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const float *alpha, const float *AP,
                           const float *x, int incx, const float *beta,
                           float *y, int incy)
cublasStatus_t cublasDspmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const double *alpha, const double *AP,
                           const double *x, int incx, const double *beta,
                           double *y, int incy)
```

This function performs the symmetric packed matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $n \times n$  symmetric matrix stored in packed format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + ((2*n-j+1)*j)/2]` for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + (j*(j+1))/2]` for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $\mathbf{A}$ lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $\mathbf{A}$ .
alpha	host or device	input	<type> scalar used for multiplication.
AP	device	input	<type> array with $\mathbf{A}$ stored in packed format.
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then $\mathbf{y}$ does not have to be a valid input.
y	device	input	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $\mathbf{y}$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sspmv](#), [dspmv](#)

## 2.6.6. `cublas<t>spr()`

```

cublasStatus_t cublasSspr(cublasHandle_t handle, cublasFillMode_t uplo,
    int n, const float *alpha,
    const float *x, int incx, float *AP)
cublasStatus_t cublasDspr(cublasHandle_t handle, cublasFillMode_t uplo,
    int n, const double *alpha,
    const double *x, int incx, double *AP)

```

This function performs the packed symmetric rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^T + A$$

where  $A$  is a  $n \times n$  symmetric matrix stored in packed format,  $\mathbf{x}$  is a vector, and  $\alpha$  is a scalar.

If  $uplo == \text{CUBLAS\_FILL\_MODE\_LOWER}$  then the elements in the lower triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + ((2*n - j + 1) * j) / 2]$  for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If  $uplo == \text{CUBLAS\_FILL\_MODE\_UPPER}$  then the elements in the upper triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + (j * (j + 1)) / 2]$  for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.

Param.	Memory	In/out	Meaning
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $A$ .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
AP	device	in/out	<type> array with $A$ stored in packed format.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ <math>alpha == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sspr](#), [dspr](#)

## 2.6.7. cublas<t>spr2()

```

cublasStatus_t cublasSspr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const float *alpha,
                           const float *x, int incx,
                           const float *y, int incy, float *AP)
cublasStatus_t cublasDspr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const double *alpha,
                           const double *x, int incx,
                           const double *y, int incy, double *AP)

```

This function performs the packed symmetric rank-2 update

$$A = \alpha(\mathbf{xy}^T + \mathbf{yx}^T) + A$$

where  $A$  is a  $n \times n$  symmetric matrix stored in packed format,  $\mathbf{x}$  is a vector, and  $\alpha$  is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + ((2*n - j + 1) * j) / 2]$  for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + (j * (j + 1)) / 2]$  for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $A$ .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .
y	device	input	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $y$ .
AP	device	in/out	<type> array with $A$ stored in packed format.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq CUBLAS_FILL_MODE\_LOWER</math>, <math>CUBLAS_FILL_MODE\_UPPER</math> or</li> <li>▶ <math>alpha == NULL</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sspr2](#), [dspr2](#)

## 2.6.8. `cublas<t>symv()`

```

cublasStatus_t cublasSsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx, const float
                           *beta,
                           float *y, int incy)
cublasStatus_t cublasDsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx, const double
                           *beta,
                           double *y, int incy)
cublasStatus_t cublasCsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex *alpha, /* host or device
                           pointer */
                           const cuComplex *A, int lda,
                           const cuComplex *x, int incx, const cuComplex
                           *beta,
                           cuComplex *y, int incy)
cublasStatus_t cublasZsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx, const cuDoubleComplex
                           *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the symmetric matrix-vector multiplication.

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $n \times n$  symmetric matrix stored in lower or upper mode,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

This function has an alternate faster implementation using atomics that can be enabled with `cublasSetAtomicMode()`.

Please see the section on the function `cublasSetAtomicMode()` for more details about the usage of atomics.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times n$ with $lda \geq \max(1, n)$ .
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	input	<type> vector with n elements.

Param.	Memory	In/out	Meaning
incx		input	Stride between consecutive elements of x.
beta	host or device	input	<type> scalar used for multiplication, if beta==0 then y does not have to be a valid input.
y	device	in/out	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ <math>lda &lt; n</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssymv](#), [dsymv](#)

## 2.6.9. `cublas<t>syr()`

```

cublasStatus_t cublasSsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const float *x, int incx, float
                          *A, int lda)
cublasStatus_t cublasDsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const double *x, int incx, double
                          *A, int lda)
cublasStatus_t cublasCsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const cuComplex *alpha,
                          const cuComplex *x, int incx, cuComplex
                          *A, int lda)
cublasStatus_t cublasZsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const cuDoubleComplex *alpha,
                          const cuDoubleComplex *x, int incx, cuDoubleComplex
                          *A, int lda)

```

This function performs the symmetric rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^T + A$$



where  $A$  is a  $n \times n$  symmetric matrix stored in column-major format,  $\mathbf{x}$  is a vector, and  $\alpha$  is a scalar.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
A	device	in/out	<type> array of dimensions lda x n, with $lda \geq \max(1, n)$ .
lda		input	Leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> or</li> <li>▶ <math>alpha == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyr](#), [dsyr](#)

## 2.6.10. `cusblas<t>syr2()`

```

cusblasStatus_t cusblasSsyr2(cusblasHandle_t handle, cusblasFillMode_t uplo, int n,
                             const float          *alpha, const float
                             *x, int incx,
                             const float          *y, int incy, float
                             *A, int lda
cusblasStatus_t cusblasDsyr2(cusblasHandle_t handle, cusblasFillMode_t uplo, int n,
```

```

    const double      *alpha, const double
    *x, int incx,      const double      *y, int incy, double
    *A, int lda
cublasStatus_t cublasCsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
    const cuComplex    *alpha, const cuComplex
    *x, int incx,      const cuComplex    *y, int incy, cuComplex
    *A, int lda
cublasStatus_t cublasZsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
    const cuDoubleComplex *alpha, const cuDoubleComplex
    *x, int incx,      const cuDoubleComplex *y, int incy, cuDoubleComplex
    *A, int lda

```

This function performs the symmetric rank-2 update

$$A = \alpha(\mathbf{xy}^T + \mathbf{yx}^T) + A$$

where  $A$  is a  $n \times n$  symmetric matrix stored in column-major format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  is a scalar.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.
A	device	in/out	<type> array of dimensions lda x n, with lda>=max(1,n).
lda		input	Leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>► If <math>n &lt; 0</math> or</li> <li>► if <math>incx = 0</math> or <math>incy = 0</math> or</li> </ul>

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if <code>uplo != CUBLAS_FILL_MODE_LOWER, CUBLAS_FILL_MODE_UPPER</code> or</li> <li>▶ if <code>alpha == NULL</code> or</li> <li>▶ <code>lda &lt; max(1, n)</code></li> </ul>
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[ssyr2](#), [dsyr2](#)

## 2.6.11. `cublas<t>tbmv()`

```

cublasStatus_t cublasStbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const float *A, int lda,
                           float *x, int incx)
cublasStatus_t cublasDtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const double *A, int lda,
                           double *x, int incx)
cublasStatus_t cublasCtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuComplex *A, int lda,
                           cuComplex *x, int incx)
cublasStatus_t cublasZtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function performs the triangular banded matrix-vector multiplication

$\mathbf{x} = \text{op}(\mathbf{A})\mathbf{x}$

where  $\mathbf{A}$  is a triangular banded matrix, and  $\mathbf{x}$  is a vector. Also, for matrix  $\mathbf{A}$

$$\text{op}(\mathbf{A}) = \begin{cases} \mathbf{A} & \text{if transa} == \text{CUBLAS_OP_N} \\ \mathbf{A}^T & \text{if transa} == \text{CUBLAS_OP_T} \\ \mathbf{A}^H & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

If `uplo == CUBLAS_FILL_MODE_LOWER` then the triangular banded matrix  $\mathbf{A}$  is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [j, \min(m, j+k)]$ . Also, the elements in the array  $\mathbf{A}$  that do not conceptually correspond to the elements in the banded matrix (the bottom right  $k \times k$  triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the triangular banded matrix  $\mathbf{A}$  is stored column by column, with the main diagonal of the matrix stored in row  $k+1$ , the first superdiagonal in row  $k$  (starting at second position), the second superdiagonal in row  $k-1$  (starting at third position), and so on. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+k+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [\max(1, j-k), j]$ . Also, the elements in the array  $\mathbf{A}$  that do

not conceptually correspond to the elements in the banded matrix (the top left  $k \times k$  triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	Number of rows and columns of matrix A.
k		input	Number of sub- and super-diagonals of matrix .
A	device	input	<type> array of dimension lda x n, with lda>=k+1.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	in/out	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>diag \neq \text{CUBLAS\_DIAG\_UNIT}, \text{CUBLAS\_DIAG\_NON\_UNIT}</math> or</li> <li>▶ <math>lda &lt; (1 + k)</math></li> </ul>
CUBLAS_STATUS_ALLOC_FAILED	The allocation of internal scratch memory failed.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[stbmv](#), [dtbmv](#), [ctbmv](#), [ztbmv](#)

## 2.6.12. cublas<t>tbsv()

```

cublasStatus_t cublasStbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const float *A, int lda,
                           float *x, int incx)
cublasStatus_t cublasDtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const double *A, int lda,
                           double *x, int incx)
cublasStatus_t cublasCtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuComplex *A, int lda,
                           cuComplex *x, int incx)
cublasStatus_t cublasZtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function solves the triangular banded linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where  $A$  is a triangular banded matrix, and  $\mathbf{x}$  and  $\mathbf{b}$  are vectors. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $\mathbf{x}$  overwrites the right-hand-sides  $\mathbf{b}$  on exit.

No test for singularity or near-singularity is included in this function.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the triangular banded matrix  $A$  is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [j, \min(m, j+k)]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the bottom right  $k \times k$  triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the triangular banded matrix  $A$  is stored column by column, with the main diagonal of the matrix stored in row  $k+1$ , the first superdiagonal in row  $k$  (starting at second position), the second superdiagonal in row  $k-1$  (starting at third position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+k+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [\max(1, j-k), j]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the top left  $k \times k$  triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.

Param.	Memory	In/out	Meaning
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	Number of rows and columns of matrix A.
k		input	Number of sub- and super-diagonals of matrix A.
A	device	input	<type> array of dimension lda x n, with lda >= k+1.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	in/out	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>trans \neq CUBLAS\_OP\_N, CUBLAS\_OP\_C, CUBLAS\_OP\_T</math> or</li> <li>▶ if <math>uplo \neq CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER</math> or</li> <li>▶ if <math>diag \neq CUBLAS\_DIAG\_UNIT, CUBLAS\_DIAG\_NON\_UNIT</math> or</li> <li>▶ <math>lda &lt; (1 + k)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[stbsv](#), [dtbsv](#), [ctbsv](#), [ztbsv](#)

### 2.6.13. `cublas<t>tpmv()`

```

cublasStatus_t cublasStpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float *A,
                           float *x, int incx)
cublasStatus_t cublasDtpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double *A,
                           double *x, int incx)
cublasStatus_t cublasCtpmv(cublasHandle_t handle, cublasFillMode_t uplo,

```

```

        cublasOperation_t trans, cublasDiagType_t diag,
        int n, const cuComplex *AP,
        cuComplex *x, int incx)
cublasStatus_t cublasZtpmv(cublasHandle_t handle, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        int n, const cuDoubleComplex *AP,
        cuDoubleComplex *x, int incx)

```

This function performs the triangular packed matrix-vector multiplication

$\mathbf{x} = \text{op}(\mathbf{A})\mathbf{x}$

where  $\mathbf{A}$  is a triangular matrix stored in packed format, and  $\mathbf{x}$  is a vector. Also, for matrix  $\mathbf{A}$

$$\text{op}(\mathbf{A}) = \begin{cases} \mathbf{A} & \text{if transa} == \text{CUBLAS\_OP\_N} \\ \mathbf{A}^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ \mathbf{A}^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the triangular matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + ((2*n-j+1)*j)/2]` for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the triangular matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + (j*(j+1))/2]` for  $A(i, j)$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $\mathbf{A}$ lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(\mathbf{A})$ that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix $\mathbf{A}$ are unity and should not be accessed.
n		input	Number of rows and columns of matrix $\mathbf{A}$ .
AP	device	input	<type> array with $\mathbf{A}$ stored in packed format.
x	device	in/out	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	► If $n < 0$ or

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if <code>incx == 0</code> or</li> <li>▶ if <code>uplo != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER</code> or</li> <li>▶ if <code>trans != CUBLAS_OP_N, CUBLAS_OP_T, CUBLAS_OP_C</code> or</li> <li>▶ <code>diag != CUBLAS_DIAG_UNIT, CUBLAS_DIAG_NON_UNIT</code></li> </ul>
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	The allocation of internal scratch memory failed.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[stpmv](#), [dtpmv](#), [ctpmv](#), [ztpmv](#)

## 2.6.14. `cublas<t>tpsv()`

```

cublasStatus_t cublasStpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float *AP,
                           float *x, int incx)
cublasStatus_t cublasDtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double *AP,
                           double *x, int incx)
cublasStatus_t cublasCtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex *AP,
                           cuComplex *x, int incx)
cublasStatus_t cublasZtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *AP,
                           cuDoubleComplex *x, int incx)

```

This function solves the packed triangular linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where  $A$  is a triangular matrix stored in packed format, and  $\mathbf{x}$  and  $\mathbf{b}$  are vectors. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $\mathbf{x}$  overwrites the right-hand-sides  $\mathbf{b}$  on exit.

No test for singularity or near-singularity is included in this function.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the triangular matrix  $A$  are packed together column by column without gaps, so that the element



$A(i, j)$  is stored in the memory location  $AP[i + ((2*n - j + 1) * j) / 2]$  for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the triangular matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + (j * (j + 1)) / 2]$  for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation $op(A)$ that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix are unity and should not be accessed.
n		input	Number of rows and columns of matrix $A$ .
AP	device	input	<type> array with $A$ stored in packed format.
x	device	in/out	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>trans \neq CUBLAS\_OP\_N, CUBLAS\_OP\_C, CUBLAS\_OP\_T</math> or</li> <li>▶ if <math>uplo \neq CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER</math> or</li> <li>▶ <math>diag \neq CUBLAS\_DIAG\_UNIT, CUBLAS\_DIAG\_NON\_UNIT</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[stpsv](#), [dtpsv](#), [ctpsv](#), [ztpsv](#)

## 2.6.15. `cublas<T>trmv()`

```

cublasStatus_t cublasStrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float *A, int lda,
                           float *x, int incx)
cublasStatus_t cublasDtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double *A, int lda,
                           double *x, int incx)
cublasStatus_t cublasCtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex *A, int lda,
                           cuComplex *x, int incx)
cublasStatus_t cublasZtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function performs the triangular matrix-vector multiplication

$$\mathbf{x} = \text{op}(\mathbf{A})\mathbf{x}$$

where  $\mathbf{A}$  is a triangular matrix stored in lower or upper mode with or without the main diagonal, and  $\mathbf{x}$  is a vector. Also, for matrix  $\mathbf{A}$

$$\text{op}(\mathbf{A}) = \begin{cases} \mathbf{A} & \text{if transa} == \text{CUBLAS\_OP\_N} \\ \mathbf{A}^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ \mathbf{A}^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $\mathbf{A}$ lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(\mathbf{A})$ (that is, non- or conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix $\mathbf{A}$ are unity and should not be accessed.
n		input	Number of rows and columns of matrix $\mathbf{A}$ .
A	device	input	<type> array of dimensions $\text{lda} \times n$ , with $\text{lda} \geq \max(1, n)$ .
lda		input	Leading dimension of two-dimensional array used to store matrix $\mathbf{A}$ .
x	device	in/out	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Error Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>diag \neq \text{CUBLAS\_DIAG\_UNIT}, \text{CUBLAS\_DIAG\_NON\_UNIT}</math> or</li> <li>▶ <math>lda &lt; \max\{1, n\}</math></li> </ul>
CUBLAS_STATUS_ALLOC_FAILED	The allocation of internal scratch memory failed.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strmv](#), [dtrmv](#), [ctrmv](#), [ztrmv](#)

## 2.6.16. `cusblas<t>trsv()`

```

cusblasStatus_t cusblasStrsv(cusblasHandle_t handle, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int n, const float *A, int lda,
                             float *x, int incx)
cusblasStatus_t cusblasDtrsv(cusblasHandle_t handle, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int n, const double *A, int lda,
                             double *x, int incx)
cusblasStatus_t cusblasCtrsv(cusblasHandle_t handle, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int n, const cuComplex *A, int lda,
                             cuComplex *x, int incx)
cusblasStatus_t cusblasZtrsv(cusblasHandle_t handle, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int n, const cuDoubleComplex *A, int lda,
                             cuDoubleComplex *x, int incx)

```

This function solves the triangular linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where  $A$  is a triangular matrix stored in lower or upper mode with or without the main diagonal, and  $\mathbf{x}$  and  $\mathbf{b}$  are vectors. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $\mathbf{x}$  overwrites the right-hand-sides  $\mathbf{b}$  on exit.

No test for singularity or near-singularity is included in this function.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	Number of rows and columns of matrix A.
A	device	input	<type> array of dimension lda x n, with lda>=max(1,n).
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	in/out	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>diag \neq \text{CUBLAS\_DIAG\_UNIT}, \text{CUBLAS\_DIAG\_NON\_UNIT}</math> or</li> <li>▶ <math>lda &lt; \max(1, n)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strsv](#), [dtrsv](#), [ctrsv](#), [ztrsv](#)

## 2.6.17. cublas<t>hemv()

```
cublasStatus_t cublasChemv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex *alpha,
```

```

        const cuComplex      *A, int lda,
        const cuComplex      *x, int incx,
        const cuComplex      *beta,
        cuComplex             *y, int incy)
cublasStatus_t cublasZhemv(cublasHandle_t handle, cublasFillMode_t uplo,
        int n, const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *beta,
        cuDoubleComplex *y, int incy)

```

This function performs the Hermitian matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $n \times n$  Hermitian matrix stored in lower or upper mode,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

This function has an alternate faster implementation using atomics that can be enabled with `cublasSetAtomicsMode()`.

Please see the section on the function `cublasSetAtomicsMode()` for more details about the usage of atomics.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times n$ , with $\text{lda} \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ then $\mathbf{y}$ does not have to be a valid input.
y	device	in/out	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $\mathbf{y}$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Error Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ <math>lda &lt; n</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[chemv](#), [zhemv](#)

## 2.6.18. `cublas<t>hbmV()`

```

cublasStatus_t cublasChbmV(cublasHandle_t handle, cublasFillMode_t uplo,
    int n, int k, const cuComplex *alpha,
    const cuComplex *A, int lda,
    const cuComplex *x, int incx,
    const cuComplex *beta,
    cuComplex *y, int incy)
cublasStatus_t cublasZhbmV(cublasHandle_t handle, cublasFillMode_t uplo,
    int n, int k, const cuDoubleComplex *alpha,
    const cuDoubleComplex *A, int lda,
    const cuDoubleComplex *x, int incx,
    const cuDoubleComplex *beta,
    cuDoubleComplex *y, int incy)

```

This function performs the Hermitian banded matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $n \times n$  Hermitian banded matrix with  $k$  subdiagonals and superdiagonals,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the Hermitian banded matrix  $\mathbf{A}$  is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [j, \min(m, j+k)]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the bottom right  $k \times k$  triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the Hermitian banded matrix  $\mathbf{A}$  is stored column by column, with the main diagonal of the matrix stored in row  $k+1$ , the first superdiagonal in row  $k$  (starting at second position), the second superdiagonal in row  $k-1$  (starting at third position), etc. So that in general, the element  $A(i, j)$  is stored in the memory location  $A(1+k+i-j, j)$  for  $j = 1, \dots, n$  and  $i \in [\max(1, j-k), j]$ . Also, the elements in the array  $A$  that do not conceptually correspond to the elements in the banded matrix (the top left  $k \times k$  triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
k		input	Number of sub- and super-diagonals of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimensions $lda \times n$ , with $lda \geq k+1$ . The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $x$ .
beta	host or device	input	<type> scalar used for multiplication, if $\beta=0$ then does not have to be a valid input.
y	device	in/out	<type> vector with $n$ elements.
incy		input	Stride between consecutive elements of $y$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>incx = 0</math> or <math>incy = 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>lda &lt; (k + 1)</math> or</li> <li>▶ <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[chbmv](#), [zhbmv](#)

## 2.6.19. cublas<t>hpmv()

```

cublasStatus_t cublasChpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex *alpha,
                           const cuComplex *AP,
                           const cuComplex *x, int incx,
                           const cuComplex *beta,
                           cuComplex *y, int incy)
cublasStatus_t cublasZhpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *AP,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the Hermitian packed matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where  $\mathbf{A}$  is a  $n \times n$  Hermitian matrix stored in packed format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  and  $\beta$  are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + ((2*n-j+1)*j)/2]$  for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix  $\mathbf{A}$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location  $AP[i + (j*(j+1))/2]$  for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $\mathbf{A}$ lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $\mathbf{A}$ .
alpha	host or device	input	<type> scalar used for multiplication.
AP	device	input	<type> array with $\mathbf{A}$ stored in packed format. The imaginary parts of the diagonal elements are assumed to be zero.
x	device	input	<type> vector with $n$ elements.
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then $\mathbf{y}$ does not have to be a valid input.
y	device	in/out	<type> vector with $n$ elements.



Param.	Memory	In/out	Meaning
incy		input	Stride between consecutive elements of $y$ .

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx == 0</math> or <math>incy == 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER, CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[chpmv](#), [zhpmv](#)

## 2.6.20. `cublas<t>her()`

```

cublasStatus_t cublasCher(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const cuComplex *x, int incx,
                          cuComplex *A, int lda)
cublasStatus_t cublasZher(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const cuDoubleComplex *x, int incx,
                          cuDoubleComplex *A, int lda)

```

This function performs the Hermitian rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^H + A$$

where  $A$  is a  $n \times n$  Hermitian matrix stored in column-major format,  $\mathbf{x}$  is a vector, and  $\alpha$  is a scalar.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix $A$ lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix $A$ .
alpha	host or device	input	<type> scalar used for multiplication.

Param.	Memory	In/out	Meaning
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
A	device	in/out	<type> array of dimensions lda x n, with lda>=max(1,n). The imaginary parts of the diagonal elements are assumed and set to zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If n &lt; 0 or</li> <li>▶ if incx == 0 or</li> <li>▶ if uplo != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER or</li> <li>▶ if lda &lt; max(1, n) or</li> <li>▶ alpha == NULL</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cher](#), [zher](#)

## 2.6.21. cublas<t>her2()

```

cublasStatus_t cublasCher2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex *alpha,
                           const cuComplex *x, int incx,
                           const cuComplex *y, int incy,
                           cuComplex *A, int lda)
cublasStatus_t cublasZher2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *y, int incy,
                           cuDoubleComplex *A, int lda)

```

This function performs the Hermitian rank-2 update

$$A = \alpha \mathbf{x} \mathbf{y}^H + \bar{\alpha} \mathbf{y} \mathbf{x}^H + A$$

where  $A$  is a  $n \times n$  Hermitian matrix stored in column-major format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  is a scalar.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.
A	device	in/out	<type> array of dimension lda x n with $lda \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx == 0</math> or <math>incy == 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER, CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> or</li> <li>▶ <math>alpha == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

cher2, zher2

## 2.6.22. cublas<t>hpr()

```
cublasStatus_t cublasChpr(cublasHandle_t handle, cublasFillMode_t uplo,
```

```

int n, const float *alpha,
const cuComplex *x, int incx,
cuComplex *AP)
cublasStatus_t cublasZhpr(cublasHandle_t handle, cublasFillMode_t uplo,
int n, const double *alpha,
const cuDoubleComplex *x, int incx,
cuDoubleComplex *AP)

```

This function performs the packed Hermitian rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^H + A$$

where  $A$  is a  $n \times n$  Hermitian matrix stored in packed format,  $\mathbf{x}$  is a vector, and  $\alpha$  is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + ((2*n-j+1)*j)/2]` for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + (j*(j+1))/2]` for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	Number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	Stride between consecutive elements of x.
AP	device	in/out	<type> array with A stored in packed format. The imaginary parts of the diagonal elements are assumed and set to zero.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx == 0</math> or</li> </ul>

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if <code>uplo != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER</code> or</li> <li>▶ <code>alpha == NULL</code></li> </ul>
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[chpr](#), [zhpr](#)

## 2.6.23. `cublas<type>hpr2()`

```

cublasStatus_t cublasChpr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex *alpha,
                           const cuComplex *x, int incx,
                           const cuComplex *y, int incy,
                           cuComplex *AP)
cublasStatus_t cublasZhpr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *y, int incy,
                           cuDoubleComplex *AP)

```

This function performs the packed Hermitian rank-2 update

$$A = \alpha \mathbf{x} \mathbf{y}^H + \bar{\alpha} \mathbf{y} \mathbf{x}^H + A$$

where  $A$  is a  $n \times n$  Hermitian matrix stored in packed format,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors, and  $\alpha$  is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + ((2*n-j+1)*j)/2]` for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + (j*(j+1))/2]` for  $j = 1, \dots, n$  and  $i \leq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>uplo</code>		input	Indicates if matrix $A$ lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
<code>n</code>		input	Number of rows and columns of matrix $A$ .
<code>alpha</code>	host or device	input	<type> scalar used for multiplication.

Param.	Memory	In/out	Meaning
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	Stride between consecutive elements of y.
AP	device	in/out	<type> array with A stored in packed format. The imaginary parts of the diagonal elements are assumed and set to zero.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>incx == 0</math> or <math>incy == 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER, CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ <math>alpha == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

chpr2, zhpr2

## 2.6.24. cublas<t>gemvBatched()

```

cublasStatus_t cublasSgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
    int m, int n,
    const float      *alpha,
    const float      *Aarray[], int lda,
    const float      *xarray[], int incx,
    const float      *beta,
    float            *yarray[], int incy,
    int batchCount)
cublasStatus_t cublasDgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
    int m, int n,
    const double     *alpha,
    const double     *Aarray[], int lda,
    const double     *xarray[], int incx,
    const double     *beta,
    double           *yarray[], int incy,
    int batchCount)
cublasStatus_t cublasCgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
    int m, int n,
    const cuComplex  *alpha,
    const cuComplex  *Aarray[], int lda,

```

```

const cuComplex      *xarray[], int incx,
const cuComplex      *beta,
cuComplex            *yarray[], int incy,
int batchCount)
cublasStatus_t cublasZgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
int m, int n,
const cuDoubleComplex *alpha,
const cuDoubleComplex *Aarray[], int lda,
const cuDoubleComplex *xarray[], int incx,
const cuDoubleComplex *beta,
cuDoubleComplex *yarray[], int incy,
int batchCount)
cublasStatus_t cublasHSHgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
int m, int n,
const float           *alpha,
const __half          *Aarray[], int lda,
const __half          *xarray[], int incx,
const float           *beta,
__half                *yarray[], int incy,
int batchCount)
cublasStatus_t cublasHSSgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
int m, int n,
const float           *alpha,
const __half          *Aarray[], int lda,
const __half          *xarray[], int incx,
const float           *beta,
float                 *yarray[], int incy,
int batchCount)
cublasStatus_t cublasTSTgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
int m, int n,
const float           *alpha,
const __nv_bfloat16   *Aarray[], int lda,
const __nv_bfloat16   *xarray[], int incx,
const float           *beta,
__nv_bfloat16         *yarray[], int incy,
int batchCount)
cublasStatus_t cublasTSSgemvBatched(cublasHandle_t handle, cublasOperation_t trans,
int m, int n,
const float           *alpha,
const __nv_bfloat16   *Aarray[], int lda,
const __nv_bfloat16   *xarray[], int incx,
const float           *beta,
float                 *yarray[], int incy,
int batchCount)

```

This function performs the matrix-vector multiplication of a batch of matrices and vectors. The batch is considered to be "uniform", i.e. all instances have the same dimensions ( $m$ ,  $n$ ), leading dimension ( $lda$ ), increments ( $incx$ ,  $incy$ ) and transposition ( $trans$ ) for their respective  $A$  matrix,  $x$  and  $y$  vectors. The address of the input matrix and vector, and the output vector of each instance of the batch are read from arrays of pointers passed to the function by the caller.

$$y[i] = \alpha op(A[i])x[i] + \beta y[i], \text{ for } i \in [0, batchCount - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$  is an array of pointers to matrices  $A[i]$  stored in column-major format with dimension  $m \times n$ , and  $x$  and  $y$  are arrays of pointers to vectors. Also, for matrix  $A[i]$ ,

$$op(A[i]) = \begin{cases} A[i] & \text{if } trans == CUBLAS\_OP\_N \\ A[i]^T & \text{if } trans == CUBLAS\_OP\_T \\ A[i]^H & \text{if } trans == CUBLAS\_OP\_C \end{cases}$$

Note:  $y[i]$  vectors must not overlap, that is, the individual gemv operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<t>gemv` in different CUDA streams, rather than use this API.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>trans</code>		input	Operation $op(A[i])$ that is non- or (conj.) transpose.
<code>m</code>		input	Number of rows of matrix $A[i]$ .
<code>n</code>		input	Number of columns of matrix $A[i]$ .
<code>alpha</code>	host or device	input	<type> scalar used for multiplication.
<code>Aarray</code>	device	input	Array of pointers to <type> array, with each array of dim. $lda \times n$ with $lda \geq \max(1, m)$ .  All pointers must meet certain alignment criteria. Please see below for details.
<code>lda</code>		input	Leading dimension of two-dimensional array used to store each matrix $A[i]$ .
<code>xarray</code>	device	input	Array of pointers to <type> array, with each dimension $n$ if <code>trans==CUBLAS_OP_N</code> and $m$ otherwise.  All pointers must meet certain alignment criteria. Please see below for details.
<code>incx</code>		input	Stride of each one-dimensional array $x[i]$ .
<code>beta</code>	host or device	input	<type> scalar used for multiplication. If <code>beta == 0</code> , $y$ does not have to be a valid input.
<code>yarray</code>	device	in/out	Array of pointers to <type> array. It has dimensions $m$ if <code>trans==CUBLAS_OP_N</code> and $n$ otherwise. Vectors $y[i]$ should not overlap; otherwise, undefined behavior is expected.  All pointers must meet certain alignment criteria. Please see below for details.
<code>incy</code>		input	Stride of each one-dimensional array $y[i]$ .
<code>batchCount</code>		input	Number of pointers contained in <code>Aarray</code> , <code>xarray</code> and <code>yarray</code> .

If math mode enables fast math modes when using `cublasSgemvBatched()`, pointers (not the pointer arrays) placed in the GPU memory must be properly aligned to avoid misaligned memory access errors. Ideally all pointers are aligned to at least 16 Bytes. Otherwise it is recommended that they meet the following rule:

- if  $k\%4==0$  then ensure `intptr_t(ptr) % 16 == 0`,



The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $m, n, \text{batchCount} < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 2.6.25. `cublas<T>gemvStridedBatched()`

```

cublasStatus_t cublasSgemvStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t trans,
                                         int m, int n,
                                         const float *alpha,
                                         const float *A, int lda,
                                         long long int strideA,
                                         const float *x, int incx,
                                         long long int stridex,
                                         const float *beta,
                                         float *y, int incy,
                                         long long int stridey,
                                         int batchCount)
cublasStatus_t cublasDgemvStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t trans,
                                         int m, int n,
                                         const double *alpha,
                                         const double *A, int lda,
                                         long long int strideA,
                                         const double *x, int incx,
                                         long long int stridex,
                                         const double *beta,
                                         double *yarray[], int incy,
                                         long long int stridey,
                                         int batchCount)
cublasStatus_t cublasCgemvStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t trans,
                                         int m, int n,
                                         const cuComplex *alpha,
                                         const cuComplex *A, int lda,
                                         long long int strideA,
                                         const cuComplex *x, int incx,
                                         long long int stridex,
                                         const cuComplex *beta,
                                         cuComplex *y, int incy,
                                         long long int stridey,
                                         int batchCount)
cublasStatus_t cublasZgemvStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t trans,
                                         int m, int n,
                                         const cuDoubleComplex *alpha,
                                         const cuDoubleComplex *A, int lda,
                                         long long int strideA,
                                         const cuDoubleComplex *x, int incx,
                                         long long int stridex,
                                         const cuDoubleComplex *beta,
                                         cuDoubleComplex *y, int incy,
                                         long long int stridey,
                                         int batchCount)
cublasStatus_t cublasHSHgemvStridedBatched(cublasHandle_t handle,

```

```

cublasOperation_t trans,
int m, int n,
const float          *alpha,
const __half         *A, int lda,
long long int        strideA,
const __half         *x, int incx,
long long int        stridex,
const float          *beta,
__half              *y, int incy,
long long int        stridey,
int batchCount)
cublasStatus_t cublasHSSgemvStridedBatched(cublasHandle_t handle,
cublasOperation_t trans,
int m, int n,
const float          *alpha,
const __half         *A, int lda,
long long int        strideA,
const __half         *x, int incx,
long long int        stridex,
const float          *beta,
float                *y, int incy,
long long int        stridey,
int batchCount)
cublasStatus_t cublasTSTgemvStridedBatched(cublasHandle_t handle,
cublasOperation_t trans,
int m, int n,
const float          *alpha,
const __nv_bfloat16  *A, int lda,
long long int        strideA,
const __nv_bfloat16  *x, int incx,
long long int        stridex,
const float          *beta,
__nv_bfloat16        *y, int incy,
long long int        stridey,
int batchCount)
cublasStatus_t cublasTSSgemvStridedBatched(cublasHandle_t handle,
cublasOperation_t trans,
int m, int n,
const float          *alpha,
const __nv_bfloat16  *A, int lda,
long long int        strideA,
const __nv_bfloat16  *x, int incx,
long long int        stridex,
const float          *beta,
float                *y, int incy,
long long int        stridey,
int batchCount)

```

This function performs the matrix-vector multiplication of a batch of matrices and vectors. The batch is considered to be "uniform", that is, all instances have the same dimensions ( $m$ ,  $n$ ), leading dimension ( $lda$ ), increments ( $incx$ ,  $incy$ ) and transposition ( $trans$ ) for their respective  $A$  matrix,  $x$  and  $y$  vectors. Input matrix  $A$  and vector  $x$ , and output vector  $y$  for each instance of the batch are located at fixed offsets in number of elements from their locations in the previous instance. Pointers to  $A$  matrix,  $x$  and  $y$  vectors for the first instance are passed to the function by the user along with offsets in number of elements -  $strideA$ ,  $stridex$  and  $stridey$  that determine the locations of input matrices and vectors, and output vectors in future instances.

$$\mathbf{y} + i * \text{stridey} = \alpha \text{op}(\mathbf{A} + i * \text{strideA})(\mathbf{x} + i * \text{stridex}) + \beta(\mathbf{y} + i * \text{stridey}), \text{ for } i \in [0, \text{batchCount} - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$  is an array of pointers to matrix stored in column-major format with dimension  $A[i] \ m \times n$ , and  $x$  and  $y$  are arrays of pointers to vectors. Also, for matrix  $A[i]$

$$\text{op}(A[i]) = \begin{cases} A[i] & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A[i]^T & \text{if trans} == \text{CUBLAS\_OP\_T} \\ A[i]^H & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

Note:  $y[i]$  matrices must not overlap, i.e. the individual gemv operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<t>gemv` in different CUDA streams, rather than use this API.

Note: In the table below, we use  $A[i]$ ,  $x[i]$ ,  $y[i]$  as notation for A matrix, and x and y vectors in the  $i$ th instance of the batch, implicitly assuming they are respectively offsets in number of elements `strideA`, `stridex`, `stridey` away from  $A[i-1]$ ,  $x[i-1]$ ,  $y[i-1]$ . The unit for the offset is number of elements and must not be zero.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
trans		input	Operation $\text{op}(A[i])$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $A[i]$ .
n		input	Number of columns of matrix $A[i]$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type>* pointer to the A matrix corresponding to the first instance of the batch, with dimensions <code>lda x n</code> with <code>lda &gt;= max(1, m)</code> .
lda		input	Leading dimension of two-dimensional array used to store each matrix $A[i]$ .
strideA		input	Value of type <code>long long int</code> that gives the offset in number of elements between $A[i]$ and $A[i+1]$
x	device	input	<type>* pointer to the x vector corresponding to the first instance of the batch, with each dimension <code>n</code> if <code>trans == CUBLAS_OP_N</code> and <code>m</code> otherwise.
incx		input	Stride of each one-dimensional array $x[i]$ .
stridex		input	Value of type <code>long long int</code> that gives the offset in number of elements between $x[i]$ and $x[i+1]$
beta	host or device	input	<type> scalar used for multiplication. If <code>beta == 0</code> , $y$ does not have to be a valid input.
y	device	in/out	<type>* pointer to the y vector corresponding to the first instance of the batch, with each dimension <code>m</code> if <code>trans == CUBLAS_OP_N</code> and <code>n</code> otherwise. Vectors $y[i]$ should not overlap; otherwise, undefined behavior is expected.
incy		input	Stride of each one-dimensional array $y[i]$ .

Param.	Memory	In/out	Meaning
stridey		input	Value of type long long int that gives the offset in number of elements between y[i] and y[i+1]
batchCount		input	Number of GEMVs to perform in the batch.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters m,n,batchCount<0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 2.7. cuBLAS Level-3 Function Reference

This section describes the Level-3 Basic Linear Algebra Subprograms (BLAS3) functions that perform matrix-matrix operations.

### 2.7.1. cublas<t>gemm()

```

cublasStatus_t cublasSgemm(cublasHandle_t handle,
                          cublasOperation_t transa, cublasOperation_t transb,
                          int m, int n, int k,
                          const float *alpha,
                          const float *A, int lda,
                          const float *B, int ldb,
                          const float *beta,
                          float *C, int ldc)
cublasStatus_t cublasDgemm(cublasHandle_t handle,
                          cublasOperation_t transa, cublasOperation_t transb,
                          int m, int n, int k,
                          const double *alpha,
                          const double *A, int lda,
                          const double *B, int ldb,
                          const double *beta,
                          double *C, int ldc)
cublasStatus_t cublasCgemm(cublasHandle_t handle,
                          cublasOperation_t transa, cublasOperation_t transb,
                          int m, int n, int k,
                          const cuComplex *alpha,
                          const cuComplex *A, int lda,
                          const cuComplex *B, int ldb,
                          const cuComplex *beta,
                          cuComplex *C, int ldc)
cublasStatus_t cublasZgemm(cublasHandle_t handle,
                          cublasOperation_t transa, cublasOperation_t transb,
                          int m, int n, int k,
                          const cuDoubleComplex *alpha,
                          const cuDoubleComplex *A, int lda,
                          const cuDoubleComplex *B, int ldb,
                          const cuDoubleComplex *beta,

```

```

        cuDoubleComplex *C, int ldc)
cublasStatus_t cublasHgemm(cublasHandle_t handle,
        cublasOperation_t transa, cublasOperation_t transb,
        int m, int n, int k,
        const __half *alpha,
        const __half *A, int lda,
        const __half *B, int ldb,
        const __half *beta,
        __half *C, int ldc)

```

This function performs the matrix-matrix multiplication

$$C = \alpha \text{op}(A)\text{op}(B) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times k$ ,  $\text{op}(B) k \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation op(A) that is non- or (conj.) transpose.
transb		input	Operation op(B) that is non- or (conj.) transpose.
m		input	Number of rows of matrix op(A) and C.
n		input	Number of columns of matrix op(B) and C.
k		input	Number of columns of op(A) and rows of op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimensions lda x k with lda>=max(1,m) if transa == CUBLAS_OP_N and lda x m with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension ldb x n with ldb>=max(1,k) if transb == CUBLAS_OP_N and ldb x k with ldb>=max(1,n) otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication. If beta==0, C does not have to be a valid input.
C	device	in/out	<type> array of dimensions ldc x n with ldc>=max(1,m).

Param.	Memory	In/out	Meaning
ldc		input	Leading dimension of a two-dimensional array used to store the matrix c.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If m, n, k &lt; 0 or</li> <li>▶ if transa, transb != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if lda &lt; max(1, m) if transa == CUBLAS_OP_N and lda &lt; max(1, k) otherwise or</li> <li>▶ if ldb &lt; max(1, k) if transb == CUBLAS_OP_N and ldb &lt; max(1, n) otherwise or</li> <li>▶ if ldc &lt; max(1, m) or</li> <li>▶ if alpha, beta == NULL or</li> <li>▶ c == NULL if c needs to be scaled</li> </ul>
CUBLAS_STATUS_ARCH_MISMATCH	In the case of cublasHgemm the device does not support math in half precision.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgemm](#), [dgemm](#), [cgemm](#), [zgemm](#)

## 2.7.2. cublas<t>gemm3m()

```

cublasStatus_t cublasCgemm3m(cublasHandle_t handle,
                             cublasOperation_t transa, cublasOperation_t transb,
                             int m, int n, int k,
                             const cuComplex *alpha,
                             const cuComplex *A, int lda,
                             const cuComplex *B, int ldb,
                             const cuComplex *beta,
                             cuComplex *C, int ldc)
cublasStatus_t cublasZgemm3m(cublasHandle_t handle,
                             cublasOperation_t transa, cublasOperation_t transb,
                             int m, int n, int k,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, int lda,
                             const cuDoubleComplex *B, int ldb,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex *C, int ldc)

```

This function performs the complex matrix-matrix multiplication, using Gauss complexity reduction algorithm. This can lead to an increase in performance up to 25%

$$C = \alpha \text{op}(A)\text{op}(B) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times k$ ,  $\text{op}(B) k \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .



**Note:** These 2 routines are only supported on GPUs with architecture capabilities equal to or greater than 5.0.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation op(A) that is non- or (conj.) transpose.
transb		input	Operation op(B) that is non- or (conj.) transpose.
m		input	Number of rows of matrix op(A) and c.
n		input	Number of columns of matrix op(B) and c.
k		input	Number of columns of op(A) and rows of op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimensions $lda \times k$ with $lda \geq \max(1, m)$ if $transa == \text{CUBLAS\_OP\_N}$ and $lda \times m$ with $lda \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, k)$ if $transb == \text{CUBLAS\_OP\_N}$ and $ldb \times k$ with $ldb \geq \max(1, n)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication. If $\beta == 0$ , C does not have to be a valid input.
C	device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of a two-dimensional array used to store the matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n, k &lt; 0</math> or</li> <li>▶ if <math>transa, transb \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if <math>transa == \text{CUBLAS\_OP\_N}</math> and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, k)</math> if <math>transb == \text{CUBLAS\_OP\_N}</math> and <math>ldb &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, m)</math> or</li> <li>▶ if <math>\alpha, \beta == \text{NULL}</math> or</li> <li>▶ <math>c == \text{NULL}</math> if <math>c</math> needs to be scaled</li> </ul>
CUBLAS_STATUS_ARCH_MISMATCH	The device has a compute capability lower than 5.0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cgemm](#), [zgemm](#)

### 2.7.3. cublas<t>gemmBatched()

```

cublasStatus_t cublasHgemmBatched(cublasHandle_t handle,
                                   cublasOperation_t transa,
                                   cublasOperation_t transb,
                                   int m, int n, int k,
                                   const __half *alpha,
                                   const __half *Aarray[], int lda,
                                   const __half *Barray[], int ldb,
                                   const __half *beta,
                                   __half *Carray[], int ldc,
                                   int batchCount)
cublasStatus_t cublasSgemmBatched(cublasHandle_t handle,
                                   cublasOperation_t transa,
                                   cublasOperation_t transb,
                                   int m, int n, int k,
                                   const float *alpha,
                                   const float *Aarray[], int lda,
                                   const float *Barray[], int ldb,
                                   const float *beta,
                                   float *Carray[], int ldc,
                                   int batchCount)
cublasStatus_t cublasDgemmBatched(cublasHandle_t handle,
                                   cublasOperation_t transa,
                                   cublasOperation_t transb,

```



```

        int m, int n, int k,
        const double      *alpha,
        const double      *Aarray[], int lda,
        const double      *Barray[], int ldb,
        const double      *beta,
        double             *Carray[], int ldc,
        int batchCount)
cublasStatus_t cublasCgemvBatched(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m, int n, int k,
        const cuComplex   *alpha,
        const cuComplex   *Aarray[], int lda,
        const cuComplex   *Barray[], int ldb,
        const cuComplex   *beta,
        cuComplex         *Carray[], int ldc,
        int batchCount)
cublasStatus_t cublasZgemvBatched(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m, int n, int k,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *Aarray[], int lda,
        const cuDoubleComplex *Barray[], int ldb,
        const cuDoubleComplex *beta,
        cuDoubleComplex       *Carray[], int ldc,
        int batchCount)

```

This function performs the matrix-matrix multiplication of a batch of matrices. The batch is considered to be "uniform", that is, all instances have the same dimensions ( $m, n, k$ ), leading dimensions ( $lda, ldb, ldc$ ) and transpositions ( $transa, transb$ ) for their respective  $A, B$  and  $C$  matrices. The address of the input matrices and the output matrix of each instance of the batch are read from arrays of pointers passed to the function by the caller.

$$C[i] = \alpha \text{op}(A[i]) \text{op}(B[i]) + \beta C[i], \text{ for } i \in [0, \text{batchCount} - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A, B$  and  $C$  are arrays of pointers to matrices stored in column-major format with dimensions  $\text{op}(A[i])\ m \times k$ ,  $\text{op}(B[i])\ k \times n$  and  $C[i]\ m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if } transa == \text{CUBLAS\_OP\_N} \\ A^T & \text{if } transa == \text{CUBLAS\_OP\_T} \\ A^H & \text{if } transa == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B[i])$  is defined similarly for matrix  $B[i]$ .

Note:  $C[i]$  matrices must not overlap, that is, the individual gemm operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<t>gemm` in different CUDA streams, rather than use this API.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}(A[i])$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}(B[i])$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}(A[i])$ and $C[i]$ .

Param.	Memory	In/out	Meaning
n		input	Number of columns of $op[B[i]]$ and $C[i]$ .
k		input	Number of columns of $op[A[i]]$ and rows of $op[B[i]]$ .
alpha	host or device	input	<type> scalar used for multiplication.
Aarray	device	input	Array of pointers to <type> array, with each array of dim. $lda \times k$ with $lda \geq \max(1, m)$ if $transa == CUBLAS\_OP\_N$ and $lda \times m$ with $lda \geq \max(1, k)$ otherwise.  All pointers must meet certain alignment criteria. Please see below for details.
lda		input	Leading dimension of two-dimensional array used to store each matrix $A[i]$ .
Barray	device	input	Array of pointers to <type> array, with each array of dim. $ldb \times n$ with $ldb \geq \max(1, k)$ if $transb == CUBLAS\_OP\_N$ and $ldb \times k$ with $ldb \geq \max(1, n)$ otherwise.  All pointers must meet certain alignment criteria. Please see below for details.
ldb		input	Leading dimension of two-dimensional array used to store each matrix $B[i]$ .
beta	host or device	input	<type> scalar used for multiplication. If $beta == 0$ , $C$ does not have to be a valid input.
Carray	device	in/out	Array of pointers to <type> array. It has dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ . Matrices $C[i]$ should not overlap; otherwise, undefined behavior is expected.  All pointers must meet certain alignment criteria. Please see below for details.
ldc		input	Leading dimension of two-dimensional array used to store each matrix $C[i]$ .
batchCount		input	Number of pointers contained in Aarray, Barray and Carray.

If math mode enables fast math modes when using `cublasSgemmBatched()`, pointers (not the pointer arrays) placed in the GPU memory must be properly aligned to avoid misaligned memory access errors. Ideally all pointers are aligned to at least 16 Bytes. Otherwise it is recommended that they meet the following rule:

- if  $k \% 4 == 0$  then ensure  $\text{intptr\_t}(ptr) \% 16 == 0$ ,

The possible error values returned by this function and their meanings are listed the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Error Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n, k, \text{batchCount} &lt; 0</math> or</li> <li>▶ if <math>\text{transa}, \text{transb} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, m)</math> if <math>\text{transa} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, k)</math> if <math>\text{transb} == \text{CUBLAS\_OP\_N}</math> and <math>\text{ldb} &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, m)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_ARCH_MISMATCH	cublasHgemmBatched is only supported for GPU with architecture capabilities equal to or greater than 5.3.

## 2.7.4. cublas<t>gemmStridedBatched()

```

cublasStatus_t cublasHgemmStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t transa,
                                         cublasOperation_t transb,
                                         int m, int n, int k,
                                         const __half *alpha,
                                         const __half *A, int lda,
                                         long long int strideA,
                                         const __half *B, int ldb,
                                         long long int strideB,
                                         const __half *beta,
                                         __half *C, int ldc,
                                         long long int strideC,
                                         int batchCount)
cublasStatus_t cublasSgemmStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t transa,
                                         cublasOperation_t transb,
                                         int m, int n, int k,
                                         const float *alpha,
                                         const float *A, int lda,
                                         long long int strideA,
                                         const float *B, int ldb,
                                         long long int strideB,
                                         const float *beta,
                                         float *C, int ldc,
                                         long long int strideC,
                                         int batchCount)
cublasStatus_t cublasDgemmStridedBatched(cublasHandle_t handle,
                                         cublasOperation_t transa,
                                         cublasOperation_t transb,
                                         int m, int n, int k,
                                         const double *alpha,
                                         const double *A, int lda,
                                         long long int strideA,
                                         const double *B, int ldb,
                                         long long int strideB,
                                         const double *beta,
                                         double *C, int ldc,
                                         long long int strideC,
                                         int batchCount)

```

```

        long long int      strideC,
        int batchCount)
cublasStatus_t cublasCgemmStridedBatched(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m, int n, int k,
        const cuComplex *alpha,
        const cuComplex *A, int lda,
        long long int strideA,
        const cuComplex *B, int ldb,
        long long int strideB,
        const cuComplex *beta,
        cuComplex *C, int ldc,
        long long int strideC,
        int batchCount)
cublasStatus_t cublasCgemm3mStridedBatched(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m, int n, int k,
        const cuComplex *alpha,
        const cuComplex *A, int lda,
        long long int strideA,
        const cuComplex *B, int ldb,
        long long int strideB,
        const cuComplex *beta,
        cuComplex *C, int ldc,
        long long int strideC,
        int batchCount)
cublasStatus_t cublasZgemmStridedBatched(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m, int n, int k,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        long long int strideA,
        const cuDoubleComplex *B, int ldb,
        long long int strideB,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc,
        long long int strideC,
        int batchCount)

```

This function performs the matrix-matrix multiplication of a batch of matrices. The batch is considered to be "uniform", i.e. all instances have the same dimensions ( $m, n, k$ ), leading dimensions ( $lda, ldb, ldc$ ) and transpositions ( $transa, transb$ ) for their respective  $A, B$  and  $C$  matrices. Input matrices  $A, B$  and output matrix  $C$  for each instance of the batch are located at fixed offsets in number of elements from their locations in the previous instance. Pointers to  $A, B$  and  $C$  matrices for the first instance are passed to the function by the user along with offsets in number of elements -  $strideA, strideB$  and  $strideC$  that determine the locations of input and output matrices in future instances.

$$C + i * strideC = \alpha op(A + i * strideA) op(B + i * strideB) + \beta (C + i * strideC), \text{ for } i \in [0, batchCount - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A, B$  and  $C$  are arrays of pointers to matrices stored in column-major format with dimensions  $op(A[i]) \ m \times k$ ,  $op(B[i]) \ k \times n$  and  $C[i] \ m \times n$ , respectively. Also, for matrix  $A$

$$op(A) = \begin{cases} A & \text{if } transa == CUBLAS\_OP\_N \\ A^T & \text{if } transa == CUBLAS\_OP\_T \\ A^H & \text{if } transa == CUBLAS\_OP\_C \end{cases}$$

and  $op(B[i])$  is defined similarly for matrix  $B[i]$ .

Note:  $C[i]$  matrices must not overlap, i.e. the individual gemm operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<t>gemm` in different CUDA streams, rather than use this API.

Note: In the table below, we use  $A[i]$ ,  $B[i]$ ,  $C[i]$  as notation for A, B and C matrices in the  $i$ th instance of the batch, implicitly assuming they are respectively offsets in number of elements `strideA`, `strideB`, `strideC` away from  $A[i-1]$ ,  $B[i-1]$ ,  $C[i-1]$ . The unit for the offset is number of elements and must not be zero.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $op(A[i])$ that is non- or (conj.) transpose.
transb		input	Operation $op(B[i])$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $op(A[i])$ and $C[i]$ .
n		input	Number of columns of $op(B[i])$ and $C[i]$ .
k		input	Number of columns of $op(A[i])$ and rows of $op(B[i])$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type>* pointer to the A matrix corresponding to the first instance of the batch, with dimensions $lda \times k$ with $lda \geq \max(1, m)$ if $transa == CUBLAS\_OP\_N$ and $lda \times m$ with $lda \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store each matrix $A[i]$ .
strideA		input	Value of type long long int that gives the offset in number of elements between $A[i]$ and $A[i+1]$
B	device	input	<type>* pointer to the B matrix corresponding to the first instance of the batch, with dimensions $ldb \times n$ with $ldb \geq \max(1, k)$ if $transb == CUBLAS\_OP\_N$ and $ldb \times k$ with $ldb \geq \max(1, n)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store each matrix $B[i]$ .
strideB		input	Value of type long long int that gives the offset in number of elements between $B[i]$ and $B[i+1]$
beta	host or device	input	<type> scalar used for multiplication. If $beta == 0$ , $C$ does not have to be a valid input.
C	device	in/out	<type>* pointer to the C matrix corresponding to the first instance of the batch, with dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ . Matrices $C[i]$ should not overlap; otherwise, undefined behavior is expected.

Param.	Memory	In/out	Meaning
ldc		input	Leading dimension of two-dimensional array used to store each matrix $C[i]$ .
strideC		input	Value of type long long int that gives the offset in number of elements between $C[i]$ and $C[i+1]$
batchCount		input	Number of GEMMs to perform in the batch.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n, k, \text{batchCount} &lt; 0</math> or</li> <li>▶ if <math>\text{transa}, \text{transb} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, m)</math> if <math>\text{transa} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, k)</math> if <math>\text{transb} == \text{CUBLAS\_OP\_N}</math> and <math>\text{ldb} &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, m)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_ARCH_MISMATCH	<code>cublasHgemvStridedBatched</code> is only supported for GPU with architecture capabilities equal to or greater than 5.3.

## 2.7.5. `cublas<t>symm()`

```

cublasStatus_t cublasSsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const float *alpha,
                           const float *A, int lda,
                           const float *B, int ldb,
                           const float *beta,
                           float *C, int ldc)
cublasStatus_t cublasDsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const double *alpha,
                           const double *A, int lda,
                           const double *B, int ldb,
                           const double *beta,
                           double *C, int ldc)
cublasStatus_t cublasCsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,

```

```

        const cuComplex      *alpha,
        const cuComplex      *A, int lda,
        const cuComplex      *B, int ldb,
        const cuComplex      *beta,
        cuComplex             *C, int ldc)
cublasStatus_t cublasZsymm(cublasHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        int m, int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *B, int ldb,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc)

```

This function performs the symmetric matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side == CUBLAS_SIDE_LEFT} \\ \alpha BA + \beta C & \text{if side == CUBLAS_SIDE_RIGHT} \end{cases}$$

where  $A$  is a symmetric matrix stored in lower or upper mode,  $B$  and  $C$  are  $m \times n$  matrices, and  $\alpha$  and  $\beta$  are scalars.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
m		input	Number of rows of matrix C and B, with matrix A sized accordingly.
n		input	Number of columns of matrix C and B, with matrix A sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times m$ with $lda \geq \max(1, m)$ if side == CUBLAS_SIDE_LEFT and $lda \times n$ with $lda \geq \max(1, n)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if beta == 0 then C does not have to be a valid input.
C	device	in/out	<type> array of dimension $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n &lt; 0</math> or</li> <li>▶ if <math>side \neq \text{CUBLAS\_SIDE\_LEFT}, \text{CUBLAS\_SIDE\_RIGHT}</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>lda &lt; \max\{1, m\}</math> if <math>side == \text{CUBLAS\_SIDE\_LEFT}</math> and <math>lda &lt; \max\{1, n\}</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max\{1, m\}</math> or</li> <li>▶ if <math>ldc &lt; \max\{1, m\}</math> or</li> <li>▶ if <math>alpha == \text{NULL}</math> or <math>beta == \text{NULL}</math> or</li> <li>▶ <math>c == \text{NULL}</math> if <math>c</math> needs to be scaled</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssymm](#), [dsymm](#), [csymm](#), [zsymm](#)

## 2.7.6. cublas<t>syrk()

```

cublasStatus_t cublasSsyrk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float          *alpha,
                           const float          *A, int lda,
                           const float          *beta,
                           float                *C, int ldc)
cublasStatus_t cublasDsyrk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const double         *alpha,
                           const double         *A, int lda,
                           const double         *beta,
                           double               *C, int ldc)
cublasStatus_t cublasCsyrk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *beta,
                           cuComplex            *C, int ldc)
cublasStatus_t cublasZsyrk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,

```



```
int n, int k,
const cuDoubleComplex *alpha,
const cuDoubleComplex *A, int lda,
const cuDoubleComplex *beta,
cuDoubleComplex *C, int ldc)
```

This function performs the symmetric rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \ n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix c lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op[A] that is non- or transpose.
n		input	Number of rows of matrix op[A] and c.
k		input	Number of columns of matrix op[A].
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension lda x k with lda>=max(1,n) if trans == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta	host or device	input	<type> scalar used for multiplication, if beta==0 then c does not have to be a valid input.
C	device	in/out	<type> array of dimension ldc x n, with ldc>=max(1,n).
ldc		input	Leading dimension of two-dimensional array used to store matrix c.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	► If n, k < 0 or

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if trans != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if uplo != CUBLAS_FILL_MODE_LOWER, CUBLAS_FILL_MODE_UPPER or</li> <li>▶ if lda &lt; max{1, n} if trans == CUBLAS_OP_N and lda &lt; max{1, k} otherwise or</li> <li>▶ if ldc &lt; max{1, n} or</li> <li>▶ if alpha == NULL or beta == NULL or</li> <li>▶ c == NULL if c needs to be scaled</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#)

## 2.7.7. `cublas<t>syr2k()`

```

cublasStatus_t cublasSsyr2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float      *alpha,
                           const float      *A, int lda,
                           const float      *B, int ldb,
                           const float      *beta,
                           float             *C, int ldc)
cublasStatus_t cublasDsyr2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const double     *alpha,
                           const double     *A, int lda,
                           const double     *B, int ldb,
                           const double     *beta,
                           double            *C, int ldc)
cublasStatus_t cublasCsyr2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex  *alpha,
                           const cuComplex  *A, int lda,
                           const cuComplex  *B, int ldb,
                           const cuComplex  *beta,
                           cuComplex        *C, int ldc)
cublasStatus_t cublasZsyr2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs the symmetric rank-  $2k$  update

$$C = \alpha(\text{op}(A)\text{op}(B)^T + \text{op}(B)\text{op}(A)^T) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^T \text{ and } B^T & \text{if trans} == \text{CUBLAS\_OP\_T} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix c lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times k$ with $\text{lda} \geq \max(1, n)$ if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $\text{lda} \times n$ with $\text{lda} \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimensions $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ if $\text{transb} == \text{CUBLAS\_OP\_N}$ and $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ , then c does not have to be a valid input.
C	device	in/out	<type> array of dimensions $\text{ldc} \times n$ with $\text{ldc} \geq \max(1, n)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	► If $n, k < 0$ or

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if trans != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if uplo != CUBLAS_FILL_MODE_LOWER, CUBLAS_FILL_MODE_UPPER or</li> <li>▶ if lda &lt; max(1, n) if trans == CUBLAS_OP_N and lda &lt; max(1, k) otherwise or</li> <li>▶ if ldb &lt; max(1, n) if trans == CUBLAS_OP_N and lda &lt; max(1, k) otherwise or</li> <li>▶ if ldc &lt; max(1, n) or</li> <li>▶ if alpha == NULL or beta == NULL or</li> <li>▶ c == NULL if c needs to be scaled</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyr2k](#), [dsyr2k](#), [csyr2k](#), [zsyr2k](#)

## 2.7.8. `cublas<t>syrkx()`

```

cublasStatus_t cublasSsyrkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float          *alpha,
                           const float          *A, int lda,
                           const float          *B, int ldb,
                           const float          *beta,
                           float                *C, int ldc)
cublasStatus_t cublasDsyrkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const double         *alpha,
                           const double         *A, int lda,
                           const double         *B, int ldb,
                           const double         *beta,
                           double               *C, int ldc)
cublasStatus_t cublasCsyrkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *B, int ldb,
                           const cuComplex      *beta,
                           cuComplex            *C, int ldc)
cublasStatus_t cublasZsyrkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs a variation of the symmetric rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(B)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrices  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^T \text{ and } B^T & \text{if trans} == \text{CUBLAS\_OP\_T} \end{cases}$$

This routine can be used when B is in such way that the result is guaranteed to be symmetric. A usual example is when the matrix B is a scaled form of the matrix A : this is equivalent to B being the product of the matrix A and a diagonal matrix. For an efficient computation of the product of a regular matrix with a diagonal matrix, refer to the routine `cublas<t>dggmm`.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix c lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension lda x k with lda>=max(1,n) if transa == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimensions ldb x k with ldb>=max(1,n) if transb == CUBLAS_OP_N and ldb x n with ldb>=max(1,k) otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if beta==0, then c does not have to be a valid input.
C	device	in/out	<type> array of dimensions ldc x n with ldc>=max(1,n).
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n, k &lt; 0</math> or</li> <li>▶ if <math>\text{trans} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{uplo} \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, n)</math> or</li> <li>▶ if <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math> or</li> <li>▶ <math>c == \text{NULL}</math> if <math>c</math> needs to be scaled</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#) and  
[ssyr2k](#), [dsyr2k](#), [csyr2k](#), [zsyr2k](#)

## 2.7.9. `cusblas<t>trmm()`

```

cusblasStatus_t cusblasStrmm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const float *alpha,
                             const float *A, int lda,
                             const float *B, int ldb,
                             float *C, int ldc)
cusblasStatus_t cusblasDtrmm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const double *alpha,
                             const double *A, int lda,
                             const double *B, int ldb,
                             double *C, int ldc)
cusblasStatus_t cusblasCtrmm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const cuComplex *alpha,
                             const cuComplex *A, int lda,
                             const cuComplex *B, int ldb,
                             cuComplex *C, int ldc)
cusblasStatus_t cusblasZtrmm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const cuComplex *alpha,
                             const cuComplex *A, int lda,
                             const cuComplex *B, int ldb,
                             cuComplex *C, int ldc)

```

```

cublasSideMode_t side, cublasFillMode_t uplo,
cublasOperation_t trans, cublasDiagType_t diag,
int m, int n,
const cuDoubleComplex *alpha,
const cuDoubleComplex *A, int lda,
const cuDoubleComplex *B, int ldb,
cuDoubleComplex *C, int ldc)

```

This function performs the triangular matrix-matrix multiplication

$$C = \begin{cases} \alpha \text{op}(A)B & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ \alpha B \text{op}(A) & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a triangular matrix stored in lower or upper mode with or without the main diagonal,  $B$  and  $C$  are  $m \times n$  matrix, and  $\alpha$  is a scalar. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

Notice that in order to achieve better parallelism cuBLAS differs from the BLAS API only for this routine. The BLAS API assumes an in-place implementation (with results written back to B), while the cuBLAS API assumes an out-of-place implementation (with results written into C). The application can obtain the in-place functionality of BLAS in the cuBLAS API by passing the address of the matrix B in place of the matrix C. No other overlapping in the input parameters is supported.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	Number of rows of matrix B, with matrix A sized accordingly.
n		input	Number of columns of matrix B, with matrix A sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication, if alpha==0 then A is not referenced and B does not have to be a valid input.
A	device	input	<type> array of dimension lda x m with lda>=max(1,m) if side == CUBLAS_SIDE_LEFT and lda x n with lda>=max(1,n) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension ldb x n with ldb>=max(1,m).

Param.	Memory	In/out	Meaning
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
C	device	in/out	<type> array of dimension ldc x n with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m, n &lt; 0</math> or</li> <li>▶ if <math>trans \neq CUBLAS\_OP\_N, CUBLAS\_OP\_C, CUBLAS\_OP\_T</math> or</li> <li>▶ if <math>uplo \neq CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER</math> or</li> <li>▶ if <math>side \neq CUBLAS\_SIDE\_LEFT, CUBLAS\_SIDE\_RIGHT</math> or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if <math>side == CUBLAS\_SIDE\_LEFT</math> and <math>lda &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, m)</math> or</li> <li>▶ <math>c == NULL</math> if <math>c</math> needs to be scaled</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strmm](#), [dtrmm](#), [ctrmm](#), [ztrmm](#)

## 2.7.10. `cusblas<t>trsm()`

```

cusblasStatus_t cusblasStrsm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const float *alpha,
                             const float *A, int lda,
                             float *B, int ldb)
cusblasStatus_t cusblasDtrsm(cusblasHandle_t handle,
                             cusblasSideMode_t side, cusblasFillMode_t uplo,
                             cusblasOperation_t trans, cusblasDiagType_t diag,
                             int m, int n,
                             const double *alpha,
```



```

const double      *A, int lda,
double            *B, int ldb)
cublasStatus_t cublasCtrsm(cublasHandle_t handle,
                          cublasSideMode_t side, cublasFillMode_t uplo,
                          cublasOperation_t trans, cublasDiagType_t diag,
                          int m, int n,
                          const cuComplex      *alpha,
                          const cuComplex      *A, int lda,
                          cuComplex            *B, int ldb)
cublasStatus_t cublasZtrsm(cublasHandle_t handle,
                          cublasSideMode_t side, cublasFillMode_t uplo,
                          cublasOperation_t trans, cublasDiagType_t diag,
                          int m, int n,
                          const cuDoubleComplex *alpha,
                          const cuDoubleComplex *A, int lda,
                          cuDoubleComplex      *B, int ldb)

```

This function solves the triangular linear system with multiple right-hand-sides

$$\begin{cases} \text{op}(A)X = \alpha B & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ X\text{op}(A) = \alpha B & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a triangular matrix stored in lower or upper mode with or without the main diagonal,  $X$  and  $B$  are  $m \times n$  matrices, and  $\alpha$  is a scalar. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $X$  overwrites the right-hand-sides  $B$  on exit.

No test for singularity or near-singularity is included in this function.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
side		input	Indicates if matrix A is on the left or right of x.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	Number of rows of matrix B, with matrix A sized accordingly.
n		input	Number of columns of matrix B, with matrix A is sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication, if alpha==0 then A is not referenced and B does not have to be a valid input.
A	device	input	<type> array of dimension lda x m with lda>=max(1,m) if side == CUBLAS_SIDE_LEFT and lda x n with lda>=max(1,n) otherwise.

Param.	Memory	In/out	Meaning
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	in/out	<type> array. It has dimensions $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or</li> <li>▶ if <math>trans \neq CUBLAS\_OP\_N, CUBLAS\_OP\_C, CUBLAS\_OP\_T</math> or</li> <li>▶ if <math>uplo \neq CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER</math> or</li> <li>▶ if <math>side \neq CUBLAS\_SIDE\_LEFT, CUBLAS\_SIDE\_RIGHT</math> or</li> <li>▶ if <math>diag \neq CUBLAS\_DIAG\_NON\_UNIT, CUBLAS\_DIAG\_UNIT</math> or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if <math>side == CUBLAS\_SIDE\_LEFT</math> and <math>lda &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, m)</math> or</li> <li>▶ <math>alpha == NULL</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strsm](#), [dtrsm](#), [ctrsm](#), [ztrsm](#)

## 2.7.11. `cusblas<t>trsmBatched()`

```

cusblasStatus_t cusblasStrsmBatched( cusblasHandle_t  handle,
                                     cublasSideMode_t  side,
                                     cublasFillMode_t  uplo,
                                     cublasOperation_t  trans,
                                     cublasDiagType_t  diag,
                                     int m,
                                     int n,
                                     const float *alpha,
                                     const float *const A[],

```

```

        int lda,
        float *const B[],
        int ldb,
        int batchCount);
cublasStatus_t cublasDtrsmBatched( cublasHandle_t handle,
        cublasSideMode_t side,
        cublasFillMode_t uplo,
        cublasOperation_t trans,
        cublasDiagType_t diag,
        int m,
        int n,
        const double *alpha,
        const double *const A[],
        int lda,
        double *const B[],
        int ldb,
        int batchCount);
cublasStatus_t cublasCtrsmBatched( cublasHandle_t handle,
        cublasSideMode_t side,
        cublasFillMode_t uplo,
        cublasOperation_t trans,
        cublasDiagType_t diag,
        int m,
        int n,
        const cuComplex *alpha,
        const cuComplex *const A[],
        int lda,
        cuComplex *const B[],
        int ldb,
        int batchCount);
cublasStatus_t cublasZtrsmBatched( cublasHandle_t handle,
        cublasSideMode_t side,
        cublasFillMode_t uplo,
        cublasOperation_t trans,
        cublasDiagType_t diag,
        int m,
        int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *const A[],
        int lda,
        cuDoubleComplex *const B[],
        int ldb,
        int batchCount);

```

This function solves an array of triangular linear systems with multiple right-hand-sides

$$\begin{cases} \text{op}(A[i])X[i] = \alpha B[i] & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ X[i]\text{op}(A[i]) = \alpha B[i] & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A[i]$  is a triangular matrix stored in lower or upper mode with or without the main diagonal,  $X[i]$  and  $B[i]$  are  $m \times n$  matrices, and  $\alpha$  is a scalar. Also, for matrix  $A$

$$\text{op}(A[i]) = \begin{cases} A[i] & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T[i] & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H[i] & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $X[i]$  overwrites the right-hand-sides  $B[i]$  on exit.

No test for singularity or near-singularity is included in this function.

This function works for any sizes but is intended to be used for matrices of small sizes where the launch overhead is a significant factor. For bigger sizes, it might be advantageous to call `batchCount` times the regular `cublas<t>trsm` within a set of CUDA streams.

The current implementation is limited to devices with compute capability above or equal to 2.0.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
side		input	Indicates if matrix $A[i]$ is on the left or right of $x[i]$ .
uplo		input	Indicates if matrix $A[i]$ lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation $op(A[i])$ that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix $A[i]$ are unity and should not be accessed.
m		input	Number of rows of matrix $B[i]$ , with matrix $A[i]$ sized accordingly.
n		input	Number of columns of matrix $B[i]$ , with matrix $A[i]$ is sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication, if $\alpha==0$ then $A[i]$ is not referenced and $B[i]$ does not have to be a valid input.
A	device	input	Array of pointers to <type> array, with each array of dim. $lda \times m$ with $lda \geq \max(1, m)$ if $side == CUBLAS\_SIDE\_LEFT$ and $lda \times n$ with $lda \geq \max(1, n)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix $A[i]$ .
B	device	in/out	Array of pointers to <type> array, with each array of dim. $ldb \times n$ with $ldb \geq \max(1, m)$ . Matrices $B[i]$ should not overlap; otherwise, undefined behavior is expected.
ldb		input	Leading dimension of two-dimensional array used to store matrix $B[i]$ .
batchCount		input	Number of pointers contained in A and B.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>batchCount &lt; 0</math> or</li> <li>▶ if <math>trans \neq CUBLAS\_OP\_N, CUBLAS\_OP\_C, CUBLAS\_OP\_T</math> or</li> <li>▶ if <math>uplo \neq CUBLAS\_FILL\_MODE\_LOWER, CUBLAS\_FILL\_MODE\_UPPER</math> or</li> </ul>

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if side != CUBLAS_SIDE_LEFT, CUBLAS_SIDE_RIGHT or</li> <li>▶ if diag != CUBLAS_DIAG_NON_UNIT, CUBLAS_DIAG_UNIT or</li> <li>▶ if lda &lt; max{1, m} if side == CUBLAS_SIDE_LEFT and lda &lt; max{1, n} otherwise or</li> <li>▶ ldb &lt; max{1, m}</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strsm](#), [dtrsm](#), [ctrsm](#), [ztrsm](#)

## 2.7.12. cublas<t>hemm()

```

cublasStatus_t cublasChemmm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *B, int ldb,
                           const cuComplex      *beta,
                           cuComplex            *C, int ldc)
cublasStatus_t cublasZhemmm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex      *C, int ldc)

```

This function performs the Hermitian matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side == CUBLAS_SIDE_LEFT} \\ \alpha BA + \beta C & \text{if side == CUBLAS_SIDE_RIGHT} \end{cases}$$

where  $A$  is a Hermitian matrix stored in lower or upper mode,  $B$  and  $C$  are  $m \times n$  matrices, and  $\alpha$  and  $\beta$  are scalars.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
m		input	Number of rows of matrix C and B, with matrix A sized accordingly.

Param.	Memory	In/out	Meaning
n		input	Number of columns of matrix C and B, with matrix A sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times m$ with $lda \geq \max(1, m)$ if <code>side == CUBLAS_SIDE_LEFT</code> and $lda \times n$ with $lda \geq \max(1, n)$ otherwise. The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta		input	<type> scalar used for multiplication, if <code>beta == 0</code> then C does not have to be a valid input.
C	device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or</li> <li>▶ if <code>side != CUBLAS_SIDE_LEFT, CUBLAS_SIDE_RIGHT</code> or</li> <li>▶ if <code>uplo != CUBLAS_FILL_MODE_LOWER, CUBLAS_FILL_MODE_UPPER</code> or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if <code>side == CUBLAS_SIDE_LEFT</code> and <math>lda &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, m)</math> or</li> <li>▶ if <math>ldc &lt; \max(1, m)</math> or</li> <li>▶ if <code>alpha == NULL</code> or <code>beta == NULL</code> or</li> <li>▶ <code>c == NULL</code></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[chemm](#), [zhemm](#)

## 2.7.13. cublas<t>herk()

```
cublasStatus_t cublasCherk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float *alpha,
                           const cuComplex *A, int lda,
                           const float *beta,
                           cuComplex *C, int ldc)
cublasStatus_t cublasZherk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const double *alpha,
                           const cuDoubleComplex *A, int lda,
                           const double *beta,
                           cuDoubleComplex *C, int ldc)
```

This function performs the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \ n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A) and c.
k		input	Number of columns of matrix op(A).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension lda x k with lda>=max(1,n) if transa == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta		input	<type> scalar used for multiplication, if beta==0 then c does not have to be a valid input.

Param.	Memory	In/out	Meaning
c	device	in/out	<type> array of dimension $ldc \times n$ , with $ldc \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix c.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> if <math>trans == \text{CUBLAS\_OP\_N}</math> and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, n)</math> or</li> <li>▶ if <math>alpha == \text{NULL}</math> or <math>beta == \text{NULL}</math> or</li> <li>▶ <math>c == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#), [zherk](#)

## 2.7.14. cublas<t>her2k()

```

cublasStatus_t cublasCher2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex *alpha,
                           const cuComplex *A, int lda,
                           const cuComplex *B, int ldb,
                           const float *beta,
                           cuComplex *C, int ldc)
cublasStatus_t cublasZher2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const double *beta,
                           cuDoubleComplex *C, int ldc)

```



This function performs the Hermitian rank-  $2k$  update

$$C = \alpha \text{op}(A)\text{op}(B)^H + \bar{\alpha} \text{op}(B)\text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^H \text{ and } B^H & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix a lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension lda x k with lda>=max(1,n) if transa == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension ldb x k with ldb>=max(1,n) if transb == CUBLAS_OP_N and ldb x n with ldb>=max(1,k) otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if beta==0 then c does not have to be a valid input.
C	device	in/out	<type> array of dimension ldc x n, with ldc>=max(1,n). The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Error Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>\text{trans} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{uplo} \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, n)</math> or</li> <li>▶ if <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math> or</li> <li>▶ <math>c == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cher2k](#), [zher2k](#)

## 2.7.15. cublas<t>herkx()

```

cublasStatus_t cublasCherkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex *alpha,
                           const cuComplex *A, int lda,
                           const cuComplex *B, int ldb,
                           const float *beta,
                           cuComplex *C, int ldc)
cublasStatus_t cublasZherkx(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const double *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs a variation of the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(B)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^H \text{ and } B^H & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

This routine can be used when the matrix B is in such way that the result is guaranteed to be Hermitian. An usual example is when the matrix B is a scaled form of the matrix A : this is equivalent to B being the product of the matrix A and a diagonal matrix. For an efficient computation of the product of a regular matrix with a diagonal matrix, refer to the routine `cublas<t>dggmm`.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times k$ with $\text{lda} \geq \max(1, n)$ if <code>transa == CUBLAS_OP_N</code> and $\text{lda} \times n$ with $\text{lda} \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ if <code>transb == CUBLAS_OP_N</code> and $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	Real scalar used for multiplication, if <code>beta==0</code> then c does not have to be a valid input.
C	device	in/out	<type> array of dimension $\text{ldc} \times n$ , with $\text{ldc} \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

Error Value	Meaning
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>\text{trans} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{uplo} \neq \text{CUBLAS\_FILL\_MODE\_LOWER}, \text{CUBLAS\_FILL\_MODE\_UPPER}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, n)</math> if <math>\text{trans} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, n)</math> or</li> <li>▶ if <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math> or</li> <li>▶ <math>\text{c} == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#), [zherk](#) and

[cher2k](#), [zher2k](#)

## 2.8. BLAS-like Extension

This section describes the BLAS-extension functions that perform matrix-matrix operations.

### 2.8.1. `cublas<t>geam()`

```

cublasStatus_t cublasSgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const float *alpha,
                           const float *A, int lda,
                           const float *beta,
                           const float *B, int ldb,
                           float *C, int ldc)
cublasStatus_t cublasDgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const double *alpha,
                           const double *A, int lda,
                           const double *beta,
                           const double *B, int ldb,
                           double *C, int ldc)
cublasStatus_t cublasCgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const cuComplex *alpha,
                           const cuComplex *A, int lda,
                           const cuComplex *beta,
                           const cuComplex *B, int ldb,
                           cuComplex *C, int ldc)

```

```

        const cuComplex      *B, int ldb,
        cuComplex      *C, int ldc)
cublasStatus_t cublasZgeam(cublasHandle_t handle,
        cublasOperation_t transa, cublasOperation_t transb,
        int m, int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *beta,
        const cuDoubleComplex *B, int ldb,
        cuDoubleComplex *C, int ldc)

```

This function performs the matrix-matrix addition/transposition

$$C = \alpha \text{op}(A) + \beta \text{op}(B)$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times n$ ,  $\text{op}(B) m \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .

The operation is out-of-place if  $C$  does not overlap  $A$  or  $B$ .

The in-place mode supports the following two operations,

$$C = \alpha * C + \beta \text{op}(B)$$

$$C = \alpha \text{op}(A) + \beta * C$$

For in-place mode, if  $C = A$ ,  $\text{ldc} = \text{lda}$  and  $\text{transa} = \text{CUBLAS\_OP\_N}$ . If  $C = B$ ,  $\text{ldc} = \text{ldb}$  and  $\text{transb} = \text{CUBLAS\_OP\_N}$ . If the user does not meet above requirements, `CUBLAS_STATUS_INVALID_VALUE` is returned.

The operation includes the following special cases:

- the user can reset matrix  $C$  to zero by setting  $*\alpha = *\beta = 0$
- the user can transpose matrix  $A$  by setting  $*\alpha = 1$  and  $*\beta = 0$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}(B)$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}(A)$ and $C$ .
n		input	Number of columns of matrix $\text{op}(B)$ and $C$ .
alpha	host or device	input	<type> scalar used for multiplication. If $*\alpha == 0$ , $A$ does not have to be a valid input.
A	device	input	<type> array of dimensions $\text{lda} \times n$ with $\text{lda} \geq \max(1, m)$ if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $\text{lda} \times m$ with $\text{lda} \geq \max(1, n)$ otherwise.

Param.	Memory	In/out	Meaning
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension ldb x n with ldb>=max(1,m) if transb == CUBLAS_OP_N and ldb x m with ldb>=max(1,n) otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication. If *beta == 0, B does not have to be a valid input.
C	device	output	<type> array of dimensions ldc x n with ldc>=max(1,m).
ldc		input	Leading dimension of a two-dimensional array used to store the matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or</li> <li>▶ if <math>\text{transa} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{transb} \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>\text{lda} &lt; \max(1, m)</math> if <math>\text{transa} == \text{CUBLAS\_OP\_N}</math> and <math>\text{lda} &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>\text{ldb} &lt; \max(1, m)</math> if <math>\text{transb} == \text{CUBLAS\_OP\_N}</math> and <math>\text{ldb} &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>\text{ldc} &lt; \max(1, m)</math> or</li> <li>▶ if <math>A == C, ((\text{CUBLAS\_OP\_N} \neq \text{transa}) \parallel (\text{lda} \neq \text{ldc}))</math> or</li> <li>▶ if <math>B == C, ((\text{CUBLAS\_OP\_N} \neq \text{transb}) \parallel (\text{ldb} \neq \text{ldc}))</math> or</li> <li>▶ <math>\alpha == \text{NULL}</math> or <math>\beta == \text{NULL}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 2.8.2. cublas<t>dgmm()

```
cublasStatust cublasSdgmm(cublasHandle_t handle, cublasSideMode_t mode,
```

```

        int m, int n,
        const float *A, int lda,
        const float *x, int incx,
        float *C, int ldc)
cublasStatus_t cublasDdggmm(cublasHandle_t handle, cublasSideMode_t mode,
        int m, int n,
        const double *A, int lda,
        const double *x, int incx,
        double *C, int ldc)
cublasStatus_t cublasCdggmm(cublasHandle_t handle, cublasSideMode_t mode,
        int m, int n,
        const cuComplex *A, int lda,
        const cuComplex *x, int incx,
        cuComplex *C, int ldc)
cublasStatus_t cublasZdggmm(cublasHandle_t handle, cublasSideMode_t mode,
        int m, int n,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *x, int incx,
        cuDoubleComplex *C, int ldc)

```

This function performs the matrix-matrix multiplication

$$C = \begin{cases} A \times \text{diag}(X) & \text{if mode} == \text{CUBLAS\_SIDE\_RIGHT} \\ \text{diag}(X) \times A & \text{if mode} == \text{CUBLAS\_SIDE\_LEFT} \end{cases}$$

where  $A$  and  $C$  are matrices stored in column-major format with dimensions  $m \times n$ .  $X$  is a vector of size  $n$  if mode == CUBLAS\_SIDE\_RIGHT and of size  $m$  if mode == CUBLAS\_SIDE\_LEFT.  $X$  is gathered from one-dimensional array  $x$  with stride  $\text{incx}$ . The absolute value of  $\text{incx}$  is the stride and the sign of  $\text{incx}$  is direction of the stride. If  $\text{incx}$  is positive, then we forward  $x$  from the first element. Otherwise, we backward  $x$  from the last element. The formula of  $X$  is

$$X[j] = \begin{cases} x[j \times \text{incx}] & \text{if } \text{incx} \geq 0 \\ x[(\chi - 1) \times |\text{incx}| - j \times |\text{incx}|] & \text{if } \text{incx} < 0 \end{cases}$$

where  $\chi = m$  if mode == CUBLAS\_SIDE\_LEFT and  $\chi = n$  if mode == CUBLAS\_SIDE\_RIGHT.

Example 1: if the user wants to perform  $\text{diag}(\text{diag}(B)) \times A$ , then  $\text{incx} = \text{ldb} + 1$  where  $\text{ldb}$  is leading dimension of matrix  $B$ , either row-major or column-major.

Example 2: if the user wants to perform  $\alpha \times A$ , then there are two choices, either cublasgeam with  $\text{*beta}=0$  and  $\text{transa} == \text{CUBLAS\_OP\_N}$  or cublasdggmm with  $\text{incx}=0$  and  $x[0]=\alpha$ .

The operation is out-of-place. The in-place only works if  $\text{lda} = \text{ldc}$ .

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
mode		input	Left multiply if mode == CUBLAS_SIDE_LEFT or right multiply if mode == CUBLAS_SIDE_RIGHT
m		input	Number of rows of matrix A and C.
n		input	Number of columns of matrix A and C.
A	device	input	<type> array of dimensions $\text{lda} \times n$ with $\text{lda} \geq \max(1, m)$
lda		input	Leading dimension of two-dimensional array used to store the matrix A.

Param.	Memory	In/out	Meaning
x	device	input	One-dimensional <type> array of size $l_{inc} \times m$ if <code>mode == CUBLAS_SIDE_LEFT</code> and $l_{inc} \times n$ if <code>mode == CUBLAS_SIDE_RIGHT</code>
incx		input	Stride of one-dimensional array x.
C	device	in/out	<type> array of dimensions <code>ldc x n</code> with <code>ldc</code> $\geq \max(1, m)$ .
ldc		input	Leading dimension of a two-dimensional array used to store the matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or</li> <li>▶ if <code>mode != CUBLAS_SIDE_LEFT</code>, <code>CUBLAS_SIDE_RIGHT</code> or</li> <li>▶ if <code>lda &lt; max(1, m)</code> or</li> <li>▶ <code>ldc &lt; max(1, m)</code></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

### 2.8.3. `cublas<t>getrfBatched()`

```

cublasStatus_t cublasSgetrfBatched(cublasHandle_t handle,
                                   int n,
                                   float *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasDgetrfBatched(cublasHandle_t handle,
                                   int n,
                                   double *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasCgetrfBatched(cublasHandle_t handle,
                                   int n,
                                   cuComplex *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasZgetrfBatched(cublasHandle_t handle,
                                   int n,

```



```
cuDoubleComplex *const Aarray[],
int lda,
int *PivotArray,
int *infoArray,
int batchSize);
```

Aarray is an array of pointers to matrices stored in column-major format with dimensions  $n \times n$  and leading dimension  $lda$ .

This function performs the LU factorization of each Aarray[i] for  $i = 0, \dots, batchSize-1$  by the following equation

$$P * Aarray[i] = L * U$$

where  $P$  is a permutation matrix which represents partial pivoting with row interchanges.  $L$  is a lower triangular matrix with unit diagonal and  $U$  is an upper triangular matrix.

Formally  $P$  is written by a product of permutation matrices  $P_j$ , for  $j = 1, 2, \dots, n$ , say  $P = P_1 * P_2 * P_3 * \dots * P_n$ .  $P_j$  is a permutation matrix which interchanges two rows of vector  $x$  when performing  $P_j * x$ .  $P_j$  can be constructed by  $j$  element of PivotArray[i] by the following Matlab code:

```
// In Matlab PivotArray[i] is an array of base-1.
// In C, PivotArray[i] is base-0.
Pj = eye(n);
swap Pj(j,:) and Pj(PivotArray[i][j] ,:)
```

$L$  and  $U$  are written back to original matrix  $A$ , and diagonal elements of  $L$  are discarded. The  $L$  and  $U$  can be constructed by the following Matlab code:

```
// A is a matrix of nxn after getrf.
L = eye(n);
for j = 1:n
    L(j+1:n,j) = A(j+1:n,j)
end
U = zeros(n);
for i = 1:n
    U(i,i:n) = A(i,i:n)
end
```

If matrix  $A(=Aarray[i])$  is singular, getrf still works and the value of info(=infoArray[i]) reports first row index that LU factorization cannot proceed. If info is  $k$ ,  $U(k,k)$  is zero. The equation  $P * A = L * U$  still holds, however  $L$  and  $U$  reconstruction needs different Matlab code as follows:

```
// A is a matrix of nxn after getrf.
// info is k, which means U(k,k) is zero.
L = eye(n);
for j = 1:k-1
    L(j+1:n,j) = A(j+1:n,j)
end
U = zeros(n);
for i = 1:k-1
    U(i,i:n) = A(i,i:n)
end
for i = k:n
    U(i,k:n) = A(i,k:n)
```

```
end
```

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor.

`cublas<t>getrfBatched` supports non-pivot LU factorization if `PivotArray` is `nil`.

`cublas<t>getrfBatched` supports arbitrary dimension.

`cublas<t>getrfBatched` only supports compute capability 2.0 or above.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>n</code>		input	Number of rows and columns of <code>Aarray[i]</code> .
<code>Aarray</code>	device	input/ output	Array of pointers to <type> array, with each array of dim. $n \times n$ with $lda \geq \max(1, n)$ . Matrices <code>Aarray[i]</code> should not overlap; otherwise, undefined behavior is expected.
<code>lda</code>		input	Leading dimension of two-dimensional array used to store each matrix <code>Aarray[i]</code> .
<code>PivotArray</code>	device	output	Array of size $n \times batchSize$ that contains the pivoting sequence of each factorization of <code>Aarray[i]</code> stored in a linear fashion. If <code>PivotArray</code> is <code>nil</code> , pivoting is disabled.
<code>infoArray</code>	device	output	Array of size <code>batchSize</code> that <code>info[=infoArray[i]]</code> contains the information of factorization of <code>Aarray[i]</code> .  If <code>info=0</code> , the execution is successful.  If <code>info = -j</code> , the <code>j</code> -th parameter had an illegal value.  If <code>info = k</code> , <code>U(k,k)</code> is 0. The factorization has been completed, but <code>U</code> is exactly singular.
<code>batchSize</code>		input	Number of pointers contained in <code>A</code> .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	The parameters <code>n</code> , <code>batchSize</code> , <code>lda</code> < 0.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[sgeqrf](#), [dgeqrf](#), [cgeqrf](#), [zgeqrf](#)

## 2.8.4. `cublas<t>getrsBatched()`

```

cublasStatus_t cublasSgetrsBatched(cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int n,
                                   int nrhs,
                                   const float *const Aarray[],
                                   int lda,
                                   const int *devIpiv,
                                   float *const Barray[],
                                   int ldb,
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasDgetrsBatched(cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int n,
                                   int nrhs,
                                   const double *const Aarray[],
                                   int lda,
                                   const int *devIpiv,
                                   double *const Barray[],
                                   int ldb,
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasCgetrsBatched(cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int n,
                                   int nrhs,
                                   const cuComplex *const Aarray[],
                                   int lda,
                                   const int *devIpiv,
                                   cuComplex *const Barray[],
                                   int ldb,
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasZgetrsBatched(cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int n,
                                   int nrhs,
                                   const cuDoubleComplex *const Aarray[],
                                   int lda,
                                   const int *devIpiv,
                                   cuDoubleComplex *const Barray[],
                                   int ldb,
                                   int *info,
                                   int batchSize);

```

This function solves an array of systems of linear equations of the form:

$$\text{op}(A[i])X[i] = B[i]$$

where  $A[i]$  is a matrix which has been LU factorized with pivoting,  $X[i]$  and  $B[i]$  are  $n \times \text{nrhs}$  matrices. Also, for matrix  $A$

$$\text{op}(A[i]) = \begin{cases} A[i] & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^T[i] & \text{if trans} == \text{CUBLAS\_OP\_T} \\ A^H[i] & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor.

`cublas<t>getrsBatched` supports non-pivot LU factorization if `devIpiv` is nil.

`cublas<t>getrsBatched` supports arbitrary dimension.

`cublas<t>getrsBatched` only supports compute capability 2.0 or above.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>trans</code>		input	Operation <code>op[A]</code> that is non- or [conj.] transpose.
<code>n</code>		input	Number of rows and columns of <code>Aarray[i]</code> .
<code>nrhs</code>		input	Number of columns of <code>Barray[i]</code> .
<code>Aarray</code>	device	input	Array of pointers to <type> array, with each array of dim. $n \times n$ with $lda \geq \max(1, n)$ .
<code>lda</code>		input	Leading dimension of two-dimensional array used to store each matrix <code>Aarray[i]</code> .
<code>devIpiv</code>	device	input	Array of size $n \times batchSize$ that contains the pivoting sequence of each factorization of <code>Aarray[i]</code> stored in a linear fashion. If <code>devIpiv</code> is nil, pivoting for all <code>Aarray[i]</code> is ignored.
<code>Barray</code>	device	input/output	Array of pointers to <type> array, with each array of dim. $n \times nrhs$ with $ldb \geq \max(1, n)$ . Matrices <code>Barray[i]</code> should not overlap; otherwise, undefined behavior is expected.
<code>ldb</code>		input	Leading dimension of two-dimensional array used to store each solution matrix <code>Barray[i]</code> .
<code>info</code>	host	output	If <code>info=0</code> , the execution is successful. If <code>info = -j</code> , the <code>j</code> -th parameter had an illegal value.
<code>batchSize</code>		input	Number of pointers contained in <code>A</code> .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<ul style="list-style-type: none"> <li>► If <math>n &lt; 0</math> or <math>nrhs &lt; 0</math> or</li> <li>► if <code>trans</code> != <code>CUBLAS_OP_N</code>, <code>CUBLAS_OP_C</code>, <code>CUBLAS_OP_T</code> or</li> <li>► if <math>lda &lt; \max(1, n)</math> or</li> <li>► <math>ldb &lt; \max(1, n)</math></li> </ul>
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

For references please refer to:

[sgeqrs](#), [dgeqrs](#), [cgeqrs](#), [zgeqrs](#)

## 2.8.5. cublas<t>getriBatched()

```

cublasStatus_t cublasSgetriBatched(cublasHandle_t handle,
                                   int n,
                                   const float *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   float *const Carray[],
                                   int ldc,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasDgetriBatched(cublasHandle_t handle,
                                   int n,
                                   const double *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   double *const Carray[],
                                   int ldc,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasCgetriBatched(cublasHandle_t handle,
                                   int n,
                                   const cuComplex *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   cuComplex *const Carray[],
                                   int ldc,
                                   int *infoArray,
                                   int batchSize);

cublasStatus_t cublasZgetriBatched(cublasHandle_t handle,
                                   int n,
                                   const cuDoubleComplex *const Aarray[],
                                   int lda,
                                   int *PivotArray,
                                   cuDoubleComplex *const Carray[],
                                   int ldc,
                                   int *infoArray,
                                   int batchSize);

```

Aarray and Carray are arrays of pointers to matrices stored in column-major format with dimensions  $n \times n$  and leading dimension `lda` and `ldc` respectively.

This function performs the inversion of matrices  $A[i]$  for  $i = 0, \dots, \text{batchSize}-1$ .

Prior to calling `cublas<t>getriBatched`, the matrix  $A[i]$  must be factorized first using the routine `cublas<t>getrfBatched`. After the call of `cublas<t>getrfBatched`, the matrix pointing by `Aarray[i]` will contain the LU factors of the matrix  $A[i]$  and the vector pointing by `(PivotArray+i)` will contain the pivoting sequence.

Following the LU factorization, `cublas<t>getriBatched` uses forward and backward triangular solvers to complete inversion of matrices  $A[i]$  for  $i = 0, \dots, \text{batchSize}-1$ . The inversion is out-of-place, so memory space of `Carray[i]` cannot overlap memory space of `Array[i]`.

Typically, all parameters in `cublas<type>getrfBatched` would be passed into `cublas<type>getriBatched`. For example,

```
// step 1: perform in-place LU decomposition, P*A = L*U.
//   Aarray[i] is n*n matrix A[i]
//   cublasDgetrfBatched(handle, n, Aarray, lda, PivotArray, infoArray, batchSize);
//   check infoArray[i] to see if factorization of A[i] is successful or not.
//   Array[i] contains LU factorization of A[i]

// step 2: perform out-of-place inversion, Carray[i] = inv(A[i])
//   cublasDgetriBatched(handle, n, Aarray, lda, PivotArray, Carray, ldc, infoArray,
//   batchSize);
//   check infoArray[i] to see if inversion of A[i] is successful or not.
```

The user can check singularity from either `cublas<type>getrfBatched` or `cublas<type>getriBatched`.

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor.

If `cublas<type>getrfBatched` is performed by non-pivoting, `PivotArray` of `cublas<type>getriBatched` should be nil.

`cublas<type>getriBatched` supports arbitrary dimension.

`cublas<type>getriBatched` only supports compute capability 2.0 or above.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of rows and columns of <code>Aarray[i]</code> .
Aarray	device	input	Array of pointers to <type> array, with each array of dimension $n \times n$ with $lda \geq \max(1, n)$ .
lda		input	Leading dimension of two-dimensional array used to store each matrix <code>Aarray[i]</code> .
PivotArray	device	output	Array of size $n \times batchSize$ that contains the pivoting sequence of each factorization of <code>Aarray[i]</code> stored in a linear fashion. If <code>PivotArray</code> is nil, pivoting is disabled.
Carray	device	output	Array of pointers to <type> array, with each array of dimension $n \times n$ with $ldc \geq \max(1, n)$ . Matrices <code>Carray[i]</code> should not overlap; otherwise, undefined behavior is expected.
ldc		input	Leading dimension of two-dimensional array used to store each matrix <code>Carray[i]</code> .
infoArray	device	output	<p>Array of size <code>batchSize</code> that <code>info[=infoArray[i]]</code> contains the information of inversion of <code>A[i]</code>.</p> <p>If <code>info=0</code>, the execution is successful.</p> <p>If <code>info = k</code>, <code>U(k,k)</code> is 0. The <code>U</code> is exactly singular and the inversion failed.</p>

Param.	Memory	In/out	Meaning
batchSize		input	number of pointers contained in A.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>lda &lt; 0</math> or <math>ldc &lt; 0</math> or <math>batchSize &lt; 0</math> or</li> <li>▶ <math>lda &lt; n</math> or <math>ldc &lt; n</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 2.8.6. `cublas<t>matinvBatched()`

```

cublasStatus_t cublasSmatinvBatched(cublasHandle_t handle,
    int n,
    const float *const A[],
    int lda,
    float *const Ainv[],
    int lda_inv,
    int *info,
    int batchSize);

cublasStatus_t cublasDmatinvBatched(cublasHandle_t handle,
    int n,
    const double *const A[],
    int lda,
    double *const Ainv[],
    int lda_inv,
    int *info,
    int batchSize);

cublasStatus_t cublasCmatinvBatched(cublasHandle_t handle,
    int n,
    const cuComplex *const A[],
    int lda,
    cuComplex *const Ainv[],
    int lda_inv,
    int *info,
    int batchSize);

cublasStatus_t cublasZmatinvBatched(cublasHandle_t handle,
    int n,
    const cuDoubleComplex *const A[],
    int lda,
    cuDoubleComplex *const Ainv[],
    int lda_inv,
    int *info,
    int batchSize);

```

`A` and `Ainv` are arrays of pointers to matrices stored in column-major format with dimensions  $n \times n$  and leading dimension `lda` and `lda_inv` respectively.

This function performs the inversion of matrices `A[i]` for  $i = 0, \dots, \text{batchSize}-1$ .

This function is a short cut of `cublas<type>getrfBatched` plus `cublas<type>getriBatched`. However it doesn't work if  $n$  is greater than 32. If not, the user has to go through `cublas<type>getrfBatched` and `cublas<type>getriBatched`.

If the matrix `A[i]` is singular, then `info[i]` reports singularity, the same as `cublas<type>getrfBatched`.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>n</code>		input	Number of rows and columns of <code>A[i]</code> .
<code>A</code>	device	input	Array of pointers to <type> array, with each array of dimension $n \times n$ with <code>lda</code> $\geq \max(1, n)$ .
<code>lda</code>		input	Leading dimension of two-dimensional array used to store each matrix <code>A[i]</code> .
<code>Ainv</code>	device	output	Array of pointers to <type> array, with each array of dimension $n \times n$ with <code>lda_inv</code> $\geq \max(1, n)$ . Matrices <code>Ainv[i]</code> should not overlap; otherwise, undefined behavior is expected.
<code>lda_inv</code>		input	Leading dimension of two-dimensional array used to store each matrix <code>Ainv[i]</code> .
<code>info</code>	device	output	Array of size <code>batchSize</code> that <code>info[i]</code> contains the information of inversion of <code>A[i]</code> .  If <code>info[i]=0</code> , the execution is successful.  If <code>info[i]=k</code> , <code>U(k,k)</code> is 0. The <code>U</code> is exactly singular and the inversion failed.
<code>batchSize</code>		input	Number of pointers contained in <code>A</code> .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>lda &lt; 0</math> or <math>lda\_inv &lt; 0</math> or <math>batchSize &lt; 0</math> or</li> <li>▶ if <math>lda &lt; n</math> or <math>lda\_inv &lt; n</math> or</li> <li>▶ <math>n &gt; 32</math></li> </ul>
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.



## 2.8.7. `cublas<t>geqrfBatched()`

```

cublasStatus_t cublasSgeqrfBatched( cublasHandle_t handle,
                                   int m,
                                   int n,
                                   float *const Aarray[],
                                   int lda,
                                   float *const TauArray[],
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasDgeqrfBatched( cublasHandle_t handle,
                                   int m,
                                   int n,
                                   double *const Aarray[],
                                   int lda,
                                   double *const TauArray[],
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasCgeqrfBatched( cublasHandle_t handle,
                                   int m,
                                   int n,
                                   cuComplex *const Aarray[],
                                   int lda,
                                   cuComplex *const TauArray[],
                                   int *info,
                                   int batchSize);

cublasStatus_t cublasZgeqrfBatched( cublasHandle_t handle,
                                   int m,
                                   int n,
                                   cuDoubleComplex *const Aarray[],
                                   int lda,
                                   cuDoubleComplex *const TauArray[],
                                   int *info,
                                   int batchSize);

```

Aarray is an array of pointers to matrices stored in column-major format with dimensions  $m \times n$  and leading dimension  $lda$ . TauArray is an array of pointers to vectors of dimension of at least  $\max(1, \min(m, n))$ .

This function performs the QR factorization of each `Aarray[i]` for  $i = 0, \dots, \text{batchSize}-1$  using Householder reflections. Each matrix  $Q[i]$  is represented as a product of elementary reflectors and is stored in the lower part of each `Aarray[i]` as follows:

```


$$Q[j] = H[j][1] H[j][2] \dots H[j](k), \text{ where } k = \min(m, n).$$

Each  $H[j][i]$  has the form

$$H[j][i] = I - \tau[j] * v * v'$$

where  $\tau[j]$  is a real scalar, and  $v$  is a real vector with
 $v(1:i-1) = 0$  and  $v(i) = 1$ ;  $v(i+1:m)$  is stored on exit in Aarray[j][i+1:m,i],
and  $\tau$  in TauArray[j][i]

```

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor.

`cublas<t>geqrfBatched` supports arbitrary dimension.

`cublas<t>geqrfBatched` only supports compute capability 2.0 or above.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
m		input	Number of rows <code>Aarray[i]</code> .
n		input	Number of columns of <code>Aarray[i]</code> .
Aarray	device	input	Array of pointers to <type> array, with each array of dim. $m \times n$ with $lda \geq \max(1, m)$ .
lda		input	Leading dimension of two-dimensional array used to store each matrix <code>Aarray[i]</code> .
TauArray	device	output	Array of pointers to <type> vector, with each vector of dim. $\max(1, \min(m, n))$ .
info	host	output	If <code>info=0</code> , the parameters passed to the function are valid. If <code>info&lt;0</code> , the parameter in position <code>-info</code> is invalid.
batchSize		input	Number of pointers contained in A.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>batchSize &lt; 0</math> or</li> <li>▶ <math>lda &lt; \max(1, m)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgeqrf](#), [dgeqrf](#), [cgeqrf](#), [zgeqrf](#)

## 2.8.8. `cublas<t>gelsBatched()`

```
cublasStatus_t cublasSgelsBatched( cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int m,
                                   int n,
                                   int nrhs,
                                   float *const Aarray[],
                                   int lda,
                                   float *const Carray[],
                                   int ldc,
                                   int *info,
                                   int *devInfoArray,
```

```

        int batchSize );

cublasStatus_t cublasDgelsBatched( cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int m,
                                   int n,
                                   int nrhs,
                                   double *const Aarray[],
                                   int lda,
                                   double *const Carray[],
                                   int ldc,
                                   int *info,
                                   int *devInfoArray,
                                   int batchSize );

cublasStatus_t cublasCgelsBatched( cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int m,
                                   int n,
                                   int nrhs,
                                   cuComplex *const Aarray[],
                                   int lda,
                                   cuComplex *const Carray[],
                                   int ldc,
                                   int *info,
                                   int *devInfoArray,
                                   int batchSize );

cublasStatus_t cublasZgelsBatched( cublasHandle_t handle,
                                   cublasOperation_t trans,
                                   int m,
                                   int n,
                                   int nrhs,
                                   cuDoubleComplex *const Aarray[],
                                   int lda,
                                   cuDoubleComplex *const Carray[],
                                   int ldc,
                                   int *info,
                                   int *devInfoArray,
                                   int batchSize );

```

Aarray is an array of pointers to matrices stored in column-major format. Carray is an array of pointers to matrices stored in column-major format.

This function find the least squares solution of a batch of overdetermined systems: it solves the least squares problem described as follows :

```

        minimize || Carray[i] - Aarray[i]*Xarray[i] || , with i =
0, ...,batchSize-1

```

On exit, each Aarray[i] is overwritten with their QR factorization and each Carray[i] is overwritten with the least square solution.

cublas<t>gelsBatched supports only the non-transpose operation and only solves overdetermined systems ( $m \geq n$ ).

cublas<t>gelsBatched only supports compute capability 2.0 or above.

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
trans		input	Operation op(Aarray[i]) that is non- or (conj.) transpose. Only non-transpose operation is currently supported.
m		input	Number of rows of each Aarray[i] and Carray[i] if trans == CUBLAS_OP_N, numbers of columns of each Aarray[i] otherwise (not supported currently).
n		input	Number of columns of each Aarray[i] if trans == CUBLAS_OP_N, and number of rows of each Aarray[i] and Carray[i] otherwise (not supported currently).
nrhs		input	Number of columns of each Carray[i].
Aarray	device	input/output	Array of pointers to <type> array, with each array of dim. $m \times n$ with $lda \geq \max(1, m)$ if trans == CUBLAS_OP_N, and $n \times m$ with $lda \geq \max(1, n)$ otherwise (not supported currently). Matrices Aarray[i] should not overlap; otherwise, undefined behavior is expected.
lda		input	Leading dimension of two-dimensional array used to store each matrix Aarray[i].
Carray	device	input/output	Array of pointers to <type> array, with each array of dim. $m \times nrhs$ with $ldc \geq \max(1, m)$ if trans == CUBLAS_OP_N, and $n \times nrhs$ with $lda \geq \max(1, n)$ otherwise (not supported currently). Matrices Carray[i] should not overlap; otherwise, undefined behavior is expected.
ldc		input	Leading dimension of two-dimensional array used to store each matrix Carray[i].
info	host	output	If info=0, the parameters passed to the function are valid. If info<0, the parameter in position -info is invalid.
devInfoArray	device	output	Optional array of integers of dimension batchsize. If non-null, every element devInfoArray[i] contain a value V with the following meaning: V = 0 : the i-th problem was sucessfully solved. V > 0 : the V-th diagonal element of the Aarray[i] is zero. Aarray[i] does not have full rank.
batchSize		input	Number of pointers contained in Aarray and Carray.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

Error Value	Meaning
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>► If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>nrhs &lt; 0</math> or <math>batchSize &lt; 0</math> or</li> <li>► <math>lda &lt; \max(1, m)</math> or <math>ldc &lt; \max(1, m)</math></li> </ul>
CUBLAS_STATUS_NOT_SUPPORTED	The parameters $m < n$ or $trans$ is different from non-transpose.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgels](#), [dgels](#), [cgels](#), [zgels](#)

## 2.8.9. `cublas<t>tptr()`

```

cublasStatus_t cublasStpttr ( cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             int n,
                             const float *AP,
                             float *A,
                             int lda );

cublasStatus_t cublasDtptr ( cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             int n,
                             const double *AP,
                             double *A,
                             int lda );

cublasStatus_t cublasCtptr ( cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             int n,
                             const cuComplex *AP,
                             cuComplex *A,
                             int lda );

cublasStatus_t cublasZtptr ( cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             int n,
                             const cuDoubleComplex *AP,
                             cuDoubleComplex *A,
                             int lda );

```

This function performs the conversion from the triangular packed format to the triangular format.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements of `AP` are copied into the lower triangular part of the triangular matrix `A` and the upper part of `A` is left untouched. If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements of `AP` are copied into the upper triangular part of the triangular matrix `A` and the lower part of `A` is left untouched.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.

Param.	Memory	In/out	Meaning
uplo		input	Indicates if matrix $AP$ contains lower or upper part of matrix $A$ .
n		input	Number of rows and columns of matrix $A$ .
$AP$	device	input	<type> array with $A$ stored in packed format.
$A$	device	output	<type> array of dimensions $lda \times n$ , with $lda \geq \max(1, n)$ . The opposite side of $A$ is left untouched.
lda		input	Leading dimension of two-dimensional array used to store matrix $A$ .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER, CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ <math>lda &lt; \max(1, n)</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[stpttr](#), [dtptr](#), [ctpttr](#), [ztptr](#)

## 2.8.10. `cusblas<t>trttp()`

```

cusblasStatus_t cusblasStrttp ( cublasHandle_t handle,
                               cublasFillMode_t uplo,
                               int n,
                               const float *A,
                               int lda,
                               float *AP );

cusblasStatus_t cusblasDtrttp ( cublasHandle_t handle,
                               cublasFillMode_t uplo,
                               int n,
                               const double *A,
                               int lda,
                               double *AP );

cusblasStatus_t cusblasCtrttp ( cublasHandle_t handle,
                               cublasFillMode_t uplo,
                               int n,
                               const cuComplex *A,
                               int lda,
                               cuComplex *AP );

cusblasStatus_t cusblasZtrttp ( cublasHandle_t handle,
```

```

cublasFillMode_t uplo,
int n,
const cuDoubleComplex *A,
int lda,
cuDoubleComplex *AP );

```

This function performs the conversion from the triangular format to the triangular packed format.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the lower triangular part of the triangular matrix **A** is copied into the array **AP**. If `uplo == CUBLAS_FILL_MODE_UPPER` then the upper triangular part of the triangular matrix **A** is copied into the array **AP**.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates which matrix <b>A</b> lower or upper part is referenced.
n		input	Number of rows and columns of matrix <b>A</b> .
A	device	input	<type> array of dimensions <code>lda x n</code> , with <code>lda ≥ max(1, n)</code> .
lda		input	Leading dimension of two-dimensional array used to store matrix <b>A</b> .
AP	device	output	<type> array with <b>A</b> stored in packed format.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <code>n &lt; 0</code> or</li> <li>▶ if <code>uplo != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER</code> or</li> <li>▶ <code>lda &lt; max(1, n)</code></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strttp](#), [dtrttp](#), [ctrttp](#), [ztrttp](#)

## 2.8.11. cublas<t>gemmEx()

```

cublasStatus_t cublasSgemmEx(cublasHandle_t handle,
cublasOperation_t transa,
cublasOperation_t transb,
int m,
int n,
int k,
const float *alpha,

```

```

        const void      *A,
        cudaDataType_t  Atype,
        int             lda,
        const void      *B,
        cudaDataType_t  Btype,
        int             ldb,
        const float      *beta,
        void             *C,
        cudaDataType_t  Ctype,
        int             ldc)
cublasStatus_t cublasCgemmEx(cublasHandle_t handle,
        cublasOperation_t transa,
        cublasOperation_t transb,
        int m,
        int n,
        int k,
        const cuComplex *alpha,
        const void      *A,
        cudaDataType_t  Atype,
        int             lda,
        const void      *B,
        cudaDataType_t  Btype,
        int             ldb,
        const cuComplex *beta,
        void             *C,
        cudaDataType_t  Ctype,
        int             ldc)

```

This function is an extension of `cublas<t>gemm`. In this function the input matrices and output matrices can have a lower precision but the computation is still done in the type `<t>`. For example, in the type `float` for `cublasSgemmEx` and in the type `cuComplex` for `cublasCgemmEx`.

$$C = \alpha \text{op}(A) \text{op}(B) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times k$ ,  $\text{op}(B) k \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}(B)$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}(A)$ and $C$ .
n		input	Number of columns of matrix $\text{op}(B)$ and $C$ .
k		input	Number of columns of $\text{op}(A)$ and rows of $\text{op}(B)$ .
alpha	host or device	input	<code>&lt;type&gt;</code> scalar used for multiplication.



Param.	Memory	In/out	Meaning
A	device	input	<type> array of dimensions $lda \times k$ with $lda \geq \max(1, m)$ if $transa == CUBLAS\_OP\_N$ and $lda \times m$ with $lda \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, k)$ if $transb == CUBLAS\_OP\_N$ and $ldb \times k$ with $ldb \geq \max(1, n)$ otherwise.
Btype		input	Enumerant specifying the datatype of matrix B.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication. If $beta == 0$ , c does not have to be a valid input.
C	device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ .
Ctype		input	Enumerant specifying the datatype of matrix C.
ldc		input	Leading dimension of a two-dimensional array used to store the matrix C.

The matrix types combinations supported for  `cublasSgemmEx`  are listed in the following table:

C	A/B
CUDA_R_16BF	CUDA_R_16BF
CUDA_R_16F	CUDA_R_16F
CUDA_R_32F	CUDA_R_8I
	CUDA_R_16BF
	CUDA_R_16F
	CUDA_R_32F

The matrix types combinations supported for  `cublasCgemmEx`  are listed in the following table :

C	A/B
CUDA_C_32F	CUDA_C_8I
	CUDA_C_32F

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ARCH_MISMATCH	cublasCgemmEx is only supported for GPU with architecture capabilities equal to or greater than 5.0.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters Atype,Btype and Ctype is not supported.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if transa or transb != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if <math>lda &lt; \max\{1, m\}</math> if transa == CUBLAS_OP_N and <math>lda &lt; \max\{1, k\}</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max\{1, k\}</math> if transb == CUBLAS_OP_N and <math>ldb &lt; \max\{1, n\}</math> otherwise or</li> <li>▶ <math>ldc &lt; \max\{1, m\}</math></li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgemm](#)

For more information about the numerical behavior of some GEMM algorithms, refer to the [GEMM Algorithms Numerical Behavior](#) section.

## 2.8.12. cublasGemmEx()

```

cublasStatus_t cublasGemmEx(cublasHandle_t handle,
                           cublasOperation_t transa,
                           cublasOperation_t transb,
                           int m,
                           int n,
                           int k,
                           const void *alpha,
                           const void *A,
                           cudaDataType_t Atype,
                           int lda,
                           const void *B,
                           cudaDataType_t Btype,
                           int ldb,
                           const void *beta,
                           void *C,
                           cudaDataType_t Ctype,
                           int ldc,
                           cublasComputeType_t computeType,
                           cublasGemmAlgo_t algo)

#ifdef __cplusplus
cublasStatus_t cublasGemmEx(cublasHandle_t handle,
                           cublasOperation_t transa,

```

```

        cublasOperation_t transb,
        int m,
        int n,
        int k,
        const void      *alpha,
        const void      *A,
        cudaDataType     Atype,
        int lda,
        const void      *B,
        cudaDataType     Btype,
        int ldb,
        const void      *beta,
        void             *C,
        cudaDataType     Ctype,
        int ldc,
        cudaDataType     computeType,
        cublasGemmAlgo_t algo)
#endif

```

This function is an extension of `cublas<t>gemm` that allows the user to individually specify the data types for each of the A, B and C matrices, the precision of computation and the GEMM algorithm to be run. Supported combinations of arguments are listed further down in this section.

Note: The second variant of `cublasGemmEx` function is provided for backward compatibility with C++ applications code, where the `computeType` parameter is of `cudaDataType` instead of `cublasComputeType_t`. C applications would still compile with the updated function signature.

This function is only supported on devices with compute capability 5.0 or later.

$$C = \alpha \text{op}(A) \text{op}(B) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times k$ ,  $\text{op}(B) k \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}(B)$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}(A)$ and $C$ .
n		input	Number of columns of matrix $\text{op}(B)$ and $C$ .
k		input	Number of columns of $\text{op}(A)$ and rows of $\text{op}(B)$ .
alpha	host or device	input	Scaling factor for $A*B$ of the type that corresponds to the <code>computeType</code> and <code>Ctype</code> , see the table below for details.

Param.	Memory	In/out	Meaning
A	device	input	<type> array of dimensions $lda \times k$ with $lda \geq \max(1, m)$ if <code>transa == CUBLAS_OP_N</code> and $lda \times m$ with $lda \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, k)$ if <code>transb == CUBLAS_OP_N</code> and $ldb \times k$ with $ldb \geq \max(1, n)$ otherwise.
Btype		input	Enumerant specifying the datatype of matrix B.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	Scaling factor for C of the type that corresponds to the computeType and Ctype, see the table below for details. If <code>beta==0</code> , C does not have to be a valid input.
C	device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ .
Ctype		input	Enumerant specifying the datatype of matrix C.
ldc		input	Leading dimension of a two-dimensional array used to store the matrix C.
computeType		input	Enumerant specifying the computation type.
algo		input	Enumerant specifying the algorithm.

`cublasGemmEx` supports the following Compute Type, Scale Type, Atype/Btype, and Ctype:

Compute Type	Scale Type (alpha and beta)	Atype/Btype	Ctype
CUBLAS_COMPUTE_16F or CUBLAS_COMPUTE_16F_PEDANTIC	CUDA_R_16F	CUDA_R_16F	CUDA_R_16F
CUBLAS_COMPUTE_32I or CUBLAS_COMPUTE_32I_PEDANTIC	CUDA_R_32I	CUDA_R_8I	CUDA_R_32I
CUBLAS_COMPUTE_32F or CUBLAS_COMPUTE_32F_PEDANTIC	CUDA_R_32F	CUDA_R_16BF	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_16F
		CUDA_R_8I	CUDA_R_32F
		CUDA_R_16BF	CUDA_R_32F

Compute Type	Scale Type (alpha and beta)	Atype/Btype	Ctype
		CUDA_R_16F	CUDA_R_32F
		CUDA_R_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_8I	CUDA_C_32F
		CUDA_C_32F	CUDA_C_32F
CUBLAS_COMPUTE_32F_FAST_16F or CUBLAS_COMPUTE_32F_FAST_16BF or CUBLAS_COMPUTE_32F_FAST_TF32	CUDA_R_32F	CUDA_R_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_32F	CUDA_C_32F
CUBLAS_COMPUTE_64F or CUBLAS_COMPUTE_64F_PEDANTIC	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
	CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

Note: CUBLAS\_COMPUTE\_32I and CUBLAS\_COMPUTE\_32I\_PEDANTIC compute types are only supported with A, B being 4-byte aligned and lda, ldb being multiples of 4. For a better performance, it is also recommended that IMMA kernels requirements for a regular data ordering are met (listed [here](#)).

cublasGemmEx routine is run for the algorithms in the following table. Note: for NVIDIA Ampere Architecture GPUs and beyond, that is, SM version  $\geq 80$ , the algorithms below are equivalent to CUBLAS\_GEMM\_DEFAULT or CUBLAS\_GEMM\_DEFAULT\_TENSOR\_OP respectively. Specifying algorithm  $\geq 99$  for a single precision operation is equivalent to using CUBLAS\_COMPUTE\_32F\_FAST\_16F compute type, even if math mode or compute type are specified to be CUBLAS\_COMPUTE\_32F or CUBLAS\_COMPUTE\_32F\_FAST\_TF32.

CublasGemmAlgo_t	Meaning
CUBLAS_GEMM_DEFAULT	Apply Heuristics to select the GEMM algorithm.
CUBLAS_GEMM_ALGO0 to CUBLAS_GEMM_ALGO23	Explicitly choose an algorithm.
CUBLAS_GEMM_DEFAULT_TENSOR_OP	Apply Heuristics to select the GEMM algorithm while allowing the use of Tensor Core operations if possible.
CUBLAS_GEMM_ALGO0_TENSOR_OP to CUBLAS_GEMM_ALGO15_TENSOR_OP	Explicitly choose a GEMM algorithm allowing it to use Tensor Core operations if possible, otherwise falls back to cublas<t>gemmBatched based on computeType.

The possible error values returned by this function and their meanings are listed in the following table.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.

Error Value	Meaning
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ARCH_MISMATCH	cublasGemmEx is only supported for GPU with architecture capabilities equal or greater than 5.0.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters Atype, Btype and Ctype or the algorithm, algo is not supported.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if transa or transb != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if transa == CUBLAS_OP_N and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, k)</math> if transb == CUBLAS_OP_N and <math>ldb &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, m)</math> or</li> <li>▶ Atype or Btype or Ctype or algo is not supported</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

Starting with release 11.2, using the typed functions instead of the extension functions (`cublas**Ex()`) helps in reducing the binary size when linking to static cuBLAS Library.

Also refer to: [sgemm](#).

For more information about the numerical behavior of some GEMM algorithms, refer to the [GEMM Algorithms Numerical Behavior](#) section.

## 2.8.13. cublasGemmBatchedEx()

```

cublasStatus_t cublasGemmBatchedEx(cublasHandle_t handle,
                                   cublasOperation_t transa,
                                   cublasOperation_t transb,
                                   int m,
                                   int n,
                                   int k,
                                   const void *alpha,
                                   const void *Aarray[],
                                   cudaDataType_t Atype,
                                   int lda,
                                   const void *Barray[],
                                   cudaDataType_t Btype,
                                   int ldb,
                                   const void *beta,
                                   void *Carray[],
                                   cudaDataType_t Ctype,
                                   int ldc,
                                   int batchSize,
                                   cublasComputeType_t computeType,
                                   cublasGemmAlgo_t algo)

```

```

#ifdef __cplusplus
cublasStatus_t cublasGemmBatchedEx(cublasHandle_t handle,
                                   cublasOperation_t transa,
                                   cublasOperation_t transb,
                                   int m,
                                   int n,
                                   int k,
                                   const void *alpha,
                                   const void *Aarray[],
                                   cudaDataType Atype,
                                   int lda,
                                   const void *Barray[],
                                   cudaDataType Btype,
                                   int ldb,
                                   const void *beta,
                                   void *Carray[],
                                   cudaDataType Ctype,
                                   int ldc,
                                   int batchCount,
                                   cudaDataType computeType,
                                   cublasGemmAlgo_t algo)
#endif

```

This function is an extension of `cublas<T>gemmBatched` that performs the matrix-matrix multiplication of a batch of matrices and allows the user to individually specify the data types for each of the A, B and C matrix arrays, the precision of computation and the GEMM algorithm to be run. Like `cublas<T>gemmBatched`, the batch is considered to be "uniform", i.e. all instances have the same dimensions (m, n, k), leading dimensions (lda, ldb, ldc) and transpositions (transa, transb) for their respective A, B and C matrices. The address of the input matrices and the output matrix of each instance of the batch are read from arrays of pointers passed to the function by the caller. Supported combinations of arguments are listed further down in this section.

Note: The second variant of `cublasGemmBatchedEx` function is provided for backward compatibility with C++ applications code, where the `computeType` parameter is of `cudaDataType` instead of `cublasComputeType_t`. C applications would still compile with the updated function signature.

$$C[i] = \alpha \text{op}(A[i]) \text{op}(B[i]) + \beta C[i], \text{ for } i \in [0, \text{batchCount} - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are arrays of pointers to matrices stored in column-major format with dimensions  $\text{op}(A[i])\ m \times k$ ,  $\text{op}(B[i])\ k \times n$  and  $C[i]\ m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B[i])$  is defined similarly for matrix  $B[i]$ .

Note:  $C[i]$  matrices must not overlap, that is, the individual gemm operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<T>gemm` in different CUDA streams, rather than use this API.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}[A[i]]$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}[B[i]]$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}[A[i]]$ and $C[i]$ .
n		input	Number of columns of matrix $\text{op}[B[i]]$ and $C[i]$ .
k		input	Number of columns of $\text{op}[A[i]]$ and rows of $\text{op}[B[i]]$ .
alpha	host or device	input	Scaling factor for $A*B$ of the type that corresponds to the <code>computeType</code> and <code>Ctype</code> , see the table below for details.
Aarray	device	input	Array of pointers to <Atype> array, with each array of dim. $\text{lda} \times k$ with $\text{lda} \geq \max(1, m)$ if <code>transa == CUBLAS_OP_N</code> and $\text{lda} \times m$ with $\text{lda} \geq \max(1, k)$ otherwise.  All pointers must meet certain alignment criteria. Please see below for details.
Atype		input	Enumerant specifying the datatype of Aarray.
lda		input	Leading dimension of two-dimensional array used to store the matrix $A[i]$ .
Barray	device	input	Array of pointers to <Btype> array, with each array of dim. $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ if <code>transb == CUBLAS_OP_N</code> and $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ otherwise.  All pointers must meet certain alignment criteria. Please see below for details.
Btype		input	Enumerant specifying the datatype of Barray.
ldb		input	Leading dimension of two-dimensional array used to store matrix $B[i]$ .
beta	host or device	input	Scaling factor for $C$ of the type that corresponds to the <code>computeType</code> and <code>Ctype</code> , see the table below for details. If <code>beta==0</code> , $C[i]$ does not have to be a valid input.
Carray	device	in/out	Array of pointers to <Ctype> array. It has dimensions $\text{ldc} \times n$ with $\text{ldc} \geq \max(1, m)$ . Matrices $C[i]$ should not overlap; otherwise, undefined behavior is expected.  All pointers must meet certain alignment criteria. Please see below for details.
Ctype		input	Enumerant specifying the datatype of Carray.
ldc		input	Leading dimension of a two-dimensional array used to store each matrix $C[i]$ .



Param.	Memory	In/out	Meaning
batchCount		input	Number of pointers contained in Aarray, Barray and Carray.
computeType		input	Enumerant specifying the computation type.
algo		input	Enumerant specifying the algorithm.

cusblasGemmBatchedEx supports the following Compute Type, Scale Type, Atype/Btype, and Ctype:

Compute Type	Scale Type (alpha and beta)	Atype/Btype	Ctype
CUBLAS_COMPUTE_16F or CUBLAS_COMPUTE_16F_PEDANTIC	CUDA_R_16F	CUDA_R_16F	CUDA_R_16F
CUBLAS_COMPUTE_32I or CUBLAS_COMPUTE_32I_PEDANTIC	CUDA_R_32I	CUDA_R_8I	CUDA_R_32I
CUBLAS_COMPUTE_32F or CUBLAS_COMPUTE_32F_PEDANTIC	CUDA_R_32F	CUDA_R_16BF	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_16F
		CUDA_R_8I	CUDA_R_32F
		CUDA_R_16BF	CUDA_R_32F
		CUDA_R_16F	CUDA_R_32F
		CUDA_R_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_8I	CUDA_C_32F
		CUDA_C_32F	CUDA_C_32F
CUBLAS_COMPUTE_32F_FAST_16F or CUBLAS_COMPUTE_32F_FAST_16BF or CUBLAS_COMPUTE_32F_FAST_TF32	CUDA_R_32F	CUDA_R_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_32F	CUDA_C_32F
CUBLAS_COMPUTE_64F or CUBLAS_COMPUTE_64F_PEDANTIC	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
	CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

If Atype is CUDA\_R\_16F or CUDA\_R\_16BF, or computeType is any of the FAST options, or when math mode or algo enable fast math modes, pointers (not the pointer arrays) placed in the GPU memory must be properly aligned to avoid misaligned memory access errors. Ideally

all pointers are aligned to at least 16 Bytes. Otherwise it is recommended that they meet the following rule:

- if  $k \% 8 == 0$  then ensure  $\text{intptr\_t}(ptr) \% 16 == 0$ ,
- if  $k \% 2 == 0$  then ensure  $\text{intptr\_t}(ptr) \% 4 == 0$ .

Compute types `CUBLAS_COMPUTE_32I` and `CUBLAS_COMPUTE_32I_PEDANTIC` are only supported with all pointers `A[i]`, `B[i]` being 4-byte aligned and `lda`, `ldb` being multiples of 4. For a better performance, it is also recommended that IMMA kernels requirements for a regular data ordering are met (listed [here](#)).

`cublasGemmBatchedEx` routine is run for the algorithms in the following table. Note: for NVIDIA Ampere Architecture GPUs and beyond, that is, SM version  $\geq 80$ , the algorithms below are equivalent to `CUBLAS_GEMM_DEFAULT` or `CUBLAS_GEMM_DEFAULT_TENSOR_OP` respectively.

CublasGemmAlgo_t	Meaning
<code>CUBLAS_GEMM_DEFAULT</code>	Apply Heuristics to select the GEMM algorithm.
<code>CUBLAS_GEMM_ALGO0</code> to <code>CUBLAS_GEMM_ALGO23</code>	Explicitly choose an algorithm.
<code>CUBLAS_GEMM_DEFAULT_TENSOR_OP</code>	Apply Heuristics to select the GEMM algorithm while allowing the use of Tensor Core operations if possible.
<code>CUBLAS_GEMM_ALGO0_TENSOR_OP</code> to <code>CUBLAS_GEMM_ALGO15_TENSOR_OP</code>	Explicitly choose a GEMM algorithm allowing it to use Tensor Core operations if possible, otherwise falls back to <code>cublas&lt;t&gt;gemmBatched</code> based on <code>computeType</code> .

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	<code>cublasGemmBatchedEx</code> is only supported for GPU with architecture capabilities equal to or greater than 5.0.
<code>CUBLAS_STATUS_NOT_SUPPORTED</code>	The combination of the parameters <code>Atype</code> , <code>Btype</code> and <code>Ctype</code> or the algorithm, <code>algo</code> is not supported.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<ul style="list-style-type: none"> <li>► If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>► if <code>transa</code> or <code>transb</code> != <code>CUBLAS_OP_N</code>, <code>CUBLAS_OP_C</code>, <code>CUBLAS_OP_T</code> or</li> <li>► if <math>lda &lt; \max(1, m)</math> if <code>transa</code> == <code>CUBLAS_OP_N</code> and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>► if <math>ldb &lt; \max(1, k)</math> if <code>transb</code> == <code>CUBLAS_OP_N</code> and <math>ldb &lt; \max(1, n)</math> otherwise or</li> </ul>

Error Value	Meaning
	<ul style="list-style-type: none"> <li>▶ if <math>ldc &lt; \max(1, m)</math> or</li> <li>▶ Atype or Btype or Ctype or algo or computeType is not supported</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

Also refer to: [sgemm](#).

## 2.8.14. cublasGemmStridedBatchedEx()

```

cublasStatus_t cublasGemmStridedBatchedEx(cublasHandle_t handle,
                                           cublasOperation_t transa,
                                           cublasOperation_t transb,
                                           int m,
                                           int n,
                                           int k,
                                           const void *alpha,
                                           const void *A,
                                           cudaDataType_t Atype,
                                           int lda,
                                           long long int strideA,
                                           const void *B,
                                           cudaDataType_t Btype,
                                           int ldb,
                                           long long int strideB,
                                           const void *beta,
                                           void *C,
                                           cudaDataType_t Ctype,
                                           int ldc,
                                           long long int strideC,
                                           int batchSize,
                                           cublasComputeType_t computeType,
                                           cublasGemmAlgo_t algo)

#ifdef __cplusplus
cublasStatus_t cublasGemmStridedBatchedEx(cublasHandle_t handle,
                                           cublasOperation_t transa,
                                           cublasOperation_t transb,
                                           int m,
                                           int n,
                                           int k,
                                           const void *alpha,
                                           const void *A,
                                           cudaDataType_t Atype,
                                           int lda,
                                           long long int strideA,
                                           const void *B,
                                           cudaDataType_t Btype,
                                           int ldb,
                                           long long int strideB,
                                           const void *beta,
                                           void *C,
                                           cudaDataType_t Ctype,
                                           int ldc,
                                           long long int strideC,
                                           int batchSize,
                                           cudaDataType_t computeType,
                                           cublasGemmAlgo_t algo)
#endif

```

This function is an extension of `cublas<t>gemmStridedBatched` that performs the matrix-matrix multiplication of a batch of matrices and allows the user to individually specify the data types for each of the A, B and C matrices, the precision of computation and the GEMM algorithm to be run. Like `cublas<t>gemmStridedBatched`, the batch is considered to be "uniform", that is, all instances have the same dimensions (m, n, k), leading dimensions (lda, ldb, ldc) and transpositions (transa, transb) for their respective A, B and C matrices. Input matrices A, B and output matrix C for each instance of the batch are located at fixed offsets in number of elements from their locations in the previous instance. Pointers to A, B and C matrices for the first instance are passed to the function by the user along with the offsets in number of elements - `strideA`, `strideB` and `strideC` that determine the locations of input and output matrices in future instances.

Note: The second variant of `cublasGemmStridedBatchedEx` function is provided for backward compatibility with C++ applications code, where the `computeType` parameter is of `cudaDataType_` instead of `cublasComputeType_t`. C applications would still compile with the updated function signature.

$$C + i * \text{strideC} = \alpha \text{op}(A + i * \text{strideA}) \text{op}(B + i * \text{strideB}) + \beta(C + i * \text{strideC}), \text{ for } i \in [0, \text{batchCount} - 1]$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are arrays of pointers to matrices stored in column-major format with dimensions  $\text{op}(A[i])\ m \times k$ ,  $\text{op}(B[i])\ k \times n$  and  $C[i]\ m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B[i])$  is defined similarly for matrix  $B[i]$ .

Note:  $C[i]$  matrices must not overlap, i.e. the individual gemm operations must be computable independently; otherwise, undefined behavior is expected.

On certain problem sizes, it might be advantageous to make multiple calls to `cublas<t>gemm` in different CUDA streams, rather than use this API.

Note: In the table below, we use  $A[i]$ ,  $B[i]$ ,  $C[i]$  as notation for A, B and C matrices in the  $i$ th instance of the batch, implicitly assuming they are respectively offsets in number of elements `strideA`, `strideB`, `strideC` away from  $A[i-1]$ ,  $B[i-1]$ ,  $C[i-1]$ . The unit for the offset is number of elements and must not be zero.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
transa		input	Operation $\text{op}(A[i])$ that is non- or (conj.) transpose.
transb		input	Operation $\text{op}(B[i])$ that is non- or (conj.) transpose.
m		input	Number of rows of matrix $\text{op}(A[i])$ and $C[i]$ .
n		input	Number of columns of matrix $\text{op}(B[i])$ and $C[i]$ .
k		input	Number of columns of $\text{op}(A[i])$ and rows of $\text{op}(B[i])$ .
alpha	host or device	input	Scaling factor for $A*B$ of the type that corresponds to the <code>computeType</code> and <code>Ctype</code> , see the table below for details.

Param.	Memory	In/out	Meaning
A	device	input	Pointer to <Atype> matrix, A, corresponds to the first instance of the batch, with dimensions $l_{da} \times k$ with $l_{da} \geq \max(1, m)$ if <code>transa == CUBLAS_OP_N</code> and $l_{da} \times m$ with $l_{da} \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of A.
lda		input	Leading dimension of two-dimensional array used to store the matrix A[i].
strideA		input	Value of type long long int that gives the offset in number of elements between A[i] and A[i+1].
B	device	input	Pointer to <Btype> matrix, B, corresponds to the first instance of the batch, with dimensions $l_{db} \times n$ with $l_{db} \geq \max(1, k)$ if <code>transb == CUBLAS_OP_N</code> and $l_{db} \times k$ with $l_{db} \geq \max(1, n)$ otherwise.
Btype		input	Enumerant specifying the datatype of B.
ldb		input	Leading dimension of two-dimensional array used to store matrix B[i].
strideB		input	Value of type long long int that gives the offset in number of elements between B[i] and B[i+1].
beta	host or device	input	Scaling factor for c of the type that corresponds to the computeType and Ctype, see the table below for details. If <code>beta==0</code> , C[i] does not have to be a valid input.
C	device	in/out	Pointer to <Ctype> matrix, c, corresponds to the first instance of the batch, with dimensions $l_{dc} \times n$ with $l_{dc} \geq \max(1, m)$ . Matrices C[i] should not overlap; otherwise, undefined behavior is expected.
Ctype		input	Enumerant specifying the datatype of c.
ldc		input	Leading dimension of a two-dimensional array used to store each matrix C[i].
strideC		input	Value of type long long int that gives the offset in number of elements between C[i] and C[i+1].
batchCount		input	Number of GEMMs to perform in the batch.
computeType		input	Enumerant specifying the computation type.
algo		input	Enumerant specifying the algorithm.

`cublasGemmStridedBatchedEx` supports the following Compute Type, Scale Type, Atype/ Btype, and Ctype:

Compute Type	Scale Type (alpha and beta)	Atype/Btype	Ctype
CUBLAS_COMPUTE_16F or CUBLAS_COMPUTE_16F_PEDANTIC	CUDA_R_16F	CUDA_R_16F	CUDA_R_16F
CUBLAS_COMPUTE_32I or CUBLAS_COMPUTE_32I_PEDANTIC	CUDA_R_32I	CUDA_R_8I	CUDA_R_32I
CUBLAS_COMPUTE_32F or CUBLAS_COMPUTE_32F_PEDANTIC	CUDA_R_32F	CUDA_R_16BF	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_16F
		CUDA_R_8I	CUDA_R_32F
		CUDA_R_16BF	CUDA_R_32F
		CUDA_R_16F	CUDA_R_32F
		CUDA_R_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_8I	CUDA_C_32F
		CUDA_C_32F	CUDA_C_32F
CUBLAS_COMPUTE_32F_FAST_16F or CUBLAS_COMPUTE_32F_FAST_16BF or CUBLAS_COMPUTE_32F_FAST_TF32	CUDA_R_32F or CUDA_C_32F	CUDA_R_32F or CUDA_C_32F	CUDA_R_32F or CUDA_C_32F
CUBLAS_COMPUTE_64F or CUBLAS_COMPUTE_64F_PEDANTIC	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
	CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

Compute types CUBLAS\_COMPUTE\_32I and CUBLAS\_COMPUTE\_32I\_PEDANTIC are only supported with all pointers  $A[i]$ ,  $B[i]$  being 4-byte aligned and  $lda$ ,  $ldb$  being multiples of 4. For a better performance, it is also recommended that IMMA kernels requirements for a regular data ordering are met (listed [here](#)).

cublasGemmStridedBatchedEx routine is run for the algorithms in the following table. Note: for NVIDIA Ampere Architecture GPUs and beyond, that is, SM version  $\geq 80$ , the algorithms below are equivalent to CUBLAS\_GEMM\_DEFAULT or CUBLAS\_GEMM\_DEFAULT\_TENSOR\_OP respectively.

CublasGemmAlgo_t	Meaning
CUBLAS_GEMM_DEFAULT	Apply Heuristics to select the GEMM algorithm.
CUBLAS_GEMM_ALGO0 to CUBLAS_GEMM_ALGO23	Explicitly choose an algorithm.

CublasGemmAlgo_t	Meaning
CUBLAS_GEMM_DEFAULT_TENSOR_OP	Apply Heuristics to select the GEMM algorithm while allowing the use of Tensor Core operations if possible.
CUBLAS_GEMM_ALGO0_TENSOR_OP to CUBLAS_GEMM_ALGO15_TENSOR_OP	Explicitly choose a GEMM algorithm allowing it to use Tensor Core operations if possible, otherwise falls back to cublas<t>gemmBatched based on computeType.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ARCH_MISMATCH	cublasGemmBatchedEx is only supported for GPU with architecture capabilities equal to or greater than 5.0.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters Atype, Btype and Ctype or the algorithm, algo is not supported.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>m &lt; 0</math> or <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if transa or transb != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if <math>lda &lt; \max(1, m)</math> if transa == CUBLAS_OP_N and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldb &lt; \max(1, k)</math> if transb == CUBLAS_OP_N and <math>ldb &lt; \max(1, n)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, m)</math> or</li> <li>▶ Atype or Btype or Ctype or algo or computeType is not supported</li> </ul>
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

Also refer to: [sgemm](#).

## 2.8.15. cublasCsyrrkEx()

```
cublasStatus_t cublasCsyrrkEx(cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             cublasOperation_t trans,
                             int n,
                             int k,
                             const float *alpha,
                             const void *A,
```

```

        cudaDataType  Atype,
        int  lda,
        const float  *beta,
        cuComplex  *C,
        cudaDataType  Ctype,
        int  ldc)

```

This function is an extension of `cublasCsyrk` where the input matrix and output matrix can have a lower precision but the computation is still done in the type `cuComplex`.

This function performs the symmetric rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \, n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \end{cases}$$



**Note:** This routine is only supported on GPUs with architecture capabilities equal to or greater than 5.0.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix c lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op[A] that is non- or transpose.
n		input	Number of rows of matrix op[A] and c.
k		input	Number of columns of matrix op[A].
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1, n)$ if $trans == \text{CUBLAS\_OP\_N}$ and $lda \times n$ with $lda \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta	host or device	input	<type> scalar used for multiplication, if $\beta == 0$ then c does not have to be a valid input.
C	device	in/out	<type> array of dimension $ldc \times n$ , with $ldc \geq \max(1, n)$ .
Ctype		input	Enumerant specifying the datatype of matrix c.
ldc		input	Leading dimension of two-dimensional array used to store matrix c.



The matrix types combinations supported for `cublasCsyrkEx` are listed below:

A	C
CUDA_C_8I	CUDA_C_32F
CUDA_C_32F	CUDA_C_32F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <code>uplo</code> != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER or</li> <li>▶ if <code>trans</code> != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> if <code>trans</code> == CUBLAS_OP_N and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, n)</math> or</li> <li>▶ <code>Atype</code> or <code>Ctype</code> is not supported</li> </ul>
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>Atype</code> and <code>Ctype</code> is not supported.
CUBLAS_STATUS_ARCH_MISMATCH	The device has a compute capability lower than 5.0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#)

## 2.8.16. cublasCsyrk3mEx()

```
cublasStatus_t cublasCsyrk3mEx(cublasHandle_t handle,
                               cublasFillMode_t uplo,
                               cublasOperation_t trans,
                               int n,
                               int k,
                               const float      *alpha,
                               const void        *A,
                               cudaDataType      Atype,
                               int lda,
                               const float      *beta,
                               cuComplex         *C,
                               cudaDataType      Ctype,
                               int ldc)
```

This function is an extension of  `cublasCsyrrk`  where the input matrix and output matrix can have a lower precision but the computation is still done in the type  `cuComplex` . This routine is implemented using the Gauss complexity reduction algorithm which can lead to an increase in performance up to 25%.

This function performs the symmetric rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \ n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \end{cases}$$



**Note:** This routine is only supported on GPUs with architecture capabilities equal to or greater than 5.0

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix c lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or transpose.
n		input	Number of rows of matrix op(A) and c.
k		input	Number of columns of matrix op(A).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1, n)$ if $trans == \text{CUBLAS\_OP\_N}$ and $lda \times n$ with $lda \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta	host or device	input	<type> scalar used for multiplication, if $beta == 0$ then c does not have to be a valid input.
C	device	in/out	<type> array of dimension $ldc \times n$ , with $ldc \geq \max(1, n)$ .
Ctype		input	Enumerant specifying the datatype of matrix c.
ldc		input	Leading dimension of two-dimensional array used to store matrix c.

The matrix types combinations supported for  `cublasCsyrrk3mEx`  are listed below:

A	C
CUDA_C_8I	CUDA_C_32F
CUDA_C_32F	CUDA_C_32F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER}, \text{CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> if <math>trans == \text{CUBLAS\_OP\_N}</math> and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, n)</math> or</li> <li>▶ Atype or Ctype is not supported</li> </ul>
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters Atype and Ctype is not supported.
CUBLAS_STATUS_ARCH_MISMATCH	The device has a compute capability lower than 5.0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#)

## 2.8.17. cublasCherkEx()

```
cublasStatus_t cublasCherkEx(cublasHandle_t handle,
                             cublasFillMode_t uplo,
                             cublasOperation_t trans,
                             int n,
                             int k,
                             const float *alpha,
                             const void *A,
                             cudaDataType Atype,
                             int lda,
                             const float *beta,
                             cuComplex *C,
                             cudaDataType Ctype,
                             int ldc)
```

This function is an extension of  `cublasCherk`  where the input matrix and output matrix can have a lower precision but the computation is still done in the type  `cuComplex` .

This function performs the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \, n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$



**Note:** This routine is only supported on GPUs with architecture capabilities equal to or greater than 5.0

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.
n		input	Number of rows of matrix $\text{op}(A)$ and $C$ .
k		input	Number of columns of matrix $\text{op}(A)$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times k$ with $\text{lda} \geq \max(1, n)$ if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $\text{lda} \times n$ with $\text{lda} \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta		input	<type> scalar used for multiplication, if $\text{beta} == 0$ then $C$ does not have to be a valid input.
C	device	in/out	<type> array of dimension $\text{ldc} \times n$ , with $\text{ldc} \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
Ctype		input	Enumerant specifying the datatype of matrix C.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The matrix types combinations supported for  `cublasCherkEx`  are listed in the following table:

A	C
CUDA_C_8I	CUDA_C_32F
CUDA_C_32F	CUDA_C_32F

The possible error values returned by this function and their meanings are listed in the following table.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <math>uplo \neq \text{CUBLAS\_FILL\_MODE\_UPPER}, \text{CUBLAS\_FILL\_MODE\_LOWER}</math> or</li> <li>▶ if <math>trans \neq \text{CUBLAS\_OP\_N}, \text{CUBLAS\_OP\_C}, \text{CUBLAS\_OP\_T}</math> or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> if <math>trans == \text{CUBLAS\_OP\_N}</math> and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, n)</math> or</li> <li>▶ <math>Atype</math> or <math>Ctype</math> is not supported</li> </ul>
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters $Atype$ and $Ctype$ is not supported.
CUBLAS_STATUS_ARCH_MISMATCH	The device has a compute capability lower than 5.0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#)

## 2.8.18. cublasCherk3mEx()

```

cublasStatus_t cublasCherk3mEx(cublasHandle_t handle,
                               cublasFillMode_t uplo,
                               cublasOperation_t trans,
                               int n,
                               int k,
                               const float *alpha,
                               const void *A,
                               cudaDataType Atype,
                               int lda,
                               const float *beta,
                               cuComplex *C,
                               cudaDataType Ctype,
                               int ldc)

```

This function is an extension of `cusblasCherk` where the input matrix and output matrix can have a lower precision but the computation is still done in the type `cuComplex`. This routine is implemented using the Gauss complexity reduction algorithm which can lead to an increase in performance up to 25%.

This function performs the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \ n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$



**Note:** This routine is only supported on GPUs with architecture capabilities equal to or greater than 5.0

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(A)$ that is non- or (conj.) transpose.
n		input	Number of rows of matrix $\text{op}(A)$ and $C$ .
k		input	Number of columns of matrix $\text{op}(A)$ .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times k$ with $\text{lda} \geq \max(1, n)$ if $\text{transa} == \text{CUBLAS\_OP\_N}$ and $\text{lda} \times n$ with $\text{lda} \geq \max(1, k)$ otherwise.
Atype		input	Enumerant specifying the datatype of matrix A.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta		input	<type> scalar used for multiplication, if $\text{beta} == 0$ then $C$ does not have to be a valid input.
C	device	in/out	<type> array of dimension $\text{ldc} \times n$ , with $\text{ldc} \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
Ctype		input	Enumerant specifying the datatype of matrix C.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The matrix types combinations supported for `cublasCherk3mEx` are listed in the following table:

A	C
CUDA_C_8I	CUDA_C_32F
CUDA_C_32F	CUDA_C_32F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <math>n &lt; 0</math> or <math>k &lt; 0</math> or</li> <li>▶ if <code>uplo</code> != CUBLAS_FILL_MODE_UPPER, CUBLAS_FILL_MODE_LOWER or</li> <li>▶ if <code>trans</code> != CUBLAS_OP_N, CUBLAS_OP_C, CUBLAS_OP_T or</li> <li>▶ if <math>lda &lt; \max(1, n)</math> if <code>trans</code> == CUBLAS_OP_N and <math>lda &lt; \max(1, k)</math> otherwise or</li> <li>▶ if <math>ldc &lt; \max(1, n)</math> or</li> <li>▶ <code>Atype</code> or <code>Ctype</code> is not supported</li> </ul>
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>Atype</code> and <code>Ctype</code> is not supported.
CUBLAS_STATUS_ARCH_MISMATCH	The device has a compute capability lower than 5.0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#)

## 2.8.19. cublasNrm2Ex()

```
cublasStatus_t cublasNrm2Ex( cublasHandle_t handle,
                             int n,
                             const void *x,
                             cudaDataType xType,
                             int incx,
                             void *result,
                             cudaDataType resultType,
                             cudaDataType executionType)
```

This function is an API generalization of the routine `cublas<t>nrm2` where input data, output data and compute type can be specified independently.

This function computes the Euclidean norm of the vector `x`. The code uses a multiphase model of accumulation to avoid intermediate underflow and overflow, with the result being equivalent

to  $\sqrt{\sum_{i=1}^n (\mathbf{x}[j] \times \mathbf{x}[j])}$  where  $j = 1 + (i - 1) * \text{incx}$  in exact arithmetic. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vector <code>x</code> .
x	device	input	<type> vector with <code>n</code> elements.
xType		input	Enumerant specifying the datatype of vector <code>x</code> .
incx		input	Stride between consecutive elements of <code>x</code> .
result	host or device	output	The resulting norm, which is 0.0 if <code>n, incx &lt;= 0</code> .
resultType		input	Enumerant specifying the datatype of the <code>result</code> .
executionType		input	Enumerant specifying the datatype in which the computation is executed.

The datatypes combinations currently supported for `cublasNrm2Ex` are listed in the following table:

x	result	execution
CUDA_R_16F	CUDA_R_16F	CUDA_R_32F
CUDA_R_16BF	CUDA_R_16BF	CUDA_R_32F
CUDA_R_32F	CUDA_R_32F	CUDA_R_32F
CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
CUDA_C_32F	CUDA_C_32F	CUDA_C_32F
CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.



Error Value	Meaning
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>xType</code> , <code>resultType</code> and <code>executionType</code> is not supported.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If <code>xType</code> or <code>resultType</code> or <code>executionType</code> is not supported or</li> <li>▶ <code>result == NULL</code></li> </ul>

For references please refer to:

`snrm2`, `snrm2`, `dnrm2`, `dnrm2`, `scnrm2`, `scnrm2`, `dznrm2`

## 2.8.20. `cublasAxyEx()`

```
cublasStatus_t cublasAxyEx (cublasHandle_t handle,
                           int n,
                           const void *alpha,
                           cudaDataType alphaType,
                           const void *x,
                           cudaDataType xType,
                           int incx,
                           void *y,
                           cudaDataType yType,
                           int incy,
                           cudaDataType executionType);
```

This function is an API generalization of the routine `cublas<t>axy` where input data, output data and compute type can be specified independently.

This function multiplies the vector  $x$  by the scalar  $\alpha$  and adds it to the vector  $y$  overwriting the latest vector with the result. Hence, the performed operation is  $y[j] = \alpha \times x[k] + y[j]$  for  $i = 1, \dots, n$ ,  $k = 1 + (i - 1) * \text{incx}$  and  $j = 1 + (i - 1) * \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
<code>handle</code>		input	Handle to the cuBLAS library context.
<code>alpha</code>	host or device	input	<type> scalar used for multiplication.
<code>n</code>		input	Number of elements in the vector $x$ and $y$ .
<code>x</code>	device	input	<type> vector with $n$ elements.
<code>xType</code>		input	Enumerant specifying the datatype of vector $x$ .
<code>incx</code>		input	Stride between consecutive elements of $x$ .
<code>y</code>	device	in/out	<type> vector with $n$ elements.
<code>yType</code>		input	Enumerant specifying the datatype of vector $y$ .

Param.	Memory	In/out	Meaning
incy		input	Stride between consecutive elements of y.
executionType		input	Enumerant specifying the datatype in which the computation is executed.

The datatypes combinations currently supported for `cublasAxpEx` are listed in the following table:

x	y	execution
CUDA_R_16F	CUDA_R_16F	CUDA_R_32F
CUDA_R_16BF	CUDA_R_16BF	CUDA_R_32F
CUDA_R_32F	CUDA_R_32F	CUDA_R_32F
CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
CUDA_C_32F	CUDA_C_32F	CUDA_C_32F
CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>xType</code> , <code>yType</code> , and <code>executionType</code> is not supported.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	<code>alphaType</code> or <code>xType</code> or <code>yType</code> or <code>executionType</code> is not supported.

For references please refer to:

[saxpy](#), [daxpy](#), [caxpy](#), [zaxpy](#)

## 2.8.21. `cublasDotEx()`

```
cublasStatus_t cublasDotEx (cublasHandle_t handle,
    int n,
    const void *x,
    cudaDataType xType,
    int incx,
    const void *y,
    cudaDataType yType,
    int incy,
    void *result,
    cudaDataType resultType,
    cudaDataType executionType);
```

```

cublasStatus_t cublasDotcEx (cublasHandle_t handle,
                             int n,
                             const void *x,
                             cudaDataType xType,
                             int incx,
                             const void *y,
                             cudaDataType yType,
                             int incy,
                             void *result,
                             cudaDataType resultType,
                             cudaDataType executionType);

```

These functions are an API generalization of the routines `cublas<t>dot` and `cublas<t>dotc` where input data, output data and compute type can be specified independently. Note: `cublas<t>dotc` is dot product conjugated, `cublas<t>dotu` is dot product unconjugated.

This function computes the dot product of vectors  $x$  and  $y$ . Hence, the result is  $\sum_{i=1}^n (x[k] \times y[j])$  where  $k = 1 + (i - 1) * \text{incx}$  and  $j = 1 + (i - 1) * \text{incy}$ . Notice that in the first equation the conjugate of the element of vector  $x$  should be used if the function name ends in character 'c' and that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vectors $x$ and $y$ .
x	device	input	<type> vector with $n$ elements.
xType		input	Enumerant specifying the datatype of vector $x$ .
incx		input	Stride between consecutive elements of $x$ .
y	device	input	<type> vector with $n$ elements.
yType		input	Enumerant specifying the datatype of vector $y$ .
incy		input	Stride between consecutive elements of $y$ .
result	host or device	output	The resulting dot product, which is 0.0 if $n \leq 0$ .
resultType		input	Enumerant specifying the datatype of the <code>result</code> .
executionType		input	Enumerant specifying the datatype in which the computation is executed.

The datatypes combinations currently supported for `cublasDotEx` and `cublasDotcEx` are listed in the following table:

x	y	result	execution
CUDA_R_16F	CUDA_R_16F	CUDA_R_16F	CUDA_R_32F
CUDA_R_16BF	CUDA_R_16BF	CUDA_R_16BF	CUDA_R_32F
CUDA_R_32F	CUDA_R_32F	CUDA_R_32F	CUDA_R_32F

x	y	result	execution
CUDA_R_64F	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
CUDA_C_32F	CUDA_C_32F	CUDA_C_32F	CUDA_C_32F
CUDA_C_64F	CUDA_C_64F	CUDA_C_64F	CUDA_C_64F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_ALLOC_FAILED	The reduction buffer could not be allocated.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>xType</code> , <code>yType</code> , <code>resultType</code> and <code>executionType</code> is not supported.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	<code>xType</code> or <code>yType</code> or <code>resultType</code> or <code>executionType</code> is not supported.

For references please refer to:

[sdot](#), [ddot](#), [cdotu](#), [cdotc](#), [zdotu](#), [zdotc](#)

## 2.8.22. cublasRotEx()

```
cublasStatus_t cublasRotEx(cublasHandle_t handle,
    int n,
    void *x,
    cudaDataType xType,
    int incx,
    void *y,
    cudaDataType yType,
    int incy,
    const void *c, /* host or device pointer */
    const void *s,
    cudaDataType csType,
    cudaDataType executionType);
```

This function is an extension to the routine `cublas<t>rot` where input data, output data, cosine/sine type, and compute type can be specified independently.

This function applies Givens rotation matrix (that is, rotation in the x,y plane counter-clockwise by angle defined by  $\cos(\alpha)=c$ ,  $\sin(\alpha)=s$ ):

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

to vectors  $x$  and  $y$ .

Hence, the result is  $\mathbf{x}[k] = c \times \mathbf{x}[k] + s \times \mathbf{y}[j]$  and  $\mathbf{y}[j] = -s \times \mathbf{x}[k] + c \times \mathbf{y}[j]$  where  $k = 1 + (i - 1) \times \text{incx}$  and  $j = 1 + (i - 1) \times \text{incy}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
n		input	Number of elements in the vectors $\mathbf{x}$ and $\mathbf{y}$ .
x	device	in/out	<type> vector with $n$ elements.
xType		input	Enumerant specifying the datatype of vector $\mathbf{x}$ .
incx		input	Stride between consecutive elements of $\mathbf{x}$ .
y	device	in/out	<type> vector with $n$ elements.
yType		input	Enumerant specifying the datatype of vector $\mathbf{y}$ .
incy		input	Stride between consecutive elements of $\mathbf{y}$ .
c	host or device	input	Cosine element of the rotation matrix.
s	host or device	input	Sine element of the rotation matrix.
csType		input	Enumerant specifying the datatype of $c$ and $s$ .
executionType		input	Enumerant specifying the datatype in which the computation is executed.

The datatypes combinations currently supported for `cublasRotEx` are listed in the following table:

executionType	xType / yType	csType
CUDA_R_32F	CUDA_R_16BF	CUDA_R_16BF
	CUDA_R_16F	CUDA_R_16F
	CUDA_R_32F	CUDA_R_32F
CUDA_R_64F	CUDA_R_64F	CUDA_R_64F
CUDA_C_32F	CUDA_C_32F	CUDA_R_32F
	CUDA_C_32F	CUDA_C_32F
CUDA_C_64F	CUDA_C_64F	CUDA_R_64F
	CUDA_C_64F	CUDA_C_64F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[srot](#), [drot](#), [crot](#), [csrot](#), [zrot](#), [zdrot](#)

## 2.8.23. cublasScalex()

```
cublasStatus_t cublasScalex(cublasHandle_t handle,
                           int n,
                           const void *alpha,
                           cudaDataType alphaType,
                           void *x,
                           cudaDataType xType,
                           int incx,
                           cudaDataType executionType);
```

This function scales the vector  $x$  by the scalar  $\alpha$  and overwrites it with the result. Hence, the performed operation is  $x[j] = \alpha \times x[j]$  for  $i = 1, \dots, n$  and  $j = 1 + (i - 1) * \text{incx}$ . Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLAS library context.
alpha	host or device	input	<type> scalar used for multiplication.
n		input	Number of elements in the vector $x$ .
x	device	in/out	<type> vector with $n$ elements.
xType		input	Enumerant specifying the datatype of vector $x$ .
incx		input	Stride between consecutive elements of $x$ .
executionType		input	Enumerant specifying the datatype in which the computation is executed.

The datatypes combinations currently supported for cublasScalex are listed in the following table:

x	execution
CUDA_R_16F	CUDA_R_32F
CUDA_R_16BF	CUDA_R_32F
CUDA_R_32F	CUDA_R_32F
CUDA_R_64F	CUDA_R_64F
CUDA_C_32F	CUDA_C_32F
CUDA_C_64F	CUDA_C_64F

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_NOT_SUPPORTED	The combination of the parameters <code>xType</code> and <code>executionType</code> is not supported.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.
CUBLAS_STATUS_INVALID_VALUE	<code>alphaType</code> or <code>xType</code> or <code>executionType</code> is not supported.

For references please refer to:

[sscal](#), [dscal](#), [csscal](#), [cscal](#), [zdscal](#), [zscal](#)

---

# Chapter 3. Using the cuBLASLt API

## 3.1. General Description

The cuBLASLt is a new lightweight library dedicated to GEneral Matrix-to-matrix Multiply (GEMM) operations with a new flexible API. This new library adds flexibility in matrix data layouts, input types, compute types, and also in choosing the algorithmic implementations and heuristics through parameter programmability.

Once a set of options for the intended GEMM operation are identified by the user, these options can be used repeatedly for different inputs. This is analogous to how cuFFT and FFTW first create a plan and reuse for same size and type FFTs with different input data.



**Note:** The cuBLASLt in general does not guarantee to support all possible sizes and configurations. The main focus of the library is to provide the most performant kernels, which might have some implied limitations. Some non-standard configurations may require a user to handle them manually, typically by decomposing the problem into smaller parts.

For example, most of the kernels do not support batch size greater than 65535 due to the limitation of the block grid's z dimension. To overcome this problem, a user may want to launch kernels for each subbatch of maximum size 65535 and combine the results.

### 3.1.1. Heuristics Cache

cuBLASLt uses heuristics to pick the most suitable matmul kernel for execution based on the problem sizes, GPU configuration, and other parameters. This requires performing some computations on the host CPU, which could take tens of microseconds. To overcome this overhead, it is recommended to query the heuristics once using [`cublasLtMatmulAlgoGetHeuristic\(\)`](#) and then reuse the result for subsequent computations using [`cublasLtMatmul\(\)`](#).

For the cases where querying heuristics once and then reusing them is not feasible, cuBLASLt implements a heuristics cache that maps matmul problems to kernels previously selected by heuristics. The heuristics cache uses an LRU-like eviction policy and is thread-safe.

The user can control the heuristics cache capacity with the `CUBLASLT_HEURISTICS_CACHE_CAPACITY` environment variable or with the [`cublasLtHeuristicsCacheSetCapacity\(\)`](#) function which has higher precedence. The



capacity is measured in number of entries. Each entry takes about 320 bytes but is subject to change. The default capacity is 8192 entries.

See also: [cublasLtHeuristicsCacheGetCapacity\(\)](#),  
[cublasLtHeuristicsCacheSetCapacity\(\)](#).

### 3.1.2. cuBLASLt Logging

The cuBLASLt logging mechanism can be enabled by setting the following environment variables before launching the target application:

CUBLASLT\_LOG\_LEVEL=<level> - while level is one of the following levels:

- ▶ "0" - Off - logging is disabled (default)
- ▶ "1" - Error - only errors will be logged
- ▶ "2" - Trace - API calls that launch CUDA kernels will log their parameters and important information
- ▶ "3" - Hints - hints that can potentially improve the application's performance
- ▶ "4" - Info - provides general information about the library execution, may contain details about heuristic status
- ▶ "5" - API Trace - API calls will log their parameter and important information

CUBLASLT\_LOG\_MASK=<mask> - while mask is a combination of the following masks:

- ▶ "0" - Off
- ▶ "1" - Error
- ▶ "2" - Trace
- ▶ "4" - Hints
- ▶ "8" - Info
- ▶ "16" - API Trace

CUBLASLT\_LOG\_FILE=<file\_name> - while file name is a path to a logging file. File name may contain %i, that will be replaced with the process ID. For example, "<file\_name>\_%i.log".

If CUBLASLT\_LOG\_FILE is not defined, the log messages are printed to stdout.

Another option is to use the experimental cublasLt logging API. See:

[cublasLtLoggerSetCallback\(\)](#), [cublasLtLoggerSetFile\(\)](#),  
[cublasLtLoggerOpenFile\(\)](#), [cublasLtLoggerSetLevel\(\)](#), [cublasLtLoggerSetMask\(\)](#),  
[cublasLtLoggerForceDisable\(\)](#)

### 3.1.3. 8-bit Floating Point Data Types (FP8) Usage

FP8 was first introduced with Hopper GPUs (compute capability  $\geq$  sm\_90) and is designed to further accelerate matrix multiplications. There are two types of FP8 available:

1. `CUDA_R_8F_E4M3` is designed to be accurate at a smaller dynamic range than half precision. The E4 and M3 represent a 4-bit exponent and a 3-bit mantissa respectively. For more details, see [\\_\\_nv\\_fp8\\_e4m3](#).
2. `CUDA_R_8F_E5M2` is designed to be accurate at a similar dynamic range as half precision. The E5 and M2 represent a 5-bit exponent and a 2-bit mantissa respectively. For more information see [\\_\\_nv\\_fp8\\_e5m2](#).



**Note:** Unless otherwise stated, FP8 refers to both `CUDA_R_8F_E4M3` and `CUDA_R_8F_E5M2`

In order to maintain accurate FP8 matrix multiplications, we define native compute FP8 matrix multiplication as follows:

$$D = \text{scaledD} * (\alpha * ((\text{scaleA} * A) * (\text{scaleB} * B)) + \beta * (\text{scaleC} * C))$$

where A, B, and C are input matrices, and scaleA, scaleB, scaleC, scaledD, alpha, and beta are input scalars. This differs from the other matrix multiplication routines because of this addition of scaling factors for each matrix. For more information on FP8, see [cublasLtMatmul\(\)](#) and [cublasLtMatmulDescAttributes\\_t](#).

For FP8 matrix multiplications, epilogues and amaxD may be computed as follows:

$$D\_temp = \text{Epilogue}(\alpha * ((\text{scaleA} * A) * (\text{scaleB} * B)) + \beta * (\text{scaleC} * C))$$

$$\text{amaxD} = \text{absmax}(D\_temp)$$

$$D = \text{scaledD} * D\_temp$$

For more information, see attribute `CUBLASLT_MATMUL_DESC_AMAX_D_POINTER` in [cublasLtMatmulDescAttributes\\_t](#).

## 3.2. cuBLASLt Code Examples

Please visit <https://github.com/NVIDIA/CUDALibrarySamples/tree/master/cuBLASLt> for updated code examples.

## 3.3. cuBLASLt Datatypes Reference

### 3.3.1. cublasLt3mMode\_t

`cublasLt3mMode_t` is an enumerated type used for computation with complex matrices. This enumerated type can be used to apply the Gaussian complexity reduction algorithm.

Value	Description
<code>CUBLASLT_3M_MODE_DISALLOWED</code>	Gaussian complexity reduction algorithm is not applied to the matrix-matrix computation.
<code>CUBLASLT_3M_MODE_ALLOWED</code>	Gaussian complexity reduction algorithm can be applied to the matrix-matrix computation.

### 3.3.2. cublasLtClusterShape\_t

`cublasLtClusterShape_t` is an enumerated type used to configure thread block cluster dimensions. Thread block clusters add an optional hierarchical level and are made up of thread blocks. Similar to thread blocks, these can be one, two, or three-dimensional. See also [Thread Block Clusters](#).

Value	Description
CUBLASLT_CLUSTER_SHAPE_AUTO	Cluster shape is automatically selected.
CUBLASLT_CLUSTER_SHAPE_1x1x1	Cluster shape is 1 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_1x2x1	Cluster shape is 1 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_1x4x1	Cluster shape is 1 x 4 x 1.
CUBLASLT_CLUSTER_SHAPE_2x1x1	Cluster shape is 2 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_2x2x1	Cluster shape is 2 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_2x4x1	Cluster shape is 2 x 4 x 1.
CUBLASLT_CLUSTER_SHAPE_4x1x1	Cluster shape is 4 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_4x2x1	Cluster shape is 4 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_4x4x1	Cluster shape is 4 x 4 x 1.
CUBLASLT_CLUSTER_SHAPE_1x8x1	Cluster shape is 1 x 8 x 1.
CUBLASLT_CLUSTER_SHAPE_8x1x1	Cluster shape is 8 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_2x8x1	Cluster shape is 2 x 8 x 1.
CUBLASLT_CLUSTER_SHAPE_8x2x1	Cluster shape is 8 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_1x16x1	Cluster shape is 1 x 16 x 1.
CUBLASLT_CLUSTER_SHAPE_16x1x1	Cluster shape is 16 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_1x3x1	Cluster shape is 1 x 3 x 1.
CUBLASLT_CLUSTER_SHAPE_1x5x1	Cluster shape is 1 x 5 x 1.
CUBLASLT_CLUSTER_SHAPE_1x6x1	Cluster shape is 1 x 6 x 1.
CUBLASLT_CLUSTER_SHAPE_1x7x1	Cluster shape is 1 x 7 x 1.
CUBLASLT_CLUSTER_SHAPE_1x9x1	Cluster shape is 1 x 9 x 1.
CUBLASLT_CLUSTER_SHAPE_1x10x1	Cluster shape is 1 x 10 x 1.
CUBLASLT_CLUSTER_SHAPE_1x11x1	Cluster shape is 1 x 11 x 1.
CUBLASLT_CLUSTER_SHAPE_1x12x1	Cluster shape is 1 x 12 x 1.
CUBLASLT_CLUSTER_SHAPE_1x13x1	Cluster shape is 1 x 13 x 1.
CUBLASLT_CLUSTER_SHAPE_1x14x1	Cluster shape is 1 x 14 x 1.
CUBLASLT_CLUSTER_SHAPE_1x15x1	Cluster shape is 1 x 15 x 1.
CUBLASLT_CLUSTER_SHAPE_2x3x1	Cluster shape is 2 x 3 x 1.
CUBLASLT_CLUSTER_SHAPE_2x5x1	Cluster shape is 2 x 5 x 1.
CUBLASLT_CLUSTER_SHAPE_2x6x1	Cluster shape is 2 x 6 x 1.
CUBLASLT_CLUSTER_SHAPE_2x7x1	Cluster shape is 2 x 7 x 1.
CUBLASLT_CLUSTER_SHAPE_3x1x1	Cluster shape is 3 x 1 x 1.

Value	Description
CUBLASLT_CLUSTER_SHAPE_3x2x1	Cluster shape is 3 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_3x3x1	Cluster shape is 3 x 3 x 1.
CUBLASLT_CLUSTER_SHAPE_3x4x1	Cluster shape is 3 x 4 x 1.
CUBLASLT_CLUSTER_SHAPE_3x5x1	Cluster shape is 3 x 5 x 1.
CUBLASLT_CLUSTER_SHAPE_4x3x1	Cluster shape is 4 x 3 x 1.
CUBLASLT_CLUSTER_SHAPE_5x1x1	Cluster shape is 5 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_5x2x1	Cluster shape is 5 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_5x3x1	Cluster shape is 5 x 3 x 1.
CUBLASLT_CLUSTER_SHAPE_6x1x1	Cluster shape is 6 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_6x2x1	Cluster shape is 6 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_7x1x1	Cluster shape is 7 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_7x2x1	Cluster shape is 7 x 2 x 1.
CUBLASLT_CLUSTER_SHAPE_9x1x1	Cluster shape is 9 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_10x1x1	Cluster shape is 10 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_11x1x1	Cluster shape is 11 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_12x1x1	Cluster shape is 12 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_13x1x1	Cluster shape is 13 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_14x1x1	Cluster shape is 14 x 1 x 1.
CUBLASLT_CLUSTER_SHAPE_15x1x1	Cluster shape is 15 x 1 x 1.

### 3.3.3. cublasLtEpilogue\_t

The `cublasLtEpilogue_t` is an enum type to set the postprocessing options for the epilogue.

Value	Description
CUBLASLT_EPILOGUE_DEFAULT = 1	No special postprocessing, just scale and quantize the results if necessary.
CUBLASLT_EPILOGUE_RELU = 2	Apply ReLU point-wise transform to the results: $(x := \max(x, 0))$
CUBLASLT_EPILOGUE_RELU_AUX = (CUBLASLT_EPILOGUE_RELU   128)	Apply ReLU point-wise transform to the results: $(x := \max(x, 0))$ . This epilogue mode produces an extra output, see <code>CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER</code> of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_BIAS = 4	Apply (broadcast) bias from the bias vector. Bias vector length must match matrix D rows, and it must be packed (such as stride between vector elements is 1). Bias vector is broadcast to all columns and added before applying the final postprocessing.
CUBLASLT_EPILOGUE_RELU_BIAS = (CUBLASLT_EPILOGUE_RELU   CUBLASLT_EPILOGUE_BIAS)	Apply bias and then ReLU transform.

Value	Description
CUBLASLT_EPILOGUE_RELU_AUX_BIAS = (CUBLASLT_EPILOGUE_RELU_AUX   CUBLASLT_EPILOGUE_BIAS)	Apply bias and then ReLU transform. This epilogue mode produces an extra output, see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_DRELU = 8   128	Apply ReLU gradient to matmul output. Store ReLU gradient in the output matrix. This epilogue mode requires an extra input, see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_DRELU_BGRAD = (CUBLASLT_EPILOGUE_DRELU   16)	Apply independently ReLU and Bias gradient to matmul output. Store ReLU gradient in the output matrix, and Bias gradient in the bias buffer (see CUBLASLT_MATMUL_DESC_BIAS_POINTER). This epilogue mode requires an extra input, see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_GELU = 32	Apply GELU point-wise transform to the results ( $x := \text{GELU}(x)$ ).
CUBLASLT_EPILOGUE_GELU_AUX = (CUBLASLT_EPILOGUE_GELU   128)	Apply GELU point-wise transform to the results ( $x := \text{GELU}(x)$ ). This epilogue mode outputs GELU input as a separate matrix (useful for training). See CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_GELU_BIAS = (CUBLASLT_EPILOGUE_GELU   CUBLASLT_EPILOGUE_BIAS)	Apply Bias and then GELU transform <sup>1</sup> .
CUBLASLT_EPILOGUE_GELU_AUX_BIAS = (CUBLASLT_EPILOGUE_GELU_AUX   CUBLASLT_EPILOGUE_BIAS)	Apply Bias and then GELU transform <sup>1</sup> . This epilogue mode outputs GELU input as a separate matrix (useful for training). See CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_DGELU = 64   128	Apply GELU gradient to matmul output. Store GELU gradient in the output matrix. This epilogue mode requires an extra input, see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_DGELU_BGRAD = (CUBLASLT_EPILOGUE_DGELU   16)	Apply independently GELU and Bias gradient to matmul output. Store GELU gradient in the output matrix, and Bias gradient in the bias buffer (see CUBLASLT_MATMUL_DESC_BIAS_POINTER). This epilogue mode requires an extra input, see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_BGRADA = 256	Apply Bias gradient to the input matrix A. The bias size corresponds to the number of rows of the matrix D. The reduction happens over the GEMM's "k" dimension. Store Bias gradient in the bias buffer,

Value	Description
	see CUBLASLT_MATMUL_DESC_BIAS_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .
CUBLASLT_EPILOGUE_BGRADB = 512	Apply Bias gradient to the input matrix B. The bias size corresponds to the number of columns of the matrix D. The reduction happens over the GEMM's "k" dimension. Store Bias gradient in the bias buffer, see CUBLASLT_MATMUL_DESC_BIAS_POINTER of <a href="#">cublasLtMatmulDescAttributes_t</a> .

**NOTES:**

1. GELU (Gaussian Error Linear Unit) is approximated by:  $0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$

### 3.3.4. [cublasLtHandle\\_t](#)

The `cublasLtHandle_t` type is a pointer type to an opaque structure holding the cuBLASLt library context. Use the following functions to manipulate this library context:

**[cublasLtCreate\(\)](#):**

To initialize the cuBLASLt library context and return a handle to an opaque structure holding the cuBLASLt library context.

**[cublasLtDestroy\(\)](#):**

To destroy a previously created cuBLASLt library context descriptor and release the resources.



**Note:** cuBLAS handle (`cublasHandle_t`) encapsulates a cuBLASLt handle. Any valid `cublasHandle_t` can be used in place of `cublasLtHandle_t` with a simple cast. However, unlike a cuBLAS handle, a cuBLASLt handle is not tied to any particular CUDA context.

### 3.3.5. [cublasLtLoggerCallback\\_t](#)

`cublasLtLoggerCallback_t` is a callback function pointer type.

**Parameters:**

Parameter	Memory	Input / Output	Description
logLevel		Output	Refer to <a href="#">cuBLASLt Logging</a> .
functionName		Output	The name of the API that logged this message.
message		Output	The log message.

Use the following function to set the callback function:

**[cublasLtLoggerSetCallback\(\)](#)**

### 3.3.6. cublasLtMatmulAlgo\_t

`cublasLtMatmulAlgo_t` is an opaque structure holding the description of the matrix multiplication algorithm. This structure can be trivially serialized and later restored for use with the same version of cuBLAS library to save on selecting the right configuration again.

### 3.3.7. cublasLtMatmulAlgoCapAttributes\_t

`cublasLtMatmulAlgoCapAttributes_t` enumerates matrix multiplication algorithm capability attributes that can be retrieved from an initialized [cublasLtMatmulAlgo\\_t](#) descriptor.

Value	Description	Data Type
CUBLASLT_ALGO_CAP_SPLITK_SUPPORT	Support for split-K. Boolean (0 or 1) to express if split-K implementation is supported. 0 means no support, and supported otherwise. See CUBLASLT_ALGO_CONFIG_SPLITK_NUM of <a href="#">cublasLtMatmulAlgoConfigAttributes_t</a> .	int32_t
CUBLASLT_ALGO_CAP_REDUCTION_SCHEME_MASK	Mask to express the types of reduction schemes supported, see <a href="#">cublasLtReductionScheme_t</a> . If the reduction scheme is not masked out then it is supported. For example: <code>int isReductionSchemeComputeTypeSupported ? (reductionSchemeMask &amp; CUBLASLT_REDUCTION_SCHEME_COMPUTE_TYPE) == CUBLASLT_REDUCTION_SCHEME_COMPUTE_TYPE ? 1 : 0;</code>	uint32_t
CUBLASLT_ALGO_CAP_CTA_SWIZZLING_SUPPORT	Support for CTA-swizzling. Boolean (0 or 1) to express if CTA-swizzling implementation is supported. 0 means no support, and 1 means supported value of 1; other values are reserved. Also refer to CUBLASLT_ALGO_CONFIG_CTA_SWIZZLING of <a href="#">cublasLtMatmulAlgoConfigAttributes_t</a> .	uint32_t
CUBLASLT_ALGO_CAP_STRIDED_BATCH_SUPPORT	Support strided batch. 0 means no support, supported otherwise.	int32_t
CUBLASLT_ALGO_CAP_OUT_OF_PLACE_RESULTS	Supports results out of place ( $D \neq C$ in $D = \alpha.A.B + \beta.C$ ).	int32_t

Value	Description	Data Type
	0 means no support, supported otherwise.	
CUBLASLT_ALGO_CAP_UPLO_SUPPORT	syrk (symmetric rank k update)/ herk (Hermitian rank k update) support (on top of regular gemm). 0 means no support, supported otherwise.	int32_t
CUBLASLT_ALGO_CAP_TILE_IDS	The tile IDs possible to use. Refer to <a href="#">cublasLtMatmulTile_t</a> . If no tile IDs are supported then use CUBLASLT_MATMUL_TILE_UNDEFINED. Use <a href="#">cublasLtMatmulAlgoCapGetAttribute()</a> with sizeInBytes=0 to query the actual count.	Array of uint32_t
CUBLASLT_ALGO_CAP_STAGES_IDS	The stages IDs possible to use. Refer to <a href="#">cublasLtMatmulStages_t</a> . If no stages IDs are supported, then use CUBLASLT_MATMUL_STAGES_UNDEFINED. Use <a href="#">cublasLtMatmulAlgoCapGetAttribute()</a> with sizeInBytes=0 to query the actual count.	Array of uint32_t
CUBLASLT_ALGO_CAP_CUSTOM_OPTION_MAX	Custom option range is from 0 to CUBLASLT_ALGO_CAP_CUSTOM_OPTION_MAX (inclusive). Refer to CUBLASLT_ALGO_CONFIG_CUSTOM_OPTION of <a href="#">cublasLtMatmulAlgoConfigAttributes_t</a> .	int32_t
CUBLASLT_ALGO_CAP_MATHMODE_IMPL	Indicates whether the algorithm is using regular compute or tensor operations. 0 means regular compute, 1 means tensor operations.  DEPRECATED	int32_t
CUBLASLT_ALGO_CAP_GAUSSIAN_IMPL	Indicates whether the algorithm implements the Gaussian optimization of complex matrix multiplication. 0 means regular compute; 1 means Gaussian. See <a href="#">cublasMath_t</a> .  DEPRECATED	int32_t
CUBLASLT_ALGO_CAP_CUSTOM_MEMORY_ORDER	Indicates whether the algorithm supports custom (not COL or ROW memory order). 0 means	int32_t



Value	Description	Data Type
	only COL and ROW memory order is allowed, non-zero means that algo might have different requirements. Refer to <a href="#">cublasLtOrder_t</a> .	
CUBLASLT_ALGO_CAP_POINTER_MODE_MASK	Bitmask enumerating the pointer modes the algorithm supports. Refer to <a href="#">cublasLtPointerModeMask_t</a> .	uint32_t
CUBLASLT_ALGO_CAP_EPILOGUE_MASK	Bitmask enumerating the kinds of postprocessing algorithm supported in the epilogue. Refer to <a href="#">cublasLtEpilogue_t</a> .	uint32_t
CUBLASLT_ALGO_CAP_LD_NEGATIVE	Support for negative leading dimension for all of the matrices. 0 means no support, supported otherwise.	uint32_t
CUBLASLT_ALGO_CAP_NUMERICAL_IMPL_FLAGS	Details about algorithm's implementation that affect it's numerical behavior. Refer to <a href="#">cublasLtNumericalImplFlags_t</a> .	uint64_t
CUBLASLT_ALGO_CAP_MIN_ALIGNMENT_A_BYTES	Minimum alignment required for A matrix in bytes.	uint32_t
CUBLASLT_ALGO_CAP_MIN_ALIGNMENT_B_BYTES	Minimum alignment required for B matrix in bytes.	uint32_t
CUBLASLT_ALGO_CAP_MIN_ALIGNMENT_C_BYTES	Minimum alignment required for C matrix in bytes.	uint32_t
CUBLASLT_ALGO_CAP_MIN_ALIGNMENT_D_BYTES	Minimum alignment required for D matrix in bytes.	uint32_t

Use the following function to manipulate this descriptor:

[cublasLtMatmulAlgoCapGetAttribute\(\)](#): To retrieve the capability attribute(s) of the descriptor.

### 3.3.8. cublasLtMatmulAlgoConfigAttributes\_t

`cublasLtMatmulAlgoConfigAttributes_t` is an enumerated type that contains the configuration attributes for the matrix multiply algorithms. These configuration attributes are algorithm-specific, and can be set. The attributes configuration of a given algorithm should lie within the boundaries expressed by its capability attributes.

Value	Description	Data Type
CUBLASLT_ALGO_CONFIG_ID	Read-only attribute. Algorithm index. Refer to <a href="#">cublasLtMatmulAlgoGetIds()</a> . Set by <a href="#">cublasLtMatmulAlgoInit()</a> .	int32_t

Value	Description	Data Type
CUBLASLT_ALGO_CONFIG_TILE_ID	Tile ID. See <a href="#">cublasLtMatmulTile_t</a> . Default: CUBLASLT_MATMUL_TILE_UNDEFINED.	uint32_t
CUBLASLT_ALGO_CONFIG_STAGES_ID	Stages ID. Refer to <a href="#">cublasLtMatmulStages_t</a> . Default: CUBLASLT_MATMUL_STAGES_UNDEFINED.	uint32_t
CUBLASLT_ALGO_CONFIG_SPLITK_NUM	Number of K splits. If the number of K splits is greater than one, SPLITK_NUM parts of matrix multiplication will be computed in parallel. The results will be accumulated according to CUBLASLT_ALGO_CONFIG_REDUCTION_SCHEME.	int32_t
CUBLASLT_ALGO_CONFIG_REDUCTION_SCHEME	Reduction scheme to use when splitK value > 1. Default: CUBLASLT_REDUCTION_SCHEME_NONE. See <a href="#">cublasLtReductionScheme_t</a> .	uint32_t
CUBLASLT_ALGO_CONFIG_CTA_SWIZZLING	Enable/Disable CTA swizzling. Change mapping from CUDA grid coordinates to parts of the matrices. Possible values: 0 and 1; other values reserved.	uint32_t
CUBLASLT_ALGO_CONFIG_CUSTOM_OPTION	Custom option value. Each algorithm can support some custom options that don't fit the description of the other configuration attributes. Refer to the CUBLASLT_ALGO_CAP_CUSTOM_OPTION_MAX of <a href="#">cublasLtMatmulAlgoCapAttributes_t</a> for the accepted range for a specific case.	uint32_t
CUBLASLT_ALGO_CONFIG_INNER_SHAPE_ID	Inner shape ID. Refer to <a href="#">cublasLtMatmulInnerShape_t</a> . Default: CUBLASLT_MATMUL_INNER_SHAPE_UNDEFINED.	uint16_t
CUBLASLT_ALGO_CONFIG_CLUSTER_SHAPE_ID	Cluster shape ID. Refer to <a href="#">cublasLtClusterShape_t</a> . Default: CUBLASLT_CLUSTER_SHAPE_AUTO.	uint16_t

Use the below function to manipulate this descriptor:

[cublasLtMatmulAlgoConfigSetAttribute\(\)](#): To retrieve the attribute(s) of the descriptor.

[cublasLtMatmulAlgoConfigGetAttribute\(\)](#): To query a previously created descriptor for the attribute(s).

### 3.3.9. `cublasLtMatmulDesc_t`

The `cublasLtMatmulDesc_t` is a pointer to an opaque structure holding the description of the matrix multiplication operation [cublasLtMatmul\(\)](#). Use the following functions to manipulate this descriptor:

**[cublasLtMatmulDescCreate\(\)](#):**

To create one instance of the descriptor.

**[cublasLtMatmulDescDestroy\(\)](#):**

To destroy a previously created descriptor and release the resources.

### 3.3.10. `cublasLtMatmulDescAttributes_t`

`cublasLtMatmulDescAttributes_t` is a descriptor structure containing the attributes that define the specifics of the matrix multiply operation.

Attribute Name	Description	Data Type
<code>CUBLASLT_MATMUL_DESC_COMPUTE_TYPE</code>	Compute type. Defines data type used for multiply and accumulate operations, and the accumulator during the matrix multiplication. Refer to <a href="#">cudaDataType_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATMUL_DESC_SCALE_TYPE</code>	Scale type. Defines the data type of the scaling factors alpha and beta. The accumulator value and the value from matrix C are typically converted to scale type before final scaling. Value is then converted from scale type to the type of matrix D before storing in memory. Default value same as <code>CUBLASLT_MATMUL_DESC_COMPUTE_TYPE</code> . See <a href="#">cudaDataType_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATMUL_DESC_POINTER_MODE</code>	Specifies alpha and beta are passed by reference, whether they are scalars on the host or on the device, or device vectors. Default value is: <code>CUBLASLT_POINTER_MODE_HOST</code> (for example, on the host). See <a href="#">cublasLtPointerMode_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATMUL_DESC_TRANSA</code>	Specifies the type of transformation operation that should be performed on matrix A. Default value is <code>CUBLAS_OP_N</code> (for example, non-transpose operation). See <a href="#">cublasOperation_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATMUL_DESC_TRANSB</code>	Specifies the type of transformation operation that should be performed on matrix B. Default value is <code>CUBLAS_OP_N</code> (i.e., non-transpose operation). Refer to <a href="#">cublasOperation_t</a> .	<code>int32_t</code>

Attribute Name	Description	Data Type
CUBLASLT_MATMUL_DESC_TRANS	Specifies the type of transformation operation that should be performed on matrix C. Currently only CUBLAS_OP_N is supported. Default value is CUBLAS_OP_N (for example, non-transpose operation). Refer to <a href="#">cublasOperation_t</a> .	int32_t
CUBLASLT_MATMUL_DESC_FILL_MODE	Indicates whether the lower or upper part of the dense matrix was filled, and consequently should be used by the function. Default value is: CUBLAS_FILL_MODE_FULL. See <a href="#">cublasFillMode_t</a> .	int32_t
CUBLASLT_MATMUL_DESC_EPILOGUE	Epilogue function. See <a href="#">cublasLtEpilogue_t</a> . Default value is: CUBLASLT_EPILOGUE_DEFAULT.	uint32_t
CUBLASLT_MATMUL_DESC_BIAS_POINTER	<p>Bias or Bias gradient vector pointer in the device memory.</p> <ul style="list-style-type: none"> <li>▶ Input vector with length that matches the number of rows of matrix D when one of the following epilogues is used: CUBLASLT_EPILOGUE_BIAS, CUBLASLT_EPILOGUE_RELU_BIAS, CUBLASLT_EPILOGUE_RELU_AUX_BIAS, CUBLASLT_EPILOGUE_GELU_BIAS, CUBLASLT_EPILOGUE_GELU_AUX_BIAS.</li> <li>▶ Output vector with length that matches the number of rows of matrix D when one of the following epilogues is used: CUBLASLT_EPILOGUE_DRELU_BGRAD, CUBLASLT_EPILOGUE_DGELU_BGRAD, CUBLASLT_EPILOGUE_BGRADA.</li> <li>▶ Output vector with length that matches the number of columns of matrix D when one of the following epilogues is used: CUBLASLT_EPILOGUE_BGRADB.</li> </ul> <p>Bias vector elements are the same type as alpha and beta (see CUBLASLT_MATMUL_DESC_SCALE_TYPE in this table) when matrix D datatype is CUDA_R_8I and same as matrix D datatype otherwise. See the datatypes table under <code>cublasLtMatmul()</code> for detailed mapping. Default value is: NULL.</p>	void */ const void *
CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER	Pointer for epilogue auxiliary buffer.	void */ const void *

Attribute Name	Description	Data Type
	<ul style="list-style-type: none"> <li>▶ Output vector for ReLu bit-mask in forward pass when CUBLASLT_EPILOGUE_RELU_AUX or CUBLASLT_EPILOGUE_RELU_AUX_BIAS are used.</li> <li>▶ Input vector for ReLu bit-mask in backward pass when CUBLASLT_EPILOGUE_DRELU or CUBLASLT_EPILOGUE_DRELU_BGRAD are used.</li> <li>▶ Output of GELU input matrix in forward pass when CUBLASLT_EPILOGUE_GELU_AUX_BIAS is used.</li> <li>▶ Input of GELU input matrix for backward pass when CUBLASLT_EPILOGUE_DGELU or CUBLASLT_EPILOGUE_DGELU_BGRAD are used.</li> </ul> <p>For aux data type see CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_DATA_TYPE. Routines that don't dereference this pointer, like cublasLtMatmulAlgoGetHeuristic() depend on its value to determine expected pointer alignment. Requires setting the CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_LD attribute.</p>	
CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_LD	<p>Leading dimension for epilogue auxiliary buffer.</p> <ul style="list-style-type: none"> <li>▶ ReLu bit-mask matrix leading dimension in elements (such as bits) when CUBLASLT_EPILOGUE_RELU_AUX, CUBLASLT_EPILOGUE_RELU_AUX_BIAS, CUBLASLT_EPILOGUE_DRELU or CUBLASLT_EPILOGUE_DRELU_BGRAD epilogue is used. Must be divisible by 128 and be no less than the number of rows in the output matrix.</li> <li>▶ GELU input matrix leading dimension in elements when CUBLASLT_EPILOGUE_GELU_AUX_BIAS, CUBLASLT_EPILOGUE_DGELU or CUBLASLT_EPILOGUE_DGELU_BGRAD epilogue is used. Must be divisible by 8 and be no less than the number of rows in the output matrix.</li> </ul>	int64_t

Attribute Name	Description	Data Type
CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_BATCH_STRIDE	<p>Batch stride for epilogue auxiliary buffer.</p> <ul style="list-style-type: none"> <li>▶ ReLu bit-mask matrix batch stride in elements (such as bits) when CUBLASLT_EPILOGUE_RELU_AUX, CUBLASLT_EPILOGUE_RELU_AUX_BIAS, CUBLASLT_EPILOGUE_DRELU or CUBLASLT_EPILOGUE_DRELU_BGRAD epilogue is used. Must be divisible by 128.</li> <li>▶ GELU input matrix batch stride in elements when CUBLASLT_EPILOGUE_GELU_AUX_BIAS, CUBLASLT_EPILOGUE_DGELU or CUBLASLT_EPILOGUE_DGELU_BGRAD epilogue is used. Must be divisible by 8.</li> </ul> <p>Default value: 0.</p>	int64_t
CUBLASLT_MATMUL_DESC_ALPHA_VECTOR_BATCH_STRIDE	<p>Batch stride for alpha vector. Used together with CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_HOST when matrix D's CUBLASLT_MATRIX_LAYOUT_BATCH_COUNT is greater than 1. If CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_ZERO is set then CUBLASLT_MATMUL_DESC_ALPHA_VECTOR_BATCH_STRIDE must be set to 0 as this mode doesn't support batched alpha vector. Default value: 0.</p>	int64_t
CUBLASLT_MATMUL_DESC_SM_COUNT_TARGET	<p>Number of SMs to target for parallel execution. Optimizes heuristics for execution on a different number of SMs when user expects a concurrent stream to be using some of the device resources. Default value: 0.</p>	int32_t
CUBLASLT_MATMUL_DESC_A_SCALE_POINTER	<p>Device pointer to the scale factor value that converts data in matrix A to the compute data type range. The scaling factor must have the same type as the compute type. If not specified, or set to NULL, the scaling factor is assumed to be 1. If set for an unsupported matrix data, scale, and compute type combination, calling <code>cublasLtMatmul()</code> will return CUBLAS_INVALID_VALUE. Default value: NULL</p>	const void*

Attribute Name	Description	Data Type
CUBLASLT_MATMUL_DESC_B_SCALE_POINTER	Equivalent to CUBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix B. Default value: NULL	const void*
CUBLASLT_MATMUL_DESC_C_SCALE_POINTER	Equivalent to CUBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix C. Default value: NULL	const void*
CUBLASLT_MATMUL_DESC_D_SCALE_POINTER	Equivalent to CUBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix D. Default value: NULL	const void*
CUBLASLT_MATMUL_DESC_AMAX_D_POINTER	Device pointer to the memory location that on completion will be set to the maximum of absolute values in the output matrix. The computed value has the same type as the compute type. If not specified, or set to NULL, the maximum absolute value is not computed. If set for an unsupported matrix data, scale, and compute type combination, calling <code>cublasLtMatmul()</code> will return <code>CUBLAS_INVALID_VALUE</code> . Default value: NULL	void *
CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_DATA_TYPE	<p>The type of the data that will be stored in CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER. If unset (or set to the default value of -1), the data type is set to be the output matrix element data type (DType) with some exceptions:</p> <ul style="list-style-type: none"> <li>► ReLu uses a bit-mask.</li> <li>► For FP8 kernels with an output type (DType) of <code>CUDA_R_8F_E4M3</code>, the data type can be set to a non-default value if: <ol style="list-style-type: none"> <li>1. AType and BType are <code>CUDA_R_8F_E4M3</code>.</li> <li>2. Bias Type is <code>CUDA_R_16F</code>.</li> <li>3. CType is <code>CUDA_R_16BF</code> or <code>CUDA_R_16F</code></li> <li>4. CUBLASLT_MATMUL_DESC_EPILOGUE is set to <code>CUBLASLT_EPILOGUE_GELU_AUX</code></li> </ol> </li> </ul> <p>When CType is <code>CUDA_R_16BF</code>, the data type may be set to <code>CUDA_R_16BF</code> or <code>CUDA_R_8F_E4M3</code>. When CType is <code>CUDA_R_16F</code>, the data type may be set to <code>CUDA_R_16F</code>. Otherwise, the data type should be left unset or set to the default value of -1.</p>	<p><code>int32_t</code> based on <code>cudaDataType</code></p>

Attribute Name	Description	Data Type
	If set for an unsupported matrix data, scale, and compute type combination, calling <code>cublasLtMatmul()</code> will return <code>CUBLAS_INVALID_VALUE</code> . Default value: -1	
<code>CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_SCALING_FACTOR_POINTER</code>	Device pointer to the scaling factor value to convert results from compute type data range to storage data range in the auxiliary matrix that is set via <code>CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER</code> . The scaling factor value must have the same type as the compute type. If not specified, or set to <code>NULL</code> , the scaling factor is assumed to be 1. If set for an unsupported matrix data, scale, and compute type combination, calling <code>cublasLtMatmul()</code> will return <code>CUBLAS_INVALID_VALUE</code> . Default value: <code>NULL</code>	<code>void *</code>
<code>CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_MAX_ABSOLUTE_POINTER</code>	Device pointer to the memory location that on completion will be set to the maximum of absolute values in the buffer that is set via <code>CUBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER</code> . The computed value has the same type as the compute type. If not specified, or set to <code>NULL</code> , the maximum absolute value is not computed. If set for an unsupported matrix data, scale, and compute type combination, calling <code>cublasLtMatmul()</code> will return <code>CUBLAS_INVALID_VALUE</code> . Default value: <code>NULL</code>	<code>void *</code>
<code>CUBLASLT_MATMUL_DESC_FAST_ACCUM</code>	Flag for managing FP8 fast accumulation mode. When enabled, problem execution might be faster but at the cost of lower accuracy because intermediate results will not periodically be promoted to a higher precision. Default value: 0 - fast accumulation mode is disabled	<code>int8_t</code>
<code>CUBLASLT_MATMUL_DESC_BIAS_DATA_TYPE</code>	Type of the bias or bias gradient vector in the device memory. Bias case: see <code>CUBLASLT_EPILOGUE_BIAS</code> . If unset (or set to the default value of -1), the bias vector elements are the same type as the elements of the output matrix (Dtype) with the following exceptions: <ul style="list-style-type: none"> <li>IMMA kernels with <code>computeType=CUDA_R_32I</code> and <code>Ctype=CUDA_R_8I</code> where the bias vector elements are</li> </ul>	<code>int32_t</code> based on <code>cudaDataType</code>



Attribute Name	Description	Data Type
	<p>the same type as alpha, beta (CUBLASLT_MATMUL_DESC_SCALE_TYPE=CUDA_R_32F)</p> <ul style="list-style-type: none"> <li>► For FP8 kernels with an output type of CUDA_R_32F, CUDA_R_8F_E4M3 or CUDA_R_8F_E5M2. See <a href="#">cublasLtMatmul()</a> for more details.</li> </ul> <p>Default value: -1</p>	

Use the following functions to manipulate this descriptor:

[cublasLtMatmulDescSetAttribute\(\)](#): To initialize the attribute(s) of the descriptor.

[cublasLtMatmulDescGetAttribute\(\)](#): To query a previously created descriptor for the attribute(s).

### 3.3.11. cublasLtMatmulHeuristicResult\_t

`cublasLtMatmulHeuristicResult_t` is a descriptor that holds the configured matrix multiplication algorithm descriptor and its runtime properties.

Member	Description
<code>cublasLtMatmulAlgo_t algo</code>	Must be initialized with <a href="#">cublasLtMatmulAlgoInit()</a> if the preference CUBLASLT_MATMUL_PERF_SEARCH_MODE is set to CUBLASLT_SEARCH_LIMITED_BY_ALGO_ID. Refer to <a href="#">cublasLtMatmulSearch_t</a> .
<code>size_t workspaceSize;</code>	Actual size of workspace memory required.
<code>cublasStatus_t state;</code>	Result status. Other fields are valid only if, after call to <a href="#">cublasLtMatmulAlgoGetHeuristic()</a> , this member is set to CUBLAS_STATUS_SUCCESS.
<code>float wavesCount;</code>	Waves count is a device utilization metric. A <code>wavesCount</code> value of 1.0f suggests that when the kernel is launched it will fully occupy the GPU.
<code>int reserved[4];</code>	Reserved.

### 3.3.12. cublasLtMatmulInnerShape\_t

`cublasLtMatmulInnerShape_t` is an enumerated type used to configure various aspects of the internal kernel design. This does not impact the CUDA grid size.

Value	Description
CUBLASLT_MATMUL_INNER_SHAPE_UNDEFINED	Inner shape is undefined.
CUBLASLT_MATMUL_INNER_SHAPE_MMA884	Inner shape is MMA884.
CUBLASLT_MATMUL_INNER_SHAPE_MMA1684	Inner shape is MMA1684.
CUBLASLT_MATMUL_INNER_SHAPE_MMA1688	Inner shape is MMA1688.
CUBLASLT_MATMUL_INNER_SHAPE_MMA16816	Inner shape is MMA16816.

### 3.3.13. `cublasLtMatmulPreference_t`

The `cublasLtMatmulPreference_t` is a pointer to an opaque structure holding the description of the preferences for `cublasLtMatmulAlgoGetHeuristic()` configuration. Use the following functions to manipulate this descriptor:

**`cublasLtMatmulPreferenceCreate()`:**

To create one instance of the descriptor.

**`cublasLtMatmulPreferenceDestroy()`:**

To destroy a previously created descriptor and release the resources.

### 3.3.14. `cublasLtMatmulPreferenceAttributes_t`

`cublasLtMatmulPreferenceAttributes_t` is an enumerated type used to apply algorithm search preferences while fine-tuning the heuristic function.

Value	Description	Data Type
<code>CUBLASLT_MATMUL_PREF_SEARCH_MODE</code>	Search mode. Refer to <code>cublasLtMatmulSearch_t</code> . Default is <code>CUBLASLT_SEARCH_BEST_FIT</code> .	<code>uint32_t</code>
<code>CUBLASLT_MATMUL_PREF_MAX_WORKSPACE_BYTES</code>	Maximum allowed workspace memory. Default is 0 (no workspace memory allowed).	<code>uint64_t</code>
<code>CUBLASLT_MATMUL_PREF_MATH_MODE_MASK</code>	Math mode mask. Refer to <code>cublasMath_t</code> . Only algorithms with <code>CUBLASLT_ALGO_CAP_MATHMODE_IMPL</code> that is not masked out by this attribute are allowed. Default is 1 (allows both default and tensor op math).  DEPRECATED; will be removed in a future release. Refer to <code>cublasLtNumericalImplFlags_t</code> for replacement.	<code>uint32_t</code>
<code>CUBLASLT_MATMUL_PREF_REDUCTION_SCHEME_MASK</code>	Reduction scheme mask. Refer to <code>cublasLtReductionScheme_t</code> . Only algorithm configurations specifying <code>CUBLASLT_ALGO_CONFIG_REDUCTION_SCHEME</code> that is not masked out by this attribute are allowed. For example, a mask value of 0x03 will allow only <code>INPLACE</code> and <code>COMPUTE_TYPE</code> reduction schemes. Default is <code>CUBLASLT_REDUCTION_SCHEME_MASK</code> (i.e., allows all reduction schemes).	<code>uint32_t</code>
<code>CUBLASLT_MATMUL_PREF_GAUSSIAN_MODE_MASK</code>	Gaussian mode mask. Refer to <code>cublasLt3mMode_t</code> . Only algorithms with <code>CUBLASLT_ALGO_CAP_GAUSSIAN_IMPL</code> that is not masked out by this attribute are allowed. Default is	<code>uint32_t</code>

Value	Description	Data Type
	CUBLASLT_3M_MODE_ALLOWED (i.e., allows both Gaussian and regular math).  DEPRECATED; will be removed in a future release, see <code>cublasLtNumericalImplFlags_t</code> for replacement.	
CUBLASLT_MATMUL_PREF_MIN_ALIGNMENT_A	Minimum buffer alignment for matrix A (in bytes). Selecting a smaller value will exclude algorithms that can not work with matrix A, which is not as strictly aligned as the algorithms need. Default is 256 bytes.	uint32_t
CUBLASLT_MATMUL_PREF_MIN_ALIGNMENT_B	Minimum buffer alignment for matrix B (in bytes). Selecting a smaller value will exclude algorithms that can not work with matrix B, which is not as strictly aligned as the algorithms need. Default is 256 bytes.	uint32_t
CUBLASLT_MATMUL_PREF_MIN_ALIGNMENT_C	Minimum buffer alignment for matrix C (in bytes). Selecting a smaller value will exclude algorithms that can not work with matrix C, which is not as strictly aligned as the algorithms need. Default is 256 bytes.	uint32_t
CUBLASLT_MATMUL_PREF_MIN_ALIGNMENT_D	Minimum buffer alignment for matrix D (in bytes). Selecting a smaller value will exclude algorithms that can not work with matrix D, which is not as strictly aligned as the algorithms need. Default is 256 bytes.	uint32_t
CUBLASLT_MATMUL_PREF_MAX_WAVES_COUNT	Maximum wave count. Refer to <code>cublasLtMatmulHeuristicResult_t::wavesCount</code> . Selecting a non-zero value will exclude algorithms that report device utilization higher than specified. Default is 0.0f.	float
CUBLASLT_MATMUL_PREF_POINTER_MODE_MASK	Pointer mode mask. Refer to <code>cublasLtPointerModeMask_t</code> . Filters the heuristic result to include only algorithms that support all required modes. Default is (CUBLASLT_POINTER_MODE_MASK_HOST   CUBLASLT_POINTER_MODE_MASK_DEVICE) (only allows algorithms that support both regular host and device pointers).	uint32_t
CUBLASLT_MATMUL_PREF_EPILOGUE_MASK	Epilogue selector mask. Refer to <code>cublasLtEpilogue_t</code> . Filters the heuristic result to include only algorithms that support all required operations. Default is CUBLASLT_EPILOGUE_DEFAULT (only allows algorithms that support default epilogue)	uint32_t
CUBLASLT_MATMUL_PREF_IMPL_MASK	Numerical implementation details mask. Refer to	uint64_t

Value	Description	Data Type
	<a href="#"><code>cublasLtNumericalImplFlags_t</code></a> . Filters heuristic result to only include algorithms that use the allowed implementations. Default is <code>uint64_t(-1)</code> (allow everything).	

Use the following functions to manipulate this descriptor:

[`cublasLtMatmulPreferenceSetAttribute\(\)`](#): To set the attribute(s) of the descriptor.

[`cublasLtMatmulPreferenceGetAttribute\(\)`](#): To query a previously created descriptor for the attribute(s).

### 3.3.15. `cublasLtMatmulSearch_t`

`cublasLtMatmulSearch_t` is an enumerated type that contains the attributes for heuristics search type.

Value	Description	Data Type
<code>CUBLASLT_SEARCH_BEST_FIT</code>	Request heuristics for the best algorithm for the given use case.	
<code>CUBLASLT_SEARCH_LIMITED_BY_ALGO_ID</code>	Request heuristics only for the pre-configured algo ID.	

### 3.3.16. `cublasLtMatmulTile_t`

`cublasLtMatmulTile_t` is an enumerated type used to set the tile size in rows x columns. See also [CUTLASS: Fast Linear Algebra in CUDA C++](#).

Value	Description
<code>CUBLASLT_MATMUL_TILE_UNDEFINED</code>	Tile size is undefined.
<code>CUBLASLT_MATMUL_TILE_8x8</code>	Tile size is 8 rows x 8 columns.
<code>CUBLASLT_MATMUL_TILE_8x16</code>	Tile size is 8 rows x 16 columns.
<code>CUBLASLT_MATMUL_TILE_16x8</code>	Tile size is 16 rows x 8 columns.
<code>CUBLASLT_MATMUL_TILE_8x32</code>	Tile size is 8 rows x 32 columns.
<code>CUBLASLT_MATMUL_TILE_16x16</code>	Tile size is 16 rows x 16 columns.
<code>CUBLASLT_MATMUL_TILE_32x8</code>	Tile size is 32 rows x 8 columns.
<code>CUBLASLT_MATMUL_TILE_8x64</code>	Tile size is 8 rows x 64 columns.
<code>CUBLASLT_MATMUL_TILE_16x32</code>	Tile size is 16 rows x 32 columns.
<code>CUBLASLT_MATMUL_TILE_32x16</code>	Tile size is 32 rows x 16 columns.
<code>CUBLASLT_MATMUL_TILE_64x8</code>	Tile size is 64 rows x 8 columns.
<code>CUBLASLT_MATMUL_TILE_32x32</code>	Tile size is 32 rows x 32 columns.
<code>CUBLASLT_MATMUL_TILE_32x64</code>	Tile size is 32 rows x 64 columns.
<code>CUBLASLT_MATMUL_TILE_64x32</code>	Tile size is 64 rows x 32 columns.
<code>CUBLASLT_MATMUL_TILE_32x128</code>	Tile size is 32 rows x 128 columns.

Value	Description
CUBLASLT_MATMUL_TILE_64x64	Tile size is 64 rows x 64 columns.
CUBLASLT_MATMUL_TILE_128x32	Tile size is 128 rows x 32 columns.
CUBLASLT_MATMUL_TILE_64x128	Tile size is 64 rows x 128 columns.
CUBLASLT_MATMUL_TILE_128x64	Tile size is 128 rows x 64 columns.
CUBLASLT_MATMUL_TILE_64x256	Tile size is 64 rows x 256 columns.
CUBLASLT_MATMUL_TILE_128x128	Tile size is 128 rows x 128 columns.
CUBLASLT_MATMUL_TILE_256x64	Tile size is 256 rows x 64 columns.
CUBLASLT_MATMUL_TILE_64x512	Tile size is 64 rows x 512 columns.
CUBLASLT_MATMUL_TILE_128x256	Tile size is 128 rows x 256 columns.
CUBLASLT_MATMUL_TILE_256x128	Tile size is 256 rows x 128 columns.
CUBLASLT_MATMUL_TILE_512x64	Tile size is 512 rows x 64 columns.
CUBLASLT_MATMUL_TILE_64x96	Tile size is 64 rows x 96 columns.
CUBLASLT_MATMUL_TILE_96x64	Tile size is 96 rows x 64 columns.
CUBLASLT_MATMUL_TILE_96x128	Tile size is 96 rows x 128 columns.
CUBLASLT_MATMUL_TILE_128x160	Tile size is 128 rows x 160 columns.
CUBLASLT_MATMUL_TILE_160x128	Tile size is 160 rows x 128 columns.
CUBLASLT_MATMUL_TILE_192x128	Tile size is 192 rows x 128 columns.
CUBLASLT_MATMUL_TILE_128x192	Tile size is 128 rows x 192 columns.
CUBLASLT_MATMUL_TILE_128x96	Tile size is 128 rows x 96 columns.

### 3.3.17. cublasLtMatmulStages\_t

`cublasLtMatmulStages_t` is an enumerated type used to configure the size and number of shared memory buffers where input elements are staged. Number of staging buffers defines kernel's pipeline depth.

Value	Description
CUBLASLT_MATMUL_STAGES_UNDEFINED	Stage size is undefined.
CUBLASLT_MATMUL_STAGES_16x1	Stage size is 16, number of stages is 1.
CUBLASLT_MATMUL_STAGES_16x2	Stage size is 16, number of stages is 2.
CUBLASLT_MATMUL_STAGES_16x3	Stage size is 16, number of stages is 3.
CUBLASLT_MATMUL_STAGES_16x4	Stage size is 16, number of stages is 4.
CUBLASLT_MATMUL_STAGES_16x5	Stage size is 16, number of stages is 5.
CUBLASLT_MATMUL_STAGES_16x6	Stage size is 16, number of stages is 6.
CUBLASLT_MATMUL_STAGES_32x1	Stage size is 32, number of stages is 1.
CUBLASLT_MATMUL_STAGES_32x2	Stage size is 32, number of stages is 2.
CUBLASLT_MATMUL_STAGES_32x3	Stage size is 32, number of stages is 3.
CUBLASLT_MATMUL_STAGES_32x4	Stage size is 32, number of stages is 4.
CUBLASLT_MATMUL_STAGES_32x5	Stage size is 32, number of stages is 5.

Value	Description
CUBLASLT_MATMUL_STAGES_32x6	Stage size is 32, number of stages is 6.
CUBLASLT_MATMUL_STAGES_64x1	Stage size is 64, number of stages is 1.
CUBLASLT_MATMUL_STAGES_64x2	Stage size is 64, number of stages is 2.
CUBLASLT_MATMUL_STAGES_64x3	Stage size is 64, number of stages is 3.
CUBLASLT_MATMUL_STAGES_64x4	Stage size is 64, number of stages is 4.
CUBLASLT_MATMUL_STAGES_64x5	Stage size is 64, number of stages is 5.
CUBLASLT_MATMUL_STAGES_64x6	Stage size is 64, number of stages is 6.
CUBLASLT_MATMUL_STAGES_128x1	Stage size is 128, number of stages is 1.
CUBLASLT_MATMUL_STAGES_128x2	Stage size is 128, number of stages is 2.
CUBLASLT_MATMUL_STAGES_128x3	Stage size is 128, number of stages is 3.
CUBLASLT_MATMUL_STAGES_128x4	Stage size is 128, number of stages is 4.
CUBLASLT_MATMUL_STAGES_128x5	Stage size is 128, number of stages is 5.
CUBLASLT_MATMUL_STAGES_128x6	Stage size is 128, number of stages is 6.
CUBLASLT_MATMUL_STAGES_32x10	Stage size is 32, number of stages is 10.
CUBLASLT_MATMUL_STAGES_8x4	Stage size is 8, number of stages is 4.
CUBLASLT_MATMUL_STAGES_16x10	Stage size is 16, number of stages is 10.
CUBLASLT_MATMUL_STAGES_8x5	Stage size is 8, number of stages is 5.
CUBLASLT_MATMUL_STAGES_8x3	Stage size is 8, number of stages is 3.
CUBLASLT_MATMUL_STAGES_8xAUTO	Stage size is 8, number of stages is selected automatically.
CUBLASLT_MATMUL_STAGES_16xAUTO	Stage size is 16, number of stages is selected automatically.
CUBLASLT_MATMUL_STAGES_32xAUTO	Stage size is 32, number of stages is selected automatically.
CUBLASLT_MATMUL_STAGES_64xAUTO	Stage size is 64, number of stages is selected automatically.
CUBLASLT_MATMUL_STAGES_128xAUTO	Stage size is 128, number of stages is selected automatically.

### 3.3.18. cublasLtNumericalImplFlags\_t

`cublasLtNumericalImplFlags_t` is a set of bit-flags that can be specified to select implementation details that may affect numerical behavior of algorithms.

The following flags can be combined using the bit OR operator "|".

Value	Description
CUBLASLT_NUMERICAL_IMPL_FLAGS_FMA	Specify that the implementation is based on [H,F,D]FMA (fused multiply-add) family instructions.
CUBLASLT_NUMERICAL_IMPL_FLAGS_HMMA	Specify that the implementation is based on HMMA (tensor operation) family instructions.

Value	Description
CUBLASLT_NUMERICAL_IMPL_FLAGS_IMMA	Specify that the implementation is based on IMMA (integer tensor operation) family instructions.
CUBLASLT_NUMERICAL_IMPL_FLAGS_DMMA	Specify that the implementation is based on DMMA (double precision tensor operation) family instructions.
CUBLASLT_NUMERICAL_IMPL_FLAGS_TENSOR_OP_MASK	Mask to filter implementations using any of the above kinds of tensor operations.
CUBLASLT_NUMERICAL_IMPL_FLAGS_OP_TYPE_MASK	Mask to filter implementation details about multiply-accumulate instructions used.
CUBLASLT_NUMERICAL_IMPL_FLAGS_ACCUMULATOR_SPEC_16	Specify that the implementation's inner dot product is using half precision accumulator.
CUBLASLT_NUMERICAL_IMPL_FLAGS_ACCUMULATOR_SPEC_32	Specify that the implementation's inner dot product is using single precision accumulator.
CUBLASLT_NUMERICAL_IMPL_FLAGS_ACCUMULATOR_SPEC_64	Specify that the implementation's inner dot product is using double precision accumulator.
CUBLASLT_NUMERICAL_IMPL_FLAGS_ACCUMULATOR_SPEC_32I	Specify that the implementation's inner dot product is using 32 bit signed integer precision accumulator.
CUBLASLT_NUMERICAL_IMPL_FLAGS_ACCUMULATOR_MASK	Mask to filter implementation details about accumulator used.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_16	Specify that the implementation's inner dot product multiply-accumulate instruction is using half-precision inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_16F	Specify that the implementation's inner dot product multiply-accumulate instruction is using bfloat16 inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_TF32	Specify that the implementation's inner dot product multiply-accumulate instruction is using TF32 inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_32	Specify that the implementation's inner dot product multiply-accumulate instruction is using single-precision inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_64	Specify that the implementation's inner dot product multiply-accumulate instruction is using double-precision inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_INPUT_8I	Specify that the implementation's inner dot product multiply-accumulate instruction is using 8-bit integer inputs.
CUBLASLT_NUMERICAL_IMPL_FLAGS_OP_INPUT_MASK	Mask to filter implementation details about accumulator input used.
CUBLASLT_NUMERICAL_IMPL_FLAGS_GAUSSIAN_COMPLEXITY_REDUCTION	Specify that the implementation applies Gauss complexity reduction algorithm to reduce arithmetic complexity of the complex matrix multiplication problem.

### 3.3.19. `cublasLtMatrixLayout_t`

The `cublasLtMatrixLayout_t` is a pointer to an opaque structure holding the description of a matrix layout. Use the following functions to manipulate this descriptor:

#### **`cublasLtMatrixLayoutCreate()`:**

To create one instance of the descriptor.

#### **`cublasLtMatrixLayoutDestroy()`:**

To destroy a previously created descriptor and release the resources.

### 3.3.20. `cublasLtMatrixLayoutAttribute_t`

`cublasLtMatrixLayoutAttribute_t` is a descriptor structure containing the attributes that define the details of the matrix operation.

Attribute Name	Description	Data Type
<code>CUBLASLT_MATRIX_LAYOUT_TYPE</code>	Specifies the data precision type. See <a href="#">cudaDataType_t</a> .	<code>uint32_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_ORDER</code>	Specifies the memory order of the data of the matrix. Default value is <code>CUBLASLT_ORDER_COL</code> . See <a href="#">cublasLtOrder_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_ROWS</code>	Describes the number of rows in the matrix. Normally only values that can be expressed as <code>int32_t</code> are supported.	<code>uint64_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_COLS</code>	Describes the number of columns in the matrix. Normally only values that can be expressed as <code>int32_t</code> are supported.	<code>uint64_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_LD</code>	The leading dimension of the matrix. For <code>CUBLASLT_ORDER_COL</code> this is the stride (in elements) of matrix column. See also <a href="#">cublasLtOrder_t</a> . <ul style="list-style-type: none"> <li>► Currently only non-negative values are supported.</li> <li>► Must be large enough so that matrix memory locations are not overlapping (for example, greater or equal to <code>CUBLASLT_MATRIX_LAYOUT_ROWS</code> in case of <code>CUBLASLT_ORDER_COL</code>).</li> </ul>	<code>int64_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_BATCH_SIZE</code>	Number of matmul operations to perform in the batch. Default value is 1. See also <code>CUBLASLT_ALGO_CAP_STRIDED_BATCH_SUPPORT</code> in <a href="#">cublasLtMatmulAlgoCapAttributes_t</a> .	<code>int32_t</code>
<code>CUBLASLT_MATRIX_LAYOUT_STRIDED_BATCH_OFFSET</code>	Stride (in elements) to the next matrix for the strided batch operation. Default value is 0. When matrix type is planar-complex	<code>int64_t</code>



Attribute Name	Description	Data Type
	<p>(CUBLASLT_MATRIX_LAYOUT_PLANE_OFFSET != 0), batch stride is interpreted by <code>cublasLtMatmul()</code> in number of real valued sub-elements. For example, for data of type <code>CUDA_C_16F</code>, offset of 1024B is encoded as a stride of value 512 (since each element of the real and imaginary matrices is a 2B (16bit) floating point type).</p> <p>NOTE: A bug in <code>cublasLtMatrixTransform()</code> causes it to interpret the batch stride for a planar-complex matrix as if it was specified in number of complex elements. Therefore an offset of 1024B must be encoded as stride value 256 when calling <code>cublasLtMatrixTransform()</code> (each complex element is 4B with real and imaginary values 2B each). This behavior is expected to be corrected in the next major cuBLAS version.</p>	
CUBLASLT_MATRIX_LAYOUT_PLANE_OFFSET	Stride (in bytes) to the imaginary plane for planar complex layout. Default value is 0, indicating that the layout is regular (real and imaginary parts of complex numbers are interleaved in memory for each element).	int64_t

Use the following functions to manipulate this descriptor:

[`cublasLtMatrixLayoutSetAttribute\(\)`](#): To initialize the descriptor.

[`cublasLtMatrixLayoutGetAttribute\(\)`](#): To query a previously created descriptor.

### 3.3.21. `cublasLtMatrixTransformDesc_t`

The `cublasLtMatrixTransformDesc_t` is a pointer to an opaque structure holding the description of a matrix transformation operation. Use the following functions to manipulate this descriptor:

[`cublasLtMatrixTransformDescCreate\(\)`](#):

To create one instance of the descriptor.

[`cublasLtMatrixTransformDescDestroy\(\)`](#):

To destroy a previously created descriptor and release the resources.

### 3.3.22. `cublasLtMatrixTransformDescAttributes_t`

`cublasLtMatrixTransformDescAttributes_t` is a descriptor structure containing the attributes that define the specifics of the matrix transform operation.

Transform Attribute Name	Description	Data Type
CUBLASLT_MATRIX_TRANSFORM_DESC_SCALETYPE	Scale type inputs are converted to the scale type for scaling and summation, and results are then converted to the output type to store in the memory. For the supported data types refer to <a href="#">cuda_datatype_t</a> .	int32_t
CUBLASLT_MATRIX_TRANSFORM_DESC_POINTER_MODE	Specifies the modes alpha and beta are passed by reference whether on the host or on the device. Default value is CUBLASLT_POINTER_MODE_HOST (i.e., on the host). Refer to <a href="#">cublasLtPointerMode_t</a> .	int32_t
CUBLASLT_MATRIX_TRANSFORM_DESC_OP_A	Specifies the type of operation that should be performed on the matrix A. Default value is CUBLAS_OP_N (such as non-transpose operation). Refer to <a href="#">cublasOperation_t</a> .	int32_t
CUBLASLT_MATRIX_TRANSFORM_DESC_OP_B	Specifies the type of operation that should be performed on the matrix B. Default value is CUBLAS_OP_N (for example, non-transpose operation). Refer to <a href="#">cublasOperation_t</a> .	int32_t

Use the following functions to manipulate this descriptor:

[cublasLtMatrixTransformDescSetAttribute\(\)](#): To set the attribute(s) of the descriptor.

[cublasLtMatrixTransformDescGetAttribute\(\)](#): To query a previously created descriptor for the attribute(s).

### 3.3.23. cublasLtOrder\_t

cublasLtOrder\_t is an enumerated type used to indicate the data ordering of the matrix.

Value	Data Order Description
CUBLASLT_ORDER_COL	Data is ordered in column-major format. The leading dimension is the stride (in elements) to the beginning of next column in memory.
CUBLASLT_ORDER_ROW	Data is ordered in row-major format. The leading dimension is the stride (in elements) to the beginning of next row in memory.
CUBLASLT_ORDER_COL32	Data is ordered in column-major ordered tiles of 32 columns. The leading dimension is the stride (in elements) to the beginning of next group of 32-columns. For example, if the matrix has 33 columns and 2 rows, then the leading dimension must be at least $(32) * 2 = 64$ .
CUBLASLT_ORDER_COL4_4R2_8C	Data is ordered in column-major ordered tiles of composite tiles with total 32 columns and 8 rows. A tile is composed of interleaved inner tiles of 4 columns within 4 even or odd rows in an alternating pattern. The leading dimension is the stride (in elements) to the beginning of the first 32 column x 8 row tile for the next 32-wide group of columns. For example, if the matrix has 33 columns and 1 row, the leading dimension must be at least $(32 * 8) * 1 = 256$ .

Value	Data Order Description
CUBLASLT_ORDER_COL32_2R_4R4	Data is ordered in column-major ordered tiles of composite tiles with total 32 columns and 32 rows. Element offset within the tile is calculated as $((\text{row}\%8)/2*4+\text{row}/8)*2+\text{row}\%2)*32+\text{col}$ . Leading dimension is the stride (in elements) to the beginning of the first 32 column x 32 row tile for the next 32-wide group of columns. For example, if the matrix has 33 columns and 1 row, the leading dimension must be at least $(32*32)*1 = 1024$ .

### 3.3.24. cublasLtPointerMode\_t

`cublasLtPointerMode_t` is an enumerated type used to set the pointer mode for the scaling factors alpha and beta.

Value	Description
CUBLASLT_POINTER_MODE_HOST = CUBLAS_POINTER_MODE_HOST	Matches CUBLAS_POINTER_MODE_HOST, and the pointer targets a single value host memory.
CUBLASLT_POINTER_MODE_DEVICE = CUBLAS_POINTER_MODE_DEVICE	Matches CUBLAS_POINTER_MODE_DEVICE, and the pointer targets a single value device memory.
CUBLASLT_POINTER_MODE_DEVICE_VECTOR = 2	Pointers target device memory vectors of length equal to the number of rows of matrix D.
CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR = 3	Alpha pointer targets a device memory vector of length equal to the number of rows of matrix D, and beta is zero.
CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR = 4	Alpha pointer targets a device memory vector of length equal to the number of rows of matrix D, and beta is a single value in host memory.

### 3.3.25. cublasLtPointerModeMask\_t

`cublasLtPointerModeMask_t` is an enumerated type used to define and query the pointer mode capability.

Value	Description
CUBLASLT_POINTER_MODE_MASK_HOST = 1	Refer to CUBLASLT_POINTER_MODE_HOST in <a href="#">cublasLtPointerMode_t</a> .
CUBLASLT_POINTER_MODE_MASK_DEVICE = 2	Refer to CUBLASLT_POINTER_MODE_DEVICE in <a href="#">cublasLtPointerMode_t</a> .
CUBLASLT_POINTER_MODE_MASK_DEVICE_VECTOR = 4	See CUBLASLT_POINTER_MODE_DEVICE_VECTOR in <a href="#">cublasLtPointerMode_t</a> .
CUBLASLT_POINTER_MODE_MASK_ALPHA_DEVICE_VECTOR = 8	Refer to BETA_ZERO CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_ZERO in <a href="#">cublasLtPointerMode_t</a> .
CUBLASLT_POINTER_MODE_MASK_ALPHA_DEVICE_VECTOR = 16	Refer to BETA_HOST CUBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_HOST in <a href="#">cublasLtPointerMode_t</a> .

### 3.3.26. cublasLtReductionScheme\_t

`cublasLtReductionScheme_t` is an enumerated type used to specify a reduction scheme for the portions of the dot-product calculated in parallel (i.e., "split - K").

Value	Description
<code>CUBLASLT_REDUCTION_SCHEME_NONE</code>	Do not apply reduction. The dot-product will be performed in one sequence.
<code>CUBLASLT_REDUCTION_SCHEME_INPLACE</code>	Reduction is performed "in place" using the output buffer, parts are added up in the output data type. Workspace is only used for counters that guarantee sequentiality.
<code>CUBLASLT_REDUCTION_SCHEME_COMPUTE_TYPE</code>	Reduction done out of place in a user-provided workspace. The intermediate results are stored in the compute type in the workspace and reduced in a separate step.
<code>CUBLASLT_REDUCTION_SCHEME_OUTPUT_TYPE</code>	Reduction done out of place in a user-provided workspace. The intermediate results are stored in the output type in the workspace and reduced in a separate step.
<code>CUBLASLT_REDUCTION_SCHEME_MASK</code>	Allows all reduction schemes.

## 3.4. cuBLASLt API Reference

### 3.4.1. cublasLtCreate()

```
cublasStatus_t
cublasLtCreate(cublasLtHandle_t *lighthandle)
```

This function initializes the cuBLASLt library and creates a handle to an opaque structure holding the cuBLASLt library context. It allocates light hardware resources on the host and device, and must be called prior to making any other cuBLASLt library calls.

The cuBLASLt library context is tied to the current CUDA device. To use the library on multiple devices, one cuBLASLt handle should be created for each device.

#### Parameters:

Parameter	Memory	Input / Output	Description
<code>lightHandle</code>		Output	Pointer to the allocated cuBLASLt handle for the created cuBLASLt context.

#### Returns:

Return Value	Description
<code>CUBLAS_STATUS_SUCCESS</code>	The allocation completed successfully.

Return Value	Description
CUBLAS_STATUS_NOT_INITIALIZED	The cuBLASLt library was not initialized. This usually happens: <ul style="list-style-type: none"> <li>- when <code>cublasLtCreate()</code> is not called first</li> <li>- an error in the CUDA Runtime API called by the cuBLASLt routine, or</li> <li>- an error in the hardware setup.</li> </ul>
CUBLAS_STATUS_ALLOC_FAILED	Resource allocation failed inside the cuBLASLt library. This is usually caused by a <code>cudaMalloc()</code> failure.  To correct: prior to the function call, deallocate the previously allocated memory as much as possible.
CUBLAS_STATUS_INVALID_VALUE	<code>lightHandle == NULL</code>

See [cublasStatus\\_t](#) for a complete list of valid return codes.

### 3.4.2. cublasLtDestroy()

```
cublasStatus_t
cublasLtDestroy(cublasLtHandle_t lightHandle)
```

This function releases hardware resources used by the cuBLASLt library. This function is usually the last call with a particular handle to the cuBLASLt library. Because `cublasLtCreate()` allocates some internal resources and the release of those resources by calling `cublasLtDestroy()` will implicitly call `cudaDeviceSynchronize()`, it is recommended to minimize the number of `cublasLtCreate()` / `cublasLtDestroy()` occurrences.

#### Parameters:

Parameter	Memory	Input / Output	Description
<code>lightHandle</code>		Input	Pointer to the cuBLASLt handle to be destroyed.

#### Returns:

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The cuBLASLt context was successfully destroyed.
CUBLAS_STATUS_NOT_INITIALIZED	The cuBLASLt library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	<code>lightHandle == NULL</code>

See [cublasStatus\\_t](#) for a complete list of valid return codes.

### 3.4.3. `cublasLtGetCudartVersion()`

```
size_t cublasLtGetCudartVersion(void);
```

This function returns the version number of the CUDA Runtime library.

**Parameters:** None.

**Returns:** `size_t` - The version number of the CUDA Runtime library.

### 3.4.4. `cublasLtGetProperty()`

```
cublasStatus_t cublasLtGetProperty(libraryPropertyType type, int *value);
```

This function returns the value of the requested property by writing it to the memory location pointed to by the value parameter.

**Parameters:**

Parameter	Memory	Input / Output	Description
type		Input	Of the type <code>libraryPropertyType</code> , whose value is requested from the property. See <a href="#">libraryPropertyType</a> .
value		Output	Pointer to the host memory location where the requested information should be written.

**Returns:**

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The requested <code>libraryPropertyType</code> information is successfully written at the provided address.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<ul style="list-style-type: none"> <li>► If invalid value of the <code>type</code> input argument or</li> <li>► <code>value == NULL</code></li> </ul>

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.5. `cublasLtGetStatusName()`

```
const char* cublasLtGetStatusName(cublasStatus_t status);
```

Returns the string representation of a given status.

**Parameters:** `cublasStatus_t` - the status.

**Returns:** `const char*` - the NULL-terminated string.

### 3.4.6. `cublasLtGetStatusString()`

```
const char* cublasLtGetStatusString(cublasStatus_t status);
```

Returns the description string for a given status.

**Parameters:** `cublasStatus_t` - the status.

**Returns:** `const char*` - the NULL-terminated string.

### 3.4.7. `cublasLtHeuristicsCacheGetCapacity()`

```
cublasStatus_t cublasLtHeuristicsCacheGetCapacity(size_t* capacity);
```

Returns the [Heuristics Cache](#) capacity.

**Parameters:**

Parameter	Description
<code>capacity</code>	The pointer to the returned capacity value.

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_SUCCESS</code>	The capacity was successfully written.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<code>lighthandle == NULL</code> .

### 3.4.8. `cublasLtHeuristicsCacheSetCapacity()`

```
cublasStatus_t cublasLtHeuristicsCacheSetCapacity(size_t capacity);
```

Sets the [Heuristics Cache](#) capacity. Set the capacity to 0 to disable the heuristics cache.

This function takes precedence over `CUBLASLT_HEURISTICS_CACHE_CAPACITY` environment variable.

**Parameters:**

Parameter	Description
<code>capacity</code>	The desirable heuristics cache capacity.

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_SUCCESS</code>	The capacity was successfully set.

### 3.4.9. cublasLtGetVersion()

```
size_t cublasLtGetVersion(void);
```

This function returns the version number of cuBLASLt library.

**Parameters:** None.

**Returns:** `size_t` - The version number of cuBLASLt library.

### 3.4.10. cublasLtLoggerSetCallback()

```
cublasStatus_t cublasLtLoggerSetCallback(cublasLtLoggerCallback_t callback);
```

Experimental: This function sets the logging callback function.

**Parameters:**

Parameter	Memory	Input / Output	Description
callback		Input	Pointer to a callback function. See <a href="#">cublasLtLoggerCallback_t</a> .

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If the callback function was successfully set.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.11. cublasLtLoggerSetFile()

```
cublasStatus_t cublasLtLoggerSetFile(FILE* file);
```

Experimental: This function sets the logging output file. Note: once registered using this function call, the provided file handle must not be closed unless the function is called again to switch to a different file handle.

**Parameters:**

Parameter	Memory	Input / Output	Description
file		Input	Pointer to an open file. File should have write permission.

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If logging file was successfully set.



See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.12. `cublasLtLoggerOpenFile()`

```
cublasStatus_t cublasLtLoggerOpenFile(const char* logFile);
```

Experimental: This function opens a logging output file in the given path.

**Parameters:**

Parameter	Memory	Input / Output	Description
logFile		Input	Path of the logging output file.

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If the logging file was successfully opened.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.13. `cublasLtLoggerSetLevel()`

```
cublasStatus_t cublasLtLoggerSetLevel(int level);
```

Experimental: This function sets the value of the logging level.

**Parameters:**

Parameter	Memory	Input / Output	Description
level		Input	Value of the logging level. See <a href="#">cuBLASLt Logging</a> .

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If the value was not a valid logging level. See <a href="#">cuBLASLt Logging</a> .
CUBLAS_STATUS_SUCCESS	If the logging level was successfully set.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.14. `cublasLtLoggerSetMask()`

```
cublasStatus_t cublasLtLoggerSetMask(int mask);
```

Experimental: This function sets the value of the logging mask.

**Parameters:**

Parameter	Memory	Input / Output	Description
mask		Input	Value of the logging mask. See <a href="#">cuBLASLt Logging</a> .

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If the logging mask was successfully set.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.15. `cublasLtLoggerForceDisable()`

```
cublasStatus_t cublasLtLoggerForceDisable();
```

Experimental: This function disables logging for the entire run.

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If logging was successfully disabled.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.16. `cublasLtMatmul()`

```
cublasStatus_t cublasLtMatmul(
    cublasLtHandle_t          lightHandle,
    cublasLtMatmulDesc_t      computeDesc,
    const void*               *alpha,
    const void*               *A,
    cublasLtMatrixLayout_t     Adesc,
    const void*               *B,
    cublasLtMatrixLayout_t     Bdesc,
    const void*               *beta,
    const void*               *C,
    cublasLtMatrixLayout_t     Cdesc,
    void*                     *D,
    cublasLtMatrixLayout_t     Ddesc,
    const cublasLtMatmulAlgo_t *algo,
    void*                     *workspace,
    size_t                    workspaceSizeInBytes,
    cudaStream_t               stream);
```

This function computes the matrix multiplication of matrices A and B to produce the output matrix D, according to the following operation:

$$D = \alpha * (A * B) + \beta * (C),$$

where  $A$ ,  $B$ , and  $C$  are input matrices, and  $\alpha$  and  $\beta$  are input scalars.

**Note:** This function supports both in-place matrix multiplication ( $C == D$  and  $Cdesc == Ddesc$ ) and out-of-place matrix multiplication ( $C != D$ , both matrices must have the same data type, number of rows, number of columns, batch size, and memory order). In the out-of-place case, the leading dimension of  $C$  can be different from the leading dimension of  $D$ . Specifically the leading dimension of  $C$  can be 0 to achieve row or column broadcast. If  $Cdesc$  is omitted, this function assumes it to be equal to  $Ddesc$ .

### Datatypes Supported:

`cusblasLtMatmul` supports the following `computeType`, `scaleType`, `Atype/Btype`, and `Ctype`:

Table 1. When  $A$ ,  $B$ ,  $C$ , and  $D$  are Regular Column- or Row-major Matrices

computeType	scaleType	Atype/Btype	Ctype	Bias Type <sup>4</sup>
CUBLAS_COMPUTE_16F or CUBLAS_COMPUTE_16F_PEDANTIC	CUDA_R_16F	CUDA_R_16F	CUDA_R_16F	CUDA_R_16F <sup>4</sup>
CUBLAS_COMPUTE_32I or CUBLAS_COMPUTE_32I_PEDANTIC	CUDA_R_32I	CUDA_R_8I	CUDA_R_32I	Non-default epilogue not supported.
	CUDA_R_32F <sup>1</sup>	CUDA_R_8I <sup>1</sup>	CUDA_R_8I <sup>1</sup>	Non-default epilogue not supported.
CUBLAS_COMPUTE_32F or CUBLAS_COMPUTE_32F_PEDANTIC	CUDA_R_32F	CUDA_R_16BF	CUDA_R_16BF	CUDA_R_16BF <sup>4</sup>
		CUDA_R_16F	CUDA_R_16F	CUDA_R_16F <sup>4</sup>
		CUDA_R_8I	CUDA_R_32F	Non-default epilogue not supported.
		CUDA_R_16BF	CUDA_R_32F	CUDA_R_32F <sup>4</sup>
		CUDA_R_16F	CUDA_R_32F	CUDA_R_32F <sup>4</sup>
		CUDA_R_32F	CUDA_R_32F	CUDA_R_32F <sup>4</sup>
	CUDA_C_32F <sup>5</sup>	CUDA_C_8I <sup>5</sup>	CUDA_C_32F <sup>5</sup>	Non-default epilogue not supported.
		CUDA_C_32F <sup>5</sup>	CUDA_C_32F <sup>5</sup>	
CUBLAS_COMPUTE_32F_FAST_16F or CUBLAS_COMPUTE_32F_FAST_16BF or	CUDA_R_32F	CUDA_R_32F	CUDA_R_32F	CUDA_R_32F <sup>4</sup>
	CUDA_C_32F <sup>5</sup>	CUDA_C_32F <sup>5</sup>	CUDA_C_32F <sup>5</sup>	Non-default epilogue not supported.

computeType	scaleType	Atype/Btype	Ctype	Bias Type <sup>4</sup>
CUBLAS_COMPUTE_32F_FAST_TF32				
CUBLAS_COMPUTE_64F or CUBLAS_COMPUTE_64F_PEDANTIC	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F	CUDA_R_64F <sup>4</sup>
	CUDA_C_64F <sup>5</sup>	CUDA_C_64F <sup>5</sup>	CUDA_C_64F <sup>5</sup>	Non-default epilogue not supported.

To use IMMA kernels, one of the following sets of requirements, with the first being the preferred one, must be met:

- Using a regular data ordering:
  - ▶ All matrix pointers must be 4-byte aligned. For even better performance, this condition should hold with 16 instead of 4.
  - ▶ Leading dimensions of matrices A, B, C must be multiples of 4.
  - ▶ Only the "TN" format is supported - A must be transposed and B non-transposed.
  - ▶ Dimensions *m* and *k* must be multiples of 4.
- Using the IMMA-specific data ordering - CUBLASLT\_ORDER\_COL32 for matrices A,C,D, and CUBLASLT\_ORDER\_COL4\_4R2\_8C (on Turing or NVIDIA Ampere GPU architecture) or CUBLASLT\_ORDER\_COL32\_2R\_4R4 (on NVIDIA Ampere GPU architecture) for matrix B:
  - ▶ Leading dimensions of matrices A, B, C must fulfill conditions specific to the memory ordering (see [cublasLtOrder\\_t](#)).
  - ▶ Matmul descriptor must specify CUBLAS\_OP\_T on matrix B and CUBLAS\_OP\_N (default) on matrix A and C.
  - ▶ If scaleType CUDA\_R\_32I is used, the only supported values for *alpha* and *beta* are 0 or 1.

See below table when using IMMA kernels. Note that IMMA does not work with CUBLASLT\_POINTER\_MODE\_MASK\_DEVICE pointer mode.

Table 2. When A, B, C, and D Use Layouts for IMMA

computeType	scaleType	Atype/Btype	Ctype	Bias Type <sup>2</sup>
CUBLAS_COMPUTE_32I or CUBLAS_COMPUTE_32I_PEDANTIC	CUDA_R_32I	CUDA_R_8I	CUDA_R_32I	Non-default epilogue not supported.
	CUDA_R_32F <sup>2</sup>	CUDA_R_8I	CUDA_R_8I	CUDA_R_32F <sup>2</sup>

To use FP8 kernels, the following set of requirements must be satisfied:

- ▶ All matrix pointers must be 16-byte aligned.

- A must be transposed and B non-transposed (The "TN" format).
- The compute type must be CUBLAS\_COMPUTE\_32F.
- The scale type must be CUDA\_R\_32F.

See the table below when using FP8 kernels:

Table 3. When A, B, C, and D Use Layouts for FP8

AType	BType	CType	DType	Bias Type
CUDA_R_8F_E4M3	CUDA_R_8F_E4M3	CUDA_R_16BF	CUDA_R_16BF	CUDA_R_16BF
			CUDA_R_8F_E4M3	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_8F_E4M3	CUDA_R_16F
			CUDA_R_16F	CUDA_R_16F
	CUDA_R_8F_E5M2	CUDA_R_16BF	CUDA_R_16BF	CUDA_R_16BF
			CUDA_R_8F_E4M3	CUDA_R_16BF
			CUDA_R_8F_E5M2	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_8F_E4M3	CUDA_R_16F
			CUDA_R_8F_E5M2	CUDA_R_16F
			CUDA_R_16F	CUDA_R_16F
		CUDA_R_32F	CUDA_R_32F	CUDA_R_16BF
CUDA_R_8F_E5M2	CUDA_R_8F_E4M3	CUDA_R_16BF	CUDA_R_16BF	CUDA_R_16BF
			CUDA_R_8F_E4M3	CUDA_R_16BF
			CUDA_R_8F_E5M2	CUDA_R_16BF
		CUDA_R_16F	CUDA_R_8F_E4M3	CUDA_R_16F
			CUDA_R_8F_E5M2	CUDA_R_16F
			CUDA_R_16F	CUDA_R_16F
		CUDA_R_32F	CUDA_R_32F	CUDA_R_16BF

And finally, see below table when A,B,C,D are planar complex matrices (see CUBLASLT\_MATRIX\_LAYOUT\_PLANE\_OFFSET) to make use of mixed precision tensor core acceleration.

Table 4. When A, B, C, and D are Planar Complex Matrices

computeType	scaleType	Atype/Btype	Ctype
CUBLAS_COMPUTE_32F	CUDA_C_32F	CUDA_C_16F <sup>3,5</sup>	CUDA_C_16F <sup>3,5</sup>
			CUDA_C_32F <sup>3,5</sup>
		CUDA_C_16BF <sup>3,5</sup>	CUDA_C_16BF <sup>3,5</sup>

computeType	scaleType	Atype/Btype	Ctype
			CUDA_C_32F <sup>3,5</sup>

**NOTES:**

1. When using regular memory order and when compute type is 32I, input type is R\_8I and output type is R\_8I, only "TN" format is supported - "A" must be transposed and "B" non-transposed.
2. IMMA kernel with computeType=32I and Ctype=CUDA\_R\_8I supports per row scaling (see CUBLASLT\_POINTER\_MODE\_DEVICE\_VECTOR and CUBLASLT\_POINTER\_MODE\_ALPHA\_DEVICE\_VECTOR\_BETA\_ZERO in [cublasLtPointerMode\\_t](#)) as well as ReLU and Bias epilogue modes (see CUBLASLT\_MATMUL\_DESC\_EPILOGUE in [cublasLtMatmulDescAttributes\\_t](#)).
3. These can only be used with planar layout (CUBLASLT\_MATRIX\_LAYOUT\_PLANE\_OFFSET != 0).
4. ReLU, dReLU, GELU, dGELU and Bias epilogue modes (see CUBLASLT\_MATMUL\_DESC\_EPILOGUE in [cublasLtMatmulDescAttributes\\_t](#)) are not supported when D matrix memory order is defined as CUBLASLT\_ORDER\_ROW. For best performance when using the bias vector, specify beta == 0 and CUBLASLT\_POINTER\_MODE\_HOST.
5. Use of CUBLAS\_ORDER\_ROW together with CUBLAS\_OP\_C (hermitian operator) is not supported unless all of A, B, C, D matrices are defined with CUBLAS\_ORDER\_ROW.

**Parameters:**

Parameter	Memory	Input / Output	Description
lightHandle		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <a href="#">cublasLtHandle_t</a> .
computeDesc		Input	Handle to a previously created matrix multiplication descriptor of type <a href="#">cublasLtMatmulDesc_t</a> .
alpha, beta	Device or host	Input	Pointers to the scalars used in the multiplication.
A, B, and C	Device	Input	Pointers to the GPU memory associated with the corresponding descriptors Adesc, Bdesc and Cdesc.
Adesc, Bdesc and Cdesc.		Input	Handles to the previous created descriptors of the type <a href="#">cublasLtMatrixLayout_t</a> .
D	Device	Output	Pointer to the GPU memory associated

Parameter	Memory	Input / Output	Description
			with the descriptor <code>Ddesc</code> .
<code>Ddesc</code>		Input	Handle to the previous created descriptor of the type <code>cublasLtMatrixLayout_t</code> .
<code>algo</code>		Input	Handle for matrix multiplication algorithm to be used. See <code>cublasLtMatmulAlgo_t</code> . When NULL, an implicit heuristics query with default search preferences will be performed to determine actual algorithm to use.
<code>workspace</code>	Device		Pointer to the workspace buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of address must be 0).
<code>workspaceSizeInBytes</code>		Input	Size of the workspace.
<code>stream</code>	Host	Input	The CUDA stream where all the GPU work will be submitted.

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	If cuBLASLt handle has not been initialized.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	If the parameters are unexpectedly NULL, in conflict or in an impossible configuration. For example, when <code>workspaceSizeInBytes</code> is less than workspace required by the configured algo.
<code>CUBLAS_STATUS_NOT_SUPPORTED</code>	If the current implementation on the selected device doesn't support the configured operation.
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	If the configured operation cannot be run using the selected device.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	If CUDA reported an execution error from the device.
<code>CUBLAS_STATUS_SUCCESS</code>	If the operation completed successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.17. cublasLtMatmulAlgoCapGetAttribute()

```
cublasStatus_t cublasLtMatmulAlgoCapGetAttribute(  
    const cublasLtMatmulAlgo_t *algo,  
    cublasLtMatmulAlgoCapAttributes_t attr,  
    void *buf,  
    size_t sizeInBytes,  
    size_t *sizeWritten);
```

This function returns the value of the queried capability attribute for an initialized `cublasLtMatmulAlgo_t` descriptor structure. The capability attribute value is retrieved from the enumerated type `cublasLtMatmulAlgoCapAttributes_t`.

For example, to get the list of supported Tile IDs:

```
cublasLtMatmulTile_t tiles[CUBLASLT_MATMUL_TILE_END];  
size_t num_tiles, size_written;  
if (cublasLtMatmulAlgoCapGetAttribute(algo, CUBLASLT_ALGO_CAP_TILE_IDS,  
    tiles, sizeof(tiles), &size_written) == CUBLAS_STATUS_SUCCESS) {  
    num_tiles = size_written / sizeof(tiles[0]);  
}
```

**Parameters:**

Parameter	Memory	Input / Output	Description
algo		Input	Pointer to the previously created opaque structure holding the matrix multiply algorithm descriptor. See <code>cublasLtMatmulAlgo_t</code> .
attr		Input	The capability attribute whose value will be retrieved by this function. See <code>cublasLtMatmulAlgoCapAttributes_t</code> .
buf		Output	The attribute value returned by this function.
sizeInBytes		Input	Size of <code>buf</code> buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is <code>CUBLAS_STATUS_SUCCESS</code> . If <code>sizeInBytes</code> is non-zero: then <code>sizeWritten</code> is the number of bytes actually written; if <code>sizeInBytes</code> is 0: then <code>sizeWritten</code> is the number of bytes needed to write full contents.

**Returns:**




Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"><li>▶ If sizeInBytes is 0 and sizeWritten is NULL, or</li><li>▶ if sizeInBytes is non-zero and buf is NULL, or</li><li>▶ sizeInBytes doesn't match size of internal storage for the selected attribute</li></ul>
CUBLAS_STATUS_SUCCESS	If attribute's value was successfully written to user memory.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.18. `cublasLtMatmulAlgoCheck()`

```
cublasStatus_t cublasLtMatmulAlgoCheck(  
    cublasLtHandle_t lightHandle,  
    cublasLtMatmulDesc_t operationDesc,  
    cublasLtMatrixLayout_t Adesc,  
    cublasLtMatrixLayout_t Bdesc,  
    cublasLtMatrixLayout_t Cdesc,  
    cublasLtMatrixLayout_t Ddesc,  
    const cublasLtMatmulAlgo_t *algo,  
    cublasLtMatmulHeuristicResult_t *result);
```

This function performs the correctness check on the matrix multiply algorithm descriptor for the matrix multiply operation `cublasLtMatmul()` function with the given input matrices A, B and C, and the output matrix D. It checks whether the descriptor is supported on the current device, and returns the result containing the required workspace and the calculated wave count.



**Note:** `CUBLAS_STATUS_SUCCESS` doesn't fully guarantee that the algo will run. The algo will fail if, for example, the buffers are not correctly aligned. However, if `cublasLtMatmulAlgoCheck` fails, the algo will not run.

**Parameters:**

Parameter	Memory	Input / Output	Description
lightHandle		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <code>cublasLtHandle_t</code> .
operationDesc		Input	Handle to a previously created matrix multiplication descriptor of type <code>cublasLtMatmulDesc_t</code> .
Adesc, Bdesc, Cdesc, and Ddesc		Input	Handles to the previously created

Parameter	Memory	Input / Output	Description
preference		Input	matrix layout descriptors of the type <code>cublasLtMatrixLayout_t</code> . Pointer to the structure holding the matrix multiply preferences descriptor. See <code>cublasLtMatrixLayout_t</code> .
algo		Input	Descriptor which specifies which matrix multiplication algorithm should be used. See <code>cublasLtMatmulAlgo_t</code> . May point to result → algo.
result		Output	Pointer to the structure holding the results returned by this function. The results comprise of the required workspace and the calculated wave count. The algo field is never updated. See <code>cublasLtMatmulHeuristicResult_t</code> .

Returns:

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If matrix layout descriptors or the operation descriptor do not match the algo descriptor.
CUBLAS_STATUS_NOT_SUPPORTED	If the algo configuration or data type combination is not currently supported on the given device.
CUBLAS_STATUS_ARCH_MISMATCH	If the algo configuration cannot be run using the selected device.
CUBLAS_STATUS_SUCCESS	If the check was successful.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.19. cublasLtMatmulAlgoConfigGetAttribute()

```
cublasStatus_t cublasLtMatmulAlgoConfigGetAttribute(  
    const cublasLtMatmulAlgo_t *algo,  
    cublasLtMatmulAlgoConfigAttributes_t attr,  
    void *buf,  
    size_t sizeInBytes,  
    size_t *sizeWritten);
```

This function returns the value of the queried configuration attribute for an initialized `cublasLtMatmulAlgo_t` descriptor. The configuration attribute value is retrieved from the enumerated type `cublasLtMatmulAlgoConfigAttributes_t`.

Parameters:

Parameter	Memory	Input / Output	Description
algo		Input	Pointer to the previously created opaque structure holding the matrix multiply algorithm descriptor. See <code>cublasLtMatmulAlgo_t</code> .
attr		Input	The configuration attribute whose value will be retrieved by this function. See <code>cublasLtMatmulAlgoConfigAttributes_t</code> .
buf		Output	The attribute value returned by this function.
sizeInBytes		Input	Size of <code>buf</code> buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is <code>CUBLAS_STATUS_SUCCESS</code> . If <code>sizeInBytes</code> is non-zero: then <code>sizeWritten</code> is the number of bytes actually written; if <code>sizeInBytes</code> is 0: then <code>sizeWritten</code> is the number of bytes needed to write full contents.

Returns:

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"><li>▶ If <code>sizeInBytes</code> is 0 and <code>sizeWritten</code> is NULL, or</li><li>▶ if <code>sizeInBytes</code> is non-zero and <code>buf</code> is NULL, or</li><li>▶ <code>sizeInBytes</code> doesn't match size of internal storage for the selected attribute</li></ul>
CUBLAS_STATUS_SUCCESS	If attribute's value was successfully written to user memory.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.20. cublasLtMatmulAlgoConfigSetAttribute()

```
cublasStatus_t cublasLtMatmulAlgoConfigSetAttribute(  
    cublasLtMatmulAlgo_t *algo,  
    cublasLtMatmulAlgoConfigAttributes_t attr,  
    const void *buf,  
    size_t sizeInBytes);
```

This function sets the value of the specified configuration attribute for an initialized cublasLtMatmulAlgo\_t descriptor. The configuration attribute is an enumerator of the type cublasLtMatmulAlgoConfigAttributes\_t.

**Parameters:**

Parameter	Memory	Input / Output	Description
algo		Input	Pointer to the previously created opaque structure holding the matrix multiply algorithm descriptor. See cublasLtMatmulAlgo_t.
attr		Input	The configuration attribute whose value will be set by this function. See cublasLtMatmulAlgoConfigAttributes_t.
buf		Input	The value to which the configuration attribute should be set.
sizeInBytes		Input	Size of buf buffer (in bytes) for verification.

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If buf is NULL or sizeInBytes doesn't match the size of the internal storage for the selected attribute.
CUBLAS_STATUS_SUCCESS	If the attribute was set successfully.

See cublasStatus\_t for a complete list of valid return codes.

### 3.4.21. cublasLtMatmulAlgoGetHeuristic()

```
cublasStatus_t cublasLtMatmulAlgoGetHeuristic(  
    cublasLtHandle_t lightHandle,  
    cublasLtMatmulDesc_t operationDesc,  
    cublasLtMatrixLayout_t Adesc,  
    cublasLtMatrixLayout_t Bdesc,  
    cublasLtMatrixLayout_t Cdesc,
```

```

cublasLtMatrixLayout_t Ddesc,
cublasLtMatmulPreference_t preference,
int requestedAlgoCount,
cublasLtMatmulHeuristicResult_t heuristicResultsArray[]
int *returnAlgoCount);

```

This function retrieves the possible algorithms for the matrix multiply operation `cublasLtMatmul()` function with the given input matrices A, B and C, and the output matrix D. The output is placed in `heuristicResultsArray[]` in the order of increasing estimated compute time.

#### Parameters:

Parameter	Memory	Input / Output	Description
<code>lightHandle</code>		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <code>cublasLtHandle_t</code> .
<code>operationDesc</code>		Input	Handle to a previously created matrix multiplication descriptor of type <code>cublasLtMatmulDesc_t</code> .
<code>Adesc, Bdesc, Cdesc, and Ddesc</code>		Input	Handles to the previously created matrix layout descriptors of the type <code>cublasLtMatrixLayout_t</code> .
<code>preference</code>		Input	Pointer to the structure holding the heuristic search preferences descriptor. See <code>cublasLtMatrixLayout_t</code> .
<code>requestedAlgoCount</code>		Input	Size of the <code>heuristicResultsArray</code> (in elements). This is the requested maximum number of algorithms to return.
<code>heuristicResultsArray[]</code>		Output	Array containing the algorithm heuristics and associated runtime characteristics, returned by this function, in the order of increasing estimated compute time.
<code>returnAlgoCount</code>		Output	Number of algorithms returned by this function. This is

Parameter	Memory	Input / Output	Description
			the number of <code>heuristicResultsArray</code> elements written.

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_INVALID_VALUE</code>	If <code>requestedAlgoCount</code> is less or equal to zero.
<code>CUBLAS_STATUS_NOT_SUPPORTED</code>	If no heuristic function available for current configuration.
<code>CUBLAS_STATUS_SUCCESS</code>	If query was successful. Inspect <code>heuristicResultsArray[0 to (returnAlgoCount - 1)].state</code> for the status of the results.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.22. `cublasLtMatmulAlgoGetIds()`

```
cublasStatus_t cublasLtMatmulAlgoGetIds(
    cublasLtHandle_t lightHandle,
    cublasComputeType_t computeType,
    cudaDataType_t scaleType,
    cudaDataType_t Atype,
    cudaDataType_t Btype,
    cudaDataType_t Ctype,
    cudaDataType_t Dtype,
    int requestedAlgoCount,
    int algoIdsArray[],
    int *returnAlgoCount);
```

This function retrieves the IDs of all the matrix multiply algorithms that are valid, and can potentially be run by the `cublasLtMatmul()` function, for given types of the input matrices A, B and C, and of the output matrix D.

Note: the IDs are returned in no particular order. To make sure the best possible algo is contained in the list, make `requestedAlgoCount` large enough to receive the full list. The list is guaranteed to be full if `returnAlgoCount < requestedAlgoCount`.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>lightHandle</code>		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <code>cublasLtHandle_t</code> .
<code>computeType,</code> <code>scaleType,</code> <code>Atype,</code> <code>Btype,</code> <code>Ctype,</code> and <code>Dtype</code>		Inputs	Data types of the computation type, scaling factors and of

Parameter	Memory	Input / Output	Description
requestedAlgoCount		Input	the operand matrices. See <a href="#">cudaDataType_t</a> . Number of algorithms requested. Must be > 0.
algoIdsArray[]		Output	Array containing the algorithm IDs returned by this function.
returnAlgoCount		Output	Number of algorithms actually returned by this function.

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If requestedAlgoCount is less or equal to zero.
CUBLAS_STATUS_SUCCESS	If query was successful. Inspect returnAlgoCount to get actual number of IDs available.

See [cublasStatus\\_t](#) for a complete list of valid return codes.

### 3.4.23. [cublasLtMatmulAlgoInit\(\)](#)

```
cublasStatus_t cublasLtMatmulAlgoInit(
    cublasLtHandle_t lightHandle,
    cublasComputeType_t computeType,
    cudaDataType_t scaleType,
    cudaDataType_t Atype,
    cudaDataType_t Btype,
    cudaDataType_t Ctype,
    cudaDataType_t Dtype,
    int algoId,
    cublasLtMatmulAlgo_t *algo);
```

This function initializes the matrix multiply algorithm structure for the [cublasLtMatmul\(\)](#), for a specified matrix multiply algorithm and input matrices A, B and C, and the output matrix D.

**Parameters:**

Parameter	Memory	Input / Output	Description
lightHandle		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <a href="#">cublasLtHandle_t</a> .
computeType		Input	Compute type. See CUBLASLT_MATMUL_DESC_COMPUTE_TYPE of <a href="#">cublasLtMatmulDescAttributes_t</a> .

Parameter	Memory	Input / Output	Description
scaleType		Input	Scale type. See <a href="#">CUBLASLT_MATMUL_DESC_SCALE_TYPE</a> of <a href="#">cublasLtMatmulDescAttributes_t</a> . Usually same as computeType.
Atype, Btype, Ctype, and Dtype		Input	Datatype precision for the input and output matrices. See <a href="#">cudaDataType_t</a> .
algoId		Input	Specifies the algorithm being initialized. Should be a valid algoId returned by the <a href="#">cublasLtMatmulAlgoGetIds()</a> function.
algo		Input	Pointer to the opaque structure to be initialized. See <a href="#">cublasLtMatmulAlgo_t</a> .

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If algo is NULL or algoId is outside the recognized range.
CUBLAS_STATUS_NOT_SUPPORTED	If algoId is not supported for given combination of data types.
CUBLAS_STATUS_SUCCESS	If the structure was successfully initialized.

See [cublasStatus\\_t](#) for a complete list of valid return codes.

### 3.4.24. [cublasLtMatmulDescCreate\(\)](#)

```
cublasStatus_t cublasLtMatmulDescCreate( cublasLtMatmulDesc_t *matmulDesc,
                                          cublasComputeType_t computeType,
                                          cudaDataType_t scaleType);
```

This function creates a matrix multiply descriptor by allocating the memory needed to hold its opaque structure.

**Parameters:**

Parameter	Memory	Input / Output	Description
matmulDesc		Output	Pointer to the structure holding the matrix multiply descriptor created



Parameter	Memory	Input / Output	Description
computeType		Input	by this function. See <code>cublasLtMatmulDesc_t</code> . Enumerant that specifies the data precision for the matrix multiply descriptor this function creates. See <code>cublasComputeType_t</code> .
scaleType		Input	Enumerant that specifies the data precision for the matrix transform descriptor this function creates. See <code>cudaDataType</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.25. `cublasLtMatmulDescInit()`

```
cublasStatus_t cublasLtMatmulDescInit( cublasLtMatmulDesc_t matmulDesc,
                                       cublasComputeType_t computeType,
                                       cudaDataType_t scaleType);
```

This function initializes a matrix multiply descriptor in a previously allocated one.

**Parameters:**

Parameter	Memory	Input / Output	Description
matmulDesc		Output	Pointer to the structure holding the matrix multiply descriptor initialized by this function. See <code>cublasLtMatmulDesc_t</code> .
computeType		Input	Enumerant that specifies the data precision for the matrix multiply descriptor this function initializes. See <code>cublasComputeType_t</code> .
scaleType		Input	Enumerant that specifies the data

Parameter	Memory	Input / Output	Description
			precision for the matrix transform descriptor this function initializes. See <code>cudaDataType</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.26. `cublasLtMatmulDescDestroy()`

```
cublasStatus_t cublasLtMatmulDescDestroy(
    cublasLtMatmulDesc_t matmulDesc);
```

This function destroys a previously created matrix multiply descriptor object.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>matmulDesc</code>		Input	Pointer to the structure holding the matrix multiply descriptor that should be destroyed by this function. See <code>cublasLtMatmulDesc_t</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_SUCCESS</code>	If operation was successful.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.27. `cublasLtMatmulDescGetAttribute()`

```
cublasStatus_t cublasLtMatmulDescGetAttribute(
    cublasLtMatmulDesc_t matmulDesc,
    cublasLtMatmulDescAttributes_t attr,
    void *buf,
    size_t sizeInBytes,
    size_t *sizeWritten);
```

This function returns the value of the queried attribute belonging to a previously created matrix multiply descriptor.

**Parameters:**

Parameter	Memory	Input / Output	Description
matmulDesc		Input	Pointer to the previously created structure holding the matrix multiply descriptor queried by this function. See <a href="#">cublasLtMatmulDesc_t</a> .
attr		Input	The attribute that will be retrieved by this function. See <a href="#">cublasLtMatmulDescAttributes_t</a> .
buf		Output	Memory address containing the attribute value retrieved by this function.
sizeInBytes		Input	Size of buf buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is CUBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

Returns:

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"><li>▶ If sizeInBytes is 0 and sizeWritten is NULL, or</li><li>▶ if sizeInBytes is non-zero and buf is NULL, or</li><li>▶ sizeInBytes doesn't match size of internal storage for the selected attribute</li></ul>
CUBLAS_STATUS_SUCCESS	If attribute's value was successfully written to user memory.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.28. `cublasLtMatmulDescSetAttribute()`

```
cublasStatus_t cublasLtMatmulDescSetAttribute(  
    cublasLtMatmulDesc_t matmulDesc,  
    cublasLtMatmulDescAttributes_t attr,
```

```
const void *buf,
size_t sizeInBytes);
```

This function sets the value of the specified attribute belonging to a previously created matrix multiply descriptor.

Parameters:

Parameter	Memory	Input / Output	Description
matmulDesc		Input	Pointer to the previously created structure holding the matrix multiply descriptor queried by this function. See <a href="#">cublasLtMatmulDesc_t</a> .
attr		Input	The attribute that will be set by this function. See <a href="#">cublasLtMatmulDescAttributes_t</a> .
buf		Input	The value to which the specified attribute should be set.
sizeInBytes		Input	Size of buf buffer (in bytes) for verification.

Returns:

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If buf is NULL or sizeInBytes doesn't match the size of the internal storage for the selected attribute.
CUBLAS_STATUS_SUCCESS	If the attribute was set successfully.

See [cublasStatus\\_t](#) for a complete list of valid return codes.

### 3.4.29. cublasLtMatmulPreferenceCreate()

```
cublasStatus_t cublasLtMatmulPreferenceCreate(
    cublasLtMatmulPreference_t *pref);
```

This function creates a matrix multiply heuristic search preferences descriptor by allocating the memory needed to hold its opaque structure.

Parameters:

Parameter	Memory	Input / Output	Description
n		Output	Pointer to the structure holding the matrix multiply preferences

Parameter	Memory	Input / Output	Description
			descriptor created by this function. See <code>cublasLtMatrixLayout_t</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.30. `cublasLtMatmulPreferenceInit()`

```
cublasStatus_t cublasLtMatmulPreferenceInit(
    cublasLtMatmulPreference_t pref);
```

This function initializes a matrix multiply heuristic search preferences descriptor in a previously allocated one.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>pref</code>		Output	Pointer to the structure holding the matrix multiply preferences descriptor created by this function. See <code>cublasLtMatrixLayout_t</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.31. `cublasLtMatmulPreferenceDestroy()`

```
cublasStatus_t cublasLtMatmulPreferenceDestroy(
    cublasLtMatmulPreference_t pref);
```

This function destroys a previously created matrix multiply preferences descriptor object.

**Parameters:**

Parameter	Memory	Input / Output	Description
pref		Input	Pointer to the structure holding the matrix multiply preferences descriptor that should be destroyed by this function. See <code>cublasLtMatmulPreference_t</code> .

Returns:

Return Value	Description
CUBLAS_STATUS_SUCCESS	If the operation was successful.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.32. `cublasLtMatmulPreferenceGetAttribute()`

```
cublasStatus_t cublasLtMatmulPreferenceGetAttribute(  
    cublasLtMatmulPreference_t pref,  
    cublasLtMatmulPreferenceAttributes_t attr,  
    void *buf,  
    size_t sizeInBytes,  
    size_t *sizeWritten);
```

This function returns the value of the queried attribute belonging to a previously created matrix multiply heuristic search preferences descriptor.

Parameters:

Parameter	Memory	Input / Output	Description
pref		Input	Pointer to the previously created structure holding the matrix multiply heuristic search preferences descriptor queried by this function. See <code>cublasLtMatmulPreference_t</code> .
attr		Input	The attribute that will be queried by this function. See <code>cublasLtMatmulPreferenceAttributes_t</code> .
buf		Output	Memory address containing the attribute value retrieved by this function.
sizeInBytes		Input	Size of <code>buf</code> buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is

Parameter	Memory	Input / Output	Description
			CUBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"> <li>▶ If sizeInBytes is 0 and sizeWritten is NULL, or</li> <li>▶ if sizeInBytes is non-zero and buf is NULL, or</li> <li>▶ sizeInBytes doesn't match size of internal storage for the selected attribute</li> </ul>
CUBLAS_STATUS_SUCCESS	If attribute's value was successfully written to user memory.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.33. `cublasLtMatmulPreferenceSetAttribute()`

```
cublasStatus_t cublasLtMatmulPreferenceSetAttribute(
    cublasLtMatmulPreference_t pref,
    cublasLtMatmulPreferenceAttributes_t attr,
    const void *buf,
    size_t sizeInBytes);
```

This function sets the value of the specified attribute belonging to a previously created matrix multiply preferences descriptor.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>pref</code>		Input	Pointer to the previously created structure holding the matrix multiply preferences descriptor queried by this function. See <code>cublasLtMatmulPreference_t</code> .
<code>attr</code>		Input	The attribute that will be set by this

Parameter	Memory	Input / Output	Description
buf		Input	function. See <code>cublasLtMatmulPreferenceAttributes</code> . The value to which the specified attribute should be set.
sizeInBytes		Input	Size of <code>buf</code> buffer (in bytes) for verification.

Returns:

Return Value	Description
<code>CUBLAS_STATUS_INVALID_VALUE</code>	If <code>buf</code> is NULL or <code>sizeInBytes</code> doesn't match the size of the internal storage for the selected attribute.
<code>CUBLAS_STATUS_SUCCESS</code>	If the attribute was set successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.34. `cublasLtMatrixLayoutCreate()`

```
cublasStatus_t cublasLtMatrixLayoutCreate( cublasLtMatrixLayout_t *matLayout,
                                           cudaDataType type,
                                           uint64_t rows,
                                           uint64_t cols,
                                           int64_t ld);
```

This function creates a matrix layout descriptor by allocating the memory needed to hold its opaque structure.

Parameters:

Parameter	Memory	Input / Output	Description
matLayout		Output	Pointer to the structure holding the matrix layout descriptor created by this function. See <code>cublasLtMatrixLayout_t</code> .
type		Input	Enumerant that specifies the data precision for the matrix layout descriptor this function creates. See <code>cudaDataType</code> .
rows, cols		Input	Number of rows and columns of the matrix.
ld		Input	The leading dimension of the matrix. In column major layout, this is the



Parameter	Memory	Input / Output	Description
			number of elements to jump to reach the next column. Thus $ld \geq m$ (number of rows).

**Returns:**

Return Value	Description
CUBLAS_STATUS_ALLOC_FAILED	If the memory could not be allocated.
CUBLAS_STATUS_SUCCESS	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.35. `cublasLtMatrixLayoutInit()`

```
cublasStatus_t cublasLtMatrixLayoutInit( cublasLtMatrixLayout_t matLayout,
                                         cudaDataType type,
                                         uint64_t rows,
                                         uint64_t cols,
                                         int64_t ld);
```

This function initializes a matrix layout descriptor in a previously allocated one.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>matLayout</code>		Output	Pointer to the structure holding the matrix layout descriptor initialized by this function. See <code>cublasLtMatrixLayout_t</code> .
<code>type</code>		Input	Enumerant that specifies the data precision for the matrix layout descriptor this function initializes. See <code>cudaDataType</code> .
<code>rows, cols</code>		Input	Number of rows and columns of the matrix.
<code>ld</code>		Input	The leading dimension of the matrix. In column major layout, this is the number of elements to jump to reach the next column. Thus $ld \geq m$ (number of rows).

**Returns:**

Return Value	Description
CUBLAS_STATUS_ALLOC_FAILED	If the memory could not be allocated.
CUBLAS_STATUS_SUCCESS	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.36. `cublasLtMatrixLayoutDestroy()`

```
cublasStatus_t cublasLtMatrixLayoutDestroy(
    cublasLtMatrixLayout_t matLayout);
```

This function destroys a previously created matrix layout descriptor object.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>matLayout</code>		Input	Pointer to the structure holding the matrix layout descriptor that should be destroyed by this function. See <code>cublasLtMatrixLayout_t</code> .

**Returns:**

Return Value	Description
CUBLAS_STATUS_SUCCESS	If the operation was successful.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.37. `cublasLtMatrixLayoutGetAttribute()`

```
cublasStatus_t cublasLtMatrixLayoutGetAttribute(
    cublasLtMatrixLayout_t matLayout,
    cublasLtMatrixLayoutAttribute_t attr,
    void *buf,
    size_t sizeInBytes,
    size_t *sizeWritten);
```

This function returns the value of the queried attribute belonging to the specified matrix layout descriptor.

**Parameters:**

Parameter	Memory	Input / Output	Description
<code>matLayout</code>		Input	Pointer to the previously created structure holding the matrix layout descriptor queried

Parameter	Memory	Input / Output	Description
attr		Input	by this function. See <code>cublasLtMatrixLayout_t</code> . The attribute being queried for. See <code>cublasLtMatrixLayoutAttribute_t</code> .
buf		Output	The attribute value returned by this function.
sizeInBytes		Input	Size of <code>buf</code> buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is <code>CUBLAS_STATUS_SUCCESS</code> . If <code>sizeInBytes</code> is non-zero: then <code>sizeWritten</code> is the number of bytes actually written; if <code>sizeInBytes</code> is 0: then <code>sizeWritten</code> is the number of bytes needed to write full contents.

Returns:

Return Value	Description
<code>CUBLAS_STATUS_INVALID_VALUE</code>	<ul style="list-style-type: none"><li>▶ If <code>sizeInBytes</code> is 0 and <code>sizeWritten</code> is NULL, or</li><li>▶ if <code>sizeInBytes</code> is non-zero and <code>buf</code> is NULL, or</li><li>▶ <code>sizeInBytes</code> doesn't match size of internal storage for the selected attribute</li></ul>
<code>CUBLAS_STATUS_SUCCESS</code>	If attribute's value was successfully written to user memory.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.38. `cublasLtMatrixLayoutSetAttribute()`

```
cublasStatus_t cublasLtMatrixLayoutSetAttribute(  
    cublasLtMatrixLayout_t matLayout,  
    cublasLtMatrixLayoutAttribute_t attr,  
    const void *buf,  
    size_t sizeInBytes);
```

This function sets the value of the specified attribute belonging to a previously created matrix layout descriptor.

Parameters:

Parameter	Memory	Input / Output	Description
matLayout		Input	Pointer to the previously created structure holding the matrix layout descriptor queried by this function. See <code>cublasLtMatrixLayout_t</code> .
attr		Input	The attribute that will be set by this function. See <code>cublasLtMatrixLayoutAttribute_t</code> .
buf		Input	The value to which the specified attribute should be set.
sizeInBytes		Input	Size of buf, the attribute buffer.

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If buf is NULL or sizeInBytes doesn't match size of internal storage for the selected attribute.
CUBLAS_STATUS_SUCCESS	If attribute was set successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.39. `cublasLtMatrixTransform()`

```

cublasStatus_t cublasLtMatrixTransform(
    cublasLtHandle_t lightHandle,
    cublasLtMatrixTransformDesc_t transformDesc,
    const void *alpha,
    const void *A,
    cublasLtMatrixLayout_t Adesc,
    const void *beta,
    const void *B,
    cublasLtMatrixLayout_t Bdesc,
    void *C,
    cublasLtMatrixLayout_t Cdesc,
    cudaStream_t stream);

```

This function computes the matrix transformation operation on the input matrices A and B, to produce the output matrix C, according to the below operation:

$$C = \text{alpha} * \text{transformation}(A) + \text{beta} * \text{transformation}(B),$$

where A, B are input matrices, and alpha and beta are input scalars. The transformation operation is defined by the `transformDesc` pointer. This function can be used to change the memory order of data or to scale and shift the values.

**Parameters:**

Parameter	Memory	Input / Output	Description
lightHandle		Input	Pointer to the allocated cuBLASLt handle for the cuBLASLt context. See <code>cublasLtHandle_t</code> .
transformDesc		Input	Pointer to the opaque descriptor holding the matrix transformation operation. See <code>cublasLtMatrixTransformDesc_t</code> .
alpha, beta	Device or host	Input	Pointers to the scalars used in the multiplication.
A, B, and C	Device	Input	Pointers to the GPU memory associated with the corresponding descriptors Adesc, Bdesc and Cdesc.
Adesc, Bdesc and Cdesc.		Input	Handles to the previous created descriptors of the type <code>cublasLtMatrixLayout_t</code> .  Adesc or Bdesc can be NULL if corresponding pointer is NULL and corresponding scalar is zero.
stream	Host	Input	The CUDA stream where all the GPU work will be submitted.

**Returns:**

Return Value	Description
CUBLAS_STATUS_NOT_INITIALIZED	If cuBLASLt handle has not been initialized.
CUBLAS_STATUS_INVALID_VALUE	If the parameters are in conflict or in an impossible configuration. For example, when A is not NULL, but Adesc is NULL.
CUBLAS_STATUS_NOT_SUPPORTED	If the current implementation on the selected device does not support the configured operation.
CUBLAS_STATUS_ARCH_MISMATCH	If the configured operation cannot be run using the selected device.
CUBLAS_STATUS_EXECUTION_FAILED	If CUDA reported an execution error from the device.
CUBLAS_STATUS_SUCCESS	If the operation completed successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.40. `cublasLtMatrixTransformDescCreate()`

```
cublasStatus_t cublasLtMatrixTransformDescCreate(
    cublasLtMatrixTransformDesc_t *transformDesc,
    cudaDataType scaleType);
```

This function creates a matrix transform descriptor by allocating the memory needed to hold its opaque structure.

#### Parameters:

Parameter	Memory	Input / Output	Description
<code>transformDesc</code>		Output	Pointer to the structure holding the matrix transform descriptor created by this function. See <code>cublasLtMatrixTransformDesc_t</code> .
<code>scaleType</code>		Input	Enumerant that specifies the data precision for the matrix transform descriptor this function creates. See <code>cudaDataType</code> .

#### Returns:

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.41. `cublasLtMatrixTransformDescInit()`

```
cublasStatus_t cublasLtMatrixTransformDescInit(
    cublasLtMatrixTransformDesc_t transformDesc,
    cudaDataType scaleType);
```

This function initializes a matrix transform descriptor in a previously allocated one.

#### Parameters:

Parameter	Memory	Input / Output	Description
<code>transformDesc</code>		Output	Pointer to the structure holding the matrix transform descriptor initialized by this function. See <code>cublasLtMatrixTransformDesc_t</code> .

Parameter	Memory	Input / Output	Description
scaleType		Input	Enumerant that specifies the data precision for the matrix transform descriptor this function initializes. See <code>cudaDataType</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	If memory could not be allocated.
<code>CUBLAS_STATUS_SUCCESS</code>	If the descriptor was created successfully.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.42. `cublasLtMatrixTransformDescDestroy()`

```
cublasStatus_t cublasLtMatrixTransformDescDestroy(
    cublasLtMatrixTransformDesc_t transformDesc);
```

This function destroys a previously created matrix transform descriptor object.

**Parameters:**

Parameter	Memory	Input / Output	Description
transformDesc		Input	Pointer to the structure holding the matrix transform descriptor that should be destroyed by this function. See <code>cublasLtMatrixTransformDesc_t</code> .

**Returns:**

Return Value	Description
<code>CUBLAS_STATUS_SUCCESS</code>	If the operation was successful.

See `cublasStatus_t` for a complete list of valid return codes.

### 3.4.43. `cublasLtMatrixTransformDescGetAttribute()`

```
cublasStatus_t cublasLtMatrixTransformDescGetAttribute(
    cublasLtMatrixTransformDesc_t transformDesc,
    cublasLtMatrixTransformDescAttributes_t attr,
    void *buf,
    size_t sizeInBytes,
    size_t *sizeWritten);
```

This function returns the value of the queried attribute belonging to a previously created matrix transform descriptor.

Parameters:

Parameter	Memory	Input / Output	Description
transformDesc		Input	Pointer to the previously created structure holding the matrix transform descriptor queried by this function. See cublasLtMatrixTransformDesc_t.
attr		Input	The attribute that will be retrieved by this function. See cublasLtMatrixTransformDescAttribute_t.
buf		Output	Memory address containing the attribute value retrieved by this function.
sizeInBytes		Input	Size of buf buffer (in bytes) for verification.
sizeWritten		Output	Valid only when the return value is CUBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

Returns:

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	<ul style="list-style-type: none"><li>▶ If sizeInBytes is 0 and sizeWritten is NULL, or</li><li>▶ if sizeInBytes is non-zero and buf is NULL, or</li><li>▶ sizeInBytes doesn't match size of internal storage for the selected attribute</li></ul>
CUBLAS_STATUS_SUCCESS	If attribute's value was successfully written to user memory.

See cublasStatus\_t for a complete list of valid return codes.



### 3.4.44. cublasLtMatrixTransformDescSetAttribute()

```
cublasStatus_t cublasLtMatrixTransformDescSetAttribute(  
    cublasLtMatrixTransformDesc_t transformDesc,  
    cublasLtMatrixTransformDescAttributes_t attr,  
    const void *buf,  
    size_t sizeInBytes);
```

This function sets the value of the specified attribute belonging to a previously created matrix transform descriptor.

**Parameters:**

Parameter	Memory	Input / Output	Description
transformDesc		Input	Pointer to the previously created structure holding the matrix transform descriptor queried by this function. See cublasLtMatrixTransformDesc_t.
attr		Input	The attribute that will be set by this function. See cublasLtMatrixTransformDescAttributes_t.
buf		Input	The value to which the specified attribute should be set.
sizeInBytes		Input	Size of buf buffer (in bytes) for verification.

**Returns:**

Return Value	Description
CUBLAS_STATUS_INVALID_VALUE	If buf is NULL or sizeInBytes does not match size of the internal storage for the selected attribute.
CUBLAS_STATUS_SUCCESS	If the attribute was set successfully.

See cublasStatus\_t for a complete list of valid return codes.

---

## Chapter 4. Using the cuBLASXt API

### 4.1. cuBLASXt General Description

The cuBLASXt API of cuBLAS exposes a multi-GPU capable Host interface: when using this API the application only needs to allocate the required matrices on the Host memory space. There are no restriction on the sizes of the matrices as long as they can fit into the Host memory. The cuBLASXt API takes care of allocating the memory across the designated GPUs and dispatched the workload between them and finally retrieves the results back to the Host. The cuBLASXt API supports only the compute-intensive BLAS3 routines (for example, matrix-matrix operations) where the PCI transfers back and forth from the GPU can be amortized. The cuBLASXt API has its own header file `cublasxt.h`.

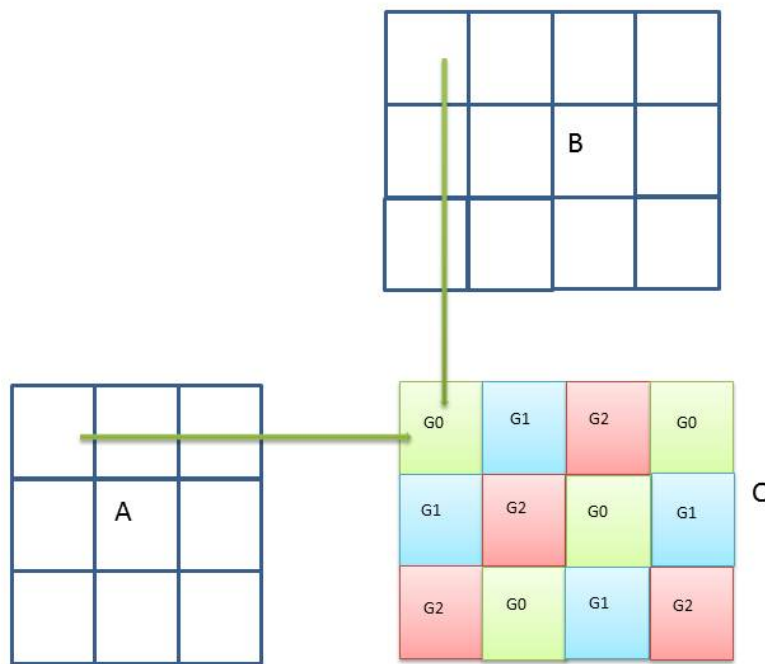
Starting with release 8.0, cuBLASXt API allows any of the matrices to be located on a GPU device.



**Note:** The cuBLASXt API is only supported on 64-bit platforms.

#### 4.1.1. Tiling design approach

To be able to share the workload between multiples GPUs, the cuBLASXt API uses a tiling strategy: every matrix is divided in square tiles of user-controllable dimension `BlockDim` x `BlockDim`. The resulting matrix tiling defines the static scheduling policy: each resulting tile is affected to a GPU in a round robin fashion. One CPU thread is created per GPU and is responsible to do the proper memory transfers and cuBLAS operations to compute all the tiles that it is responsible for. From a performance point of view, due to this static scheduling strategy, it is better that compute capabilities and PCI bandwidth are the same for every GPU. The figure below illustrates the tiles distribution between 3 GPUs. To compute the first tile `G0` from `C`, the CPU thread 0 responsible of GPU0, have to load 3 tiles from the first row of `A` and tiles from the first column of `B` in a pipeline fashion in order to overlap memory transfer and computations and sum the results into the first tile `G0` of `C` before to move on to the next tile `G0`.

Figure 1. Example of `cublasXt<t>gemm()` tiling for 3 Gpus

When the tile dimension is not an exact multiple of the dimensions of C, some tiles are partially filled on the right border and/or the bottom border. The current implementation does not pad the incomplete tiles but simply keeps track of those incomplete tiles by doing the right reduced cuBLAS operations; this way, no extra computation is done. However it still can lead to some load unbalance when all GPUS do not have the same number of incomplete tiles to work on.

When one or more matrices are located on some GPU devices, the same tiling approach and workload sharing is applied. The memory transfers are in this case done between devices. However, when the computation of a tile and some data are located on the same GPU device, the memory transfer to/from the local data into tiles is bypassed and the GPU operates directly on the local data. This can lead to a significant performance increase, especially when only one GPU is used for the computation.

The matrices can be located on any GPU device, and do not have to be located on the same GPU device. Furthermore, the matrices can even be located on a GPU device that do not participate to the computation.

Unlike the cuBLAS API, even if all matrices are located on the same device, the cuBLASXt API is still a blocking API from the Host point of view: the data results wherever located will be valid on the call return and no device synchronization is required.

## 4.1.2. Hybrid CPU-GPU Computation

In the case of very large problems, the cuBLASXt API offers the possibility to offload some of the computation to the Host CPU. This feature can be set up with the routines `cublasXtSetCpuRoutine()` and `cublasXtSetCpuRatio()`. The workload affected to the CPU is put aside; it is simply a percentage of the resulting matrix taken from the bottom and the right side whichever dimension is bigger. The GPU tiling is done after that on the reduced resulting matrix.

If any of the matrices is located on a GPU device, the feature is ignored, and all computations will be done only on the GPUs.

This feature should be used with caution because it could interfere with the CPU threads responsible of feeding the GPUs.

Currently, only the routine `cublasXt<t>gemm()` supports this feature.

## 4.1.3. Results Reproducibility

Currently all cuBLASXt API routines from a given toolkit version, generate the same bit-wise results when the following conditions are respected:

- ▶ all GPUs participating to the computation have the same compute-capabilities and the same number of SMs.
- ▶ the tiles size is kept the same between runs.
- ▶ either the CPU hybrid computation is not used or the CPU Blas provided is also guaranteed to produce reproducible results.

# 4.2. cuBLASXt API Datatypes Reference

## 4.2.1. cublasXtHandle\_t

The `cublasXtHandle_t` type is a pointer type to an opaque structure holding the cuBLASXt API context. The cuBLASXt API context must be initialized using `cublasXtCreate()` and the returned handle must be passed to all subsequent cuBLASXt API function calls. The context should be destroyed at the end using `cublasXtDestroy()`.

## 4.2.2. cublasXtOpType\_t

The `cublasOpType_t` enumerates the four possible types supported by BLAS routines. This enum is used as parameters of the routines `cublasXtSetCpuRoutine` and `cublasXtSetCpuRatio` to set up the hybrid configuration.

Value	Meaning
CUBLASXT_FLOAT	Float or single precision type.

Value	Meaning
CUBLASXT_DOUBLE	Double precision type.
CUBLASXT_COMPLEX	Single precision complex.
CUBLASXT_DOUBLECOMPLEX	Double precision complex.

### 4.2.3. cublasXtBlasOp\_t

The `cublasXtBlasOp_t` type enumerates the BLAS3 or BLAS-like routine supported by cuBLASxt API. This enum is used as parameters of the routines `cublasXtSetCpuRoutine` and `cublasXtSetCpuRatio` to set up the hybrid configuration.

Value	Meaning
CUBLASXT_GEMM	GEMM routine
CUBLASXT_SYRK	SYRK routine
CUBLASXT_HERK	HERK routine
CUBLASXT_SYMM	SYMM routine
CUBLASXT_HEMM	HEMM routine
CUBLASXT_TRSM	TRSM routine
CUBLASXT_SYR2K	SYR2K routine
CUBLASXT_HER2K	HER2K routine
CUBLASXT_SPM	SPMM routine
CUBLASXT_SYRKX	SYRKX routine
CUBLASXT_HERKX	HERKX routine

### 4.2.4. cublasXtPinningMemMode\_t

This type is used to enable or disable the Pinning Memory mode through the routine `cublasMgSetPinningMemMode`.

Value	Meaning
CUBLASXT_PINNING_DISABLED	The Pinning Memory mode is disabled.
CUBLASXT_PINNING_ENABLED	The Pinning Memory mode is enabled.

## 4.3. cuBLASXt API Helper Function Reference

### 4.3.1. cublasXtCreate()

```
cublasStatus_t
cublasXtCreate(cublasXtHandle_t *handle)
```

This function initializes the cuBLASXt API and creates a handle to an opaque structure holding the cuBLASXt API context. It allocates hardware resources on the host and device and must be called prior to making any other cuBLASXt API calls.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The initialization succeeded.
CUBLAS_STATUS_ALLOC_FAILED	The resources could not be allocated.
CUBLAS_STATUS_NOT_SUPPORTED	cuBLASXt API is only supported on 64-bit platform.

### 4.3.2. cublasXtDestroy()

```
cublasStatus_t
cublasXtDestroy(cublasXtHandle_t handle)
```

This function releases hardware resources used by the cuBLASXt API context. The release of GPU resources may be deferred until the application exits. This function is usually the last call with a particular handle to the cuBLASXt API.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The shut down succeeded.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.

### 4.3.3. cublasXtDeviceSelect()

```
cublasXtDeviceSelect(cublasXtHandle_t handle, int nbDevices, int deviceId[])
```

This function allows the user to provide the number of GPU devices and their respective IDs that will participate to the subsequent cuBLASXt API Math function calls. This function will create a cuBLAS context for every GPU provided in that list. Currently the device configuration is static and cannot be changed between Math function calls. In that regard, this function should be called only once after `cublasXtCreate`. To be able to run multiple configurations, multiple cuBLASXt API contexts should be created.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	User call was successful.
CUBLAS_STATUS_INVALID_VALUE	Access to at least one of the devices could not be done or a cuBLAS context could not be created on at least one of the devices.
CUBLAS_STATUS_ALLOC_FAILED	Some resources could not be allocated.

#### 4.3.4. cublasXtSetBlockDim()

```
cublasXtSetBlockDim(cublasXtHandle_t handle, int blockDim)
```

This function allows the user to set the block dimension used for the tiling of the matrices for the subsequent Math function calls. Matrices are split in square tiles of blockDim x blockDim dimension. This function can be called any time and will take effect for the following math function calls. The block dimension should be chosen in a way to optimize the math operation and to make sure that the PCI transfers are well-overlapped with the computation.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call has been successful.
CUBLAS_STATUS_INVALID_VALUE	blockDim <= 0

#### 4.3.5. cublasXtGetBlockDim()

```
cublasXtGetBlockDim(cublasXtHandle_t handle, int *blockDim)
```

This function allows the user to query the block dimension used for the tiling of the matrices.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call was successful.

#### 4.3.6. cublasXtSetCpuRoutine()

```
cublasXtSetCpuRoutine(cublasXtHandle_t handle, cublasXtBlasOp_t blasOp,
cublasXtOpType_t type, void *blasFuncPtr)
```

This function allows the user to provide a CPU implementation of the corresponding BLAS routine. This function can be used with the function cublasXtSetCpuRatio() to define a hybrid computation between the CPU and the GPUs. Currently the hybrid feature is only supported for the xGEMM routines.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call was successful.
CUBLAS_STATUS_INVALID_VALUE	blasOp or type define an invalid combination.

Return Value	Meaning
CUBLAS_STATUS_NOT_SUPPORTED	CPU-GPU Hybridization for that routine is not supported.

### 4.3.7. `cublasXtSetCpuRatio()`

```
cublasXtSetCpuRatio(cublasXtHandle_t handle, cublasXtBlasOp_t blasOp,
cublasXtOpType_t type, float ratio )
```

This function allows the user to define the percentage of workload that should be done on a CPU in the context of a hybrid computation. This function can be used with the function `cublasXtSetCpuRoutine()` to define a hybrid computation between the CPU and the GPUs. Currently the hybrid feature is only supported for the xGEMM routines.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call was successful.
CUBLAS_STATUS_INVALID_VALUE	<code>blasOp</code> or <code>type</code> define an invalid combination.
CUBLAS_STATUS_NOT_SUPPORTED	CPU-GPU Hybridization for that routine is not supported.

### 4.3.8. `cublasXtSetPinningMemMode()`

```
cublasXtSetPinningMemMode(cublasXtHandle_t handle, cublasXtPinningMemMode_t mode)
```

This function allows the user to enable or disable the Pinning Memory mode. When enabled, the matrices passed in subsequent cuBLASXt API calls will be pinned/unpinned using the CUDART routine `cudaHostRegister` and `cudaHostUnregister` respectively if the matrices are not already pinned. If a matrix happened to be pinned partially, it will also not be pinned. Pinning the memory improve PCI transfer performance and allows overlapping PCI memory transfer with computation. However pinning/unpinning the memory take some time, which might not be amortized. It is advised that the user pins the memory on its own using `cudaMallocHost` or `cudaHostRegister` and unpin it when the computation sequence is completed. By default, the Pinning Memory mode is disabled.



**Note:** The Pinning Memory mode should not be enabled when matrices used for different calls to cuBLASXt API overlap. cuBLASXt determines that a matrix is pinned or not if the first address of that matrix is pinned using `cudaHostGetFlags`, thus cannot know if the matrix is already partially pinned or not. This is especially true in multithreaded application where memory could be partially or totally pinned or unpinned while another thread is accessing that memory.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call was successful.



Return Value	Meaning
CUBLAS_STATUS_INVALID_VALUE	The mode value is different from CUBLASXT_PINNING_DISABLED and CUBLASXT_PINNING_ENABLED.

### 4.3.9. cublasXtGetPinningMemMode()

```
cublasXtGetPinningMemMode(cublasXtHandle_t handle, cublasXtPinningMemMode_t *mode)
```

This function allows the user to query the Pinning Memory mode. By default, the Pinning Memory mode is disabled.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	The call was successful.

## 4.4. cuBLASXt API Math Functions Reference

This section describes the actual linear algebra routines that cuBLASXt API supports. We will use abbreviations *<type>* for type and *<t>* for the corresponding short type to make a more concise and clear presentation of the implemented functions. Unless otherwise specified, *<type>* and *<t>* have the following meanings:

<type>	<t>	Meaning
float	's' or 'S'	Real single-precision.
double	'd' or 'D'	Real double-precision.
cuComplex	'c' or 'C'	Complex single-precision.
cuDoubleComplex	'z' or 'Z'	Complex double-precision.

The abbreviation **Re(.)** and **Im(.)** stand for the real and imaginary part of a number, respectively. Since imaginary part of a real number does not exist, we will consider it to be zero, and can usually discard it from the equation where it is being used. Also, the  $\bar{\alpha}$  will denote the complex conjugate of  $\alpha$ .

In general throughout the documentation, the lowercase Greek symbols  $\alpha$  and  $\beta$  will denote scalars, lowercase English letters in bold type **x** and **y** will denote vectors and capital English letters **A**, **B** and **C** will denote matrices.

### 4.4.1. cublasXt<t>gemm()

```
cublasStatus_t cublasXtSgemm(cublasXtHandle_t handle,
                             cublasOperation_t transa, cublasOperation_t transb,
                             size_t m, size_t n, size_t k,
                             const float *alpha,
                             const float *A, int lda,
```

```

        const float      *B, int ldb,
        const float      *beta,
        float            *C, int ldc)
cublasStatus_t cublasXtDgemm(cublasXtHandle_t handle,
        cublasOperation_t transa, cublasOperation_t transb,
        int m, int n, int k,
        const double     *alpha,
        const double     *A, int lda,
        const double     *B, int ldb,
        const double     *beta,
        double           *C, int ldc)
cublasStatus_t cublasXtCgemm(cublasXtHandle_t handle,
        cublasOperation_t transa, cublasOperation_t transb,
        int m, int n, int k,
        const cuComplex  *alpha,
        const cuComplex  *A, int lda,
        const cuComplex  *B, int ldb,
        const cuComplex  *beta,
        cuComplex        *C, int ldc)
cublasStatus_t cublasXtZgemm(cublasXtHandle_t handle,
        cublasOperation_t transa, cublasOperation_t transb,
        int m, int n, int k,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *B, int ldb,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc)

```

This function performs the matrix-matrix multiplication

$$C = \alpha \text{op}(A) \text{op}(B) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars, and  $A$ ,  $B$  and  $C$  are matrices stored in column-major format with dimensions  $\text{op}(A) m \times k$ ,  $\text{op}(B) k \times n$  and  $C m \times n$ , respectively. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

and  $\text{op}(B)$  is defined similarly for matrix  $B$ .

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASxT API context.
transa		input	Operation op(A) that is non- or (conj.) transpose.
transb		input	Operation op(B) that is non- or (conj.) transpose.
m		input	Number of rows of matrix op(A) and c.
n		input	Number of columns of matrix op(B) and c.
k		input	Number of columns of op(A) and rows of op(B).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimensions lda x k with lda>=max(1,m) if transa == CUBLAS_OP_N and lda x m with lda>=max(1,k) otherwise.

Param.	Memory	In/out	Meaning
lda		input	Leading dimension of two-dimensional array used to store the matrix A.
B	host or device	input	<type> array of dimension ldb x n with ldb>=max(1,k) if transb == CUBLAS_OP_N and ldb x k with ldb>=max(1,n) otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication. If beta==0, c does not have to be a valid input.
C	host or device	in/out	<type> array of dimensions ldc x n with ldc>=max(1,m).
ldc		input	Leading dimension of a two-dimensional array used to store the matrix c.

The possible error values returned by this function and their meanings are listed in the following table.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters m,n,k<0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[sgemm](#), [dgemm](#), [cgemm](#), [zgemm](#)

## 4.4.2. cublasXt<t>hemm()

```

cublasStatus_t cublasXtChemm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const cuComplex      *alpha,
                             const cuComplex      *A, size_t lda,
                             const cuComplex      *B, size_t ldb,
                             const cuComplex      *beta,
                             cuComplex            *C, size_t ldc)
cublasStatus_t cublasXtZhemm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, size_t lda,
                             const cuDoubleComplex *B, size_t ldb,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex *C, size_t ldc)

```

This function performs the Hermitian matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ \alpha BA + \beta C & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a Hermitian matrix stored in lower or upper mode,  $B$  and  $C$  are  $m \times n$  matrices, and  $\alpha$  and  $\beta$  are scalars.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
m		input	Number of rows of matrix C and B, with matrix A sized accordingly.
n		input	Number of columns of matrix C and B, with matrix A sized accordingly.
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension $lda \times m$ with $lda \geq \max(1, m)$ if <code>side==CUBLAS_SIDE_LEFT</code> and $lda \times n$ with $lda \geq \max(1, n)$ otherwise. The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if <code>beta==0</code> then C does not have to be a valid input.
C	host or device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $m, n < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[chemm](#), [zhemm](#)

### 4.4.3. cublasXt<t>symm()

```
cublasStatus_t cublasXtSsymm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const float          *alpha,
                             const float          *A, size_t lda,
                             const float          *B, size_t ldb,
                             const float          *beta,
                             float                *C, size_t ldc)
cublasStatus_t cublasXtDsymm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const double         *alpha,
                             const double         *A, size_t lda,
                             const double         *B, size_t ldb,
                             const double         *beta,
                             double               *C, size_t ldc)
cublasStatus_t cublasXtCsymm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const cuComplex      *alpha,
                             const cuComplex      *A, size_t lda,
                             const cuComplex      *B, size_t ldb,
                             const cuComplex      *beta,
                             cuComplex            *C, size_t ldc)
cublasStatus_t cublasXtZsymm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             size_t m, size_t n,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, size_t lda,
                             const cuDoubleComplex *B, size_t ldb,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex      *C, size_t ldc)
```

This function performs the symmetric matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ \alpha BA + \beta C & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a symmetric matrix stored in lower or upper mode,  $A$  and  $A$  are  $m \times n$  matrices, and  $\alpha$  and  $\beta$  are scalars.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
m		input	Number of rows of matrix A and B, with matrix A sized accordingly.

Param.	Memory	In/out	Meaning
n		input	Number of columns of matrix C and A, with matrix A sized accordingly.
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension $lda \times m$ with $lda \geq \max(1, m)$ if <code>side == CUBLAS_SIDE_LEFT</code> and $lda \times n$ with $lda \geq \max(1, n)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if <code>beta == 0</code> then C does not have to be a valid input.
C	host or device	in/out	<type> array of dimension $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $m, n < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssymm](#), [dsymm](#), [csymm](#), [zsymm](#)

#### 4.4.4. `cublasXt<t>syrk()`

```

cublasStatus_t cublasXtSsyrk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const float          *alpha,
                             const float          *A, int lda,
                             const float          *beta,
                             float                *C, int ldc)
cublasStatus_t cublasXtDsyrk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const double         *alpha,
                             const double         *A, int lda,
                             const double         *beta,

```

```

double *C, int ldc)
cublasStatus_t cublasXtCsyrk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const cuComplex *alpha,
                             const cuComplex *A, int lda,
                             const cuComplex *beta,
                             cuComplex *C, int ldc)
cublasStatus_t cublasXtZsyrk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, int lda,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex *C, int ldc)

```

This function performs the symmetric rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \ n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASxT API context.
uplo		input	Indicates if matrix c lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or transpose.
n		input	Number of rows of matrix op(A) and c.
k		input	Number of columns of matrix op(A).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension lda x k with lda>=max(1,n) if trans == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta	host	input	<type> scalar used for multiplication, if beta==0 then c does not have to be a valid input.
C	host or device	in/out	<type> array of dimension ldc x n, with ldc>=max(1,n).
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#)

### 4.4.5. `cublasXt<t>syr2k()`

```

cublasStatus_t cublasXtSsyr2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const float      *alpha,
                             const float      *A, size_t lda,
                             const float      *B, size_t ldb,
                             const float      *beta,
                             float            *C, size_t ldc)
cublasStatus_t cublasXtDsyr2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const double     *alpha,
                             const double     *A, size_t lda,
                             const double     *B, size_t ldb,
                             const double     *beta,
                             double           *C, size_t ldc)
cublasStatus_t cublasXtCsyr2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const cuComplex  *alpha,
                             const cuComplex  *A, size_t lda,
                             const cuComplex  *B, size_t ldb,
                             const cuComplex  *beta,
                             cuComplex        *C, size_t ldc)
cublasStatus_t cublasXtZsyr2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, size_t lda,
                             const cuDoubleComplex *B, size_t ldb,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex *C, size_t ldc)

```

This function performs the symmetric rank-  $2k$  update

$$C = \alpha(\text{op}(A)\text{op}(B)^T + \text{op}(B)\text{op}(A)^T) + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^T \text{ and } B^T & \text{if trans} == \text{CUBLAS\_OP\_T} \end{cases}$$



Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
uplo		input	Indicates if matrix c lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1, n)$ if $transa == CUBLAS\_OP\_N$ and $lda \times n$ with $lda \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimensions $ldb \times k$ with $ldb \geq \max(1, n)$ if $transb == CUBLAS\_OP\_N$ and $ldb \times n$ with $ldb \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if $beta == 0$ , then c does not have to be a valid input.
C	host or device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, n)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyr2k](#), [dsyr2k](#), [csyr2k](#), [zsyr2k](#)

## 4.4.6. `cublasXt<t>syrkx()`

```

cublasStatus_t cublasXtSsyrkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const float      *alpha,
                              const float      *A, size_t lda,
                              const float      *B, size_t ldb,
                              const float      *beta,
                              float            *C, size_t ldc)
cublasStatus_t cublasXtDsyrkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const double     *alpha,
                              const double     *A, size_t lda,
                              const double     *B, size_t ldb,
                              const double     *beta,
                              double           *C, size_t ldc)
cublasStatus_t cublasXtCsyrkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const cuComplex  *alpha,
                              const cuComplex  *A, size_t lda,
                              const cuComplex  *B, size_t ldb,
                              const cuComplex  *beta,
                              cuComplex        *C, size_t ldc)
cublasStatus_t cublasXtZsyrkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const cuDoubleComplex *alpha,
                              const cuDoubleComplex *A, size_t lda,
                              const cuDoubleComplex *B, size_t ldb,
                              const cuDoubleComplex *beta,
                              cuDoubleComplex *C, size_t ldc)

```

This function performs a variation of the symmetric rank-  $k$  update

$$C = \alpha(\text{op}(A)\text{op}(B)^T + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a symmetric matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^T \text{ and } B^T & \text{if trans} == \text{CUBLAS\_OP\_T} \end{cases}$$

This routine can be used when  $B$  is in such way that the result is guaranteed to be symmetric. An usual example is when the matrix  $B$  is a scaled form of the matrix  $A$ : this is equivalent to  $B$  being the product of the matrix  $A$  and a diagonal matrix.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
uplo		input	Indicates if matrix $c$ lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	Operation $\text{op}(A)$ that is non- or transpose.
n		input	Number of rows of matrix $\text{op}(A)$ , $\text{op}(B)$ and $c$ .
k		input	Number of columns of matrix $\text{op}(A)$ and $\text{op}(B)$ .

Param.	Memory	In/out	Meaning
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1, n)$ if $transa == \text{CUBLAS\_OP\_N}$ and $lda \times n$ with $lda \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimensions $ldb \times k$ with $ldb \geq \max(1, n)$ if $transb == \text{CUBLAS\_OP\_N}$ and $ldb \times n$ with $ldb \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if $\beta = 0$ , then C does not have to be a valid input.
C	host or device	in/out	<type> array of dimensions $ldc \times n$ with $ldc \geq \max(1, n)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed in the following table:

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#) and  
[ssyr2k](#), [dsyr2k](#), [csyr2k](#), [zsyr2k](#)

## 4.4.7. cublasXt<t>herk()

```

cublasStatus_t cublasXtCherk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const float *alpha,
                             const cuComplex *A, int lda,
                             const float *beta,
                             cuComplex *C, int ldc)
cublasStatus_t cublasXtZherk(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             int n, int k,
                             const double *alpha,

```

```
const cuDoubleComplex *A, int lda,
const double *beta,
cuDoubleComplex *C, int ldc)
```

This function performs the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A) \text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  is a matrix with dimensions  $\text{op}(A) \, n \times k$ . Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A) and c.
k		input	Number of columns of matrix op(A).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension lda x k with lda>=max(1,n) if transa == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
beta	host	input	<type> scalar used for multiplication, if beta==0 then c does not have to be a valid input.
C	host or device	in/out	<type> array of dimension ldc x n, with ldc>=max(1,n). The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix c.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#), [zherk](#)

## 4.4.8. cublasXt<t>her2k()

```
cublasStatus_t cublasXtCher2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const cuComplex *alpha,
                             const cuComplex *A, size_t lda,
                             const cuComplex *B, size_t ldb,
                             const float *beta,
                             cuComplex *C, size_t ldc)
cublasStatus_t cublasXtZher2k(cublasXtHandle_t handle,
                             cublasFillMode_t uplo, cublasOperation_t trans,
                             size_t n, size_t k,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, size_t lda,
                             const cuDoubleComplex *B, size_t ldb,
                             const double *beta,
                             cuDoubleComplex *C, size_t ldc)
```

This function performs the Hermitian rank-  $2k$  update

$$C = \alpha \text{op}(A)\text{op}(B)^H + \bar{\alpha} \text{op}(B)\text{op}(A)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^H \text{ and } B^H & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension lda x k with lda>=max(1,n) if transa == CUBLAS_OP_N and lda x n with lda>=max(1,k) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.

Param.	Memory	In/out	Meaning
B	host or device	input	<type> array of dimension $ldb \times k$ with $ldb \geq \max(1, n)$ if $transb == CUBLAS\_OP\_N$ and $ldb \times n$ with $ldb \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if $\beta == 0$ then c does not have to be a valid input.
C	host or device	in/out	<type> array of dimension $ldc \times n$ , with $ldc \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cher2k](#), [zher2k](#)

## 4.4.9. cublasXt<t>herkx()

```

cublasStatus_t cublasXtCherkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const cuComplex *alpha,
                              const cuComplex *A, size_t lda,
                              const cuComplex *B, size_t ldb,
                              const float *beta,
                              cuComplex *C, size_t ldc)
cublasStatus_t cublasXtZherkx(cublasXtHandle_t handle,
                              cublasFillMode_t uplo, cublasOperation_t trans,
                              size_t n, size_t k,
                              const cuDoubleComplex *alpha,
                              const cuDoubleComplex *A, size_t lda,
                              const cuDoubleComplex *B, size_t ldb,
                              const double *beta,
                              cuDoubleComplex *C, size_t ldc)

```

This function performs a variation of the Hermitian rank-  $k$  update

$$C = \alpha \text{op}(A)\text{op}(B)^H + \beta C$$

where  $\alpha$  and  $\beta$  are scalars,  $C$  is a Hermitian matrix stored in lower or upper mode, and  $A$  and  $B$  are matrices with dimensions  $\text{op}(A) \ n \times k$  and  $\text{op}(B) \ n \times k$ , respectively. Also, for matrix  $A$  and  $B$

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if trans} == \text{CUBLAS\_OP\_N} \\ A^H \text{ and } B^H & \text{if trans} == \text{CUBLAS\_OP\_C} \end{cases}$$

This routine can be used when the matrix B is in such way that the result is guaranteed to be hermitian. An usual example is when the matrix B is a scaled form of the matrix A: this is equivalent to B being the product of the matrix A and a diagonal matrix. For an efficient computation of the product of a regular matrix with a diagonal matrix, refer to the routine `cusblasXt<t>dgmm`.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASxt API context.
uplo		input	Indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
n		input	Number of rows of matrix op(A), op(B) and c.
k		input	Number of columns of matrix op(A) and op(B).
alpha	host	input	<type> scalar used for multiplication.
A	host or device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1, n)$ if $transa == \text{CUBLAS\_OP\_N}$ and $lda \times n$ with $lda \geq \max(1, k)$ otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimension $ldb \times k$ with $ldb \geq \max(1, n)$ if $transb == \text{CUBLAS\_OP\_N}$ and $ldb \times n$ with $ldb \geq \max(1, k)$ otherwise.
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	Real scalar used for multiplication, if $\beta=0$ then c does not have to be a valid input.
C	host or device	in/out	<type> array of dimension $ldc \times n$ , with $ldc \geq \max(1, n)$ . The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $n, k < 0$ .
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[cherk](#), [zherk](#) and

[cher2k](#), [zher2k](#)

## 4.4.10. cublasXt<t>trsm()

```

cublasStatus_t cublasXtStrsm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             cublasOperation_t trans, cublasXtDiagType_t diag,
                             size_t m, size_t n,
                             const float *alpha,
                             const float *A, size_t lda,
                             float *B, size_t ldb)
cublasStatus_t cublasXtDtrsm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             cublasOperation_t trans, cublasXtDiagType_t diag,
                             size_t m, size_t n,
                             const double *alpha,
                             const double *A, size_t lda,
                             double *B, size_t ldb)
cublasStatus_t cublasXtCtrsm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             cublasOperation_t trans, cublasXtDiagType_t diag,
                             size_t m, size_t n,
                             const cuComplex *alpha,
                             const cuComplex *A, size_t lda,
                             cuComplex *B, size_t ldb)
cublasStatus_t cublasXtZtrsm(cublasXtHandle_t handle,
                             cublasSideMode_t side, cublasFillMode_t uplo,
                             cublasOperation_t trans, cublasXtDiagType_t diag,
                             size_t m, size_t n,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *A, size_t lda,
                             cuDoubleComplex *B, size_t ldb)

```

This function solves the triangular linear system with multiple right-hand-sides

$$\begin{cases} \text{op}(A)X = \alpha B & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ X\text{op}(A) = \alpha B & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a triangular matrix stored in lower or upper mode with or without the main diagonal,  $X$  and  $B$  are  $m \times n$  matrices, and  $\alpha$  is a scalar. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

The solution  $X$  overwrites the right-hand-sides  $B$  on exit.



No test for singularity or near-singularity is included in this function.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
side		input	Indicates if matrix A is on the left or right of x.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	Number of rows of matrix B, with matrix A sized accordingly.
n		input	Number of columns of matrix B, with matrix A is sized accordingly.
alpha	host	input	<type> scalar used for multiplication, if alpha==0 then A is not referenced and B does not have to be a valid input.
A	host or device	input	<type> array of dimension lda x m with lda>=max(1,m) if side == CUBLAS_SIDE_LEFT and lda x n with lda>=max(1,n) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	in/out	<type> array. It has dimensions ldb x n with ldb>=max(1,m).
ldb		input	Leading dimension of two-dimensional array used to store matrix B.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters m, n<0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strsm](#), [dtrsm](#), [ctrsm](#), [ztrsm](#)

## 4.4.11. `cublasXt<t>trmm()`

```
cublasStatus_t cublasXtStrmm(cublasXtHandle_t handle,
```

```

        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        size_t m, size_t n,
        const float          *alpha,
        const float          *A, size_t lda,
        const float          *B, size_t ldb,
        float                *C, size_t ldc)
cublasStatus_t cublasXtDtrmm(cublasXtHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        size_t m, size_t n,
        const double         *alpha,
        const double         *A, size_t lda,
        const double         *B, size_t ldb,
        double               *C, size_t ldc)
cublasStatus_t cublasXtCtrmm(cublasXtHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        size_t m, size_t n,
        const cuComplex      *alpha,
        const cuComplex      *A, size_t lda,
        const cuComplex      *B, size_t ldb,
        cuComplex            *C, size_t ldc)
cublasStatus_t cublasXtZtrmm(cublasXtHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        size_t m, size_t n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, size_t lda,
        const cuDoubleComplex *B, size_t ldb,
        cuDoubleComplex       *C, size_t ldc)

```

This function performs the triangular matrix-matrix multiplication

$$C = \begin{cases} \alpha \text{op}(A)B & \text{if side} == \text{CUBLAS\_SIDE\_LEFT} \\ \alpha B \text{op}(A) & \text{if side} == \text{CUBLAS\_SIDE\_RIGHT} \end{cases}$$

where  $A$  is a triangular matrix stored in lower or upper mode with or without the main diagonal,  $B$  and  $C$  are  $m \times n$  matrix, and  $\alpha$  is a scalar. Also, for matrix  $A$

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS\_OP\_N} \\ A^T & \text{if transa} == \text{CUBLAS\_OP\_T} \\ A^H & \text{if transa} == \text{CUBLAS\_OP\_C} \end{cases}$$

Note that in order to achieve better parallelism, similar to the cuBLAS API, cuBLASXt API differs from the BLAS API for this routine. The BLAS API assumes an in-place implementation (with results written back to B), while the cuBLASXt API assumes an out-of-place implementation (with results written into C). The application can still obtain the in-place functionality of BLAS in the cuBLASXt API by passing the address of the matrix B in place of the matrix C. No other overlapping in the input parameters is supported.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASXt API context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.

Param.	Memory	In/out	Meaning
trans		input	Operation op(A) that is non- or (conj.) transpose.
diag		input	Indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	Number of rows of matrix B, with matrix A sized accordingly.
n		input	Number of columns of matrix B, with matrix A sized accordingly.
alpha	host	input	<type> scalar used for multiplication, if alpha==0 then A is not referenced and B does not have to be a valid input.
A	host or device	input	<type> array of dimension lda x m with lda>=max(1,m) if side == CUBLAS_SIDE_LEFT and lda x n with lda>=max(1,n) otherwise.
lda		input	Leading dimension of two-dimensional array used to store matrix A.
B	host or device	input	<type> array of dimension ldb x n with ldb>=max(1,m).
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
C	host or device	in/out	<type> array of dimension ldc x n with ldc>=max(1,m).
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters m, n < 0.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[strmm](#), [dtrmm](#), [ctrmm](#), [ztrmm](#)

## 4.4.12. cublasXt<t>spmm()

```
cublasStatus_t cublasXtSspmm( cublasXtHandle_t handle,
                               cublasSideMode_t side,
                               cublasFillMode_t uplo,
                               size_t m,
                               size_t n,
                               const float *alpha,
                               const float *AP,
                               const float *B,
```

```

        size_t ldb,
        const float *beta,
        float *C,
        size_t ldc );

cublasStatus_t cublasXtDspmm( cublasXtHandle_t handle,
                             cublasSideMode_t side,
                             cublasFillMode_t uplo,
                             size_t m,
                             size_t n,
                             const double *alpha,
                             const double *AP,
                             const double *B,
                             size_t ldb,
                             const double *beta,
                             double *C,
                             size_t ldc );

cublasStatus_t cublasXtCspmm( cublasXtHandle_t handle,
                             cublasSideMode_t side,
                             cublasFillMode_t uplo,
                             size_t m,
                             size_t n,
                             const cuComplex *alpha,
                             const cuComplex *AP,
                             const cuComplex *B,
                             size_t ldb,
                             const cuComplex *beta,
                             cuComplex *C,
                             size_t ldc );

cublasStatus_t cublasXtZspmm( cublasXtHandle_t handle,
                             cublasSideMode_t side,
                             cublasFillMode_t uplo,
                             size_t m,
                             size_t n,
                             const cuDoubleComplex *alpha,
                             const cuDoubleComplex *AP,
                             const cuDoubleComplex *B,
                             size_t ldb,
                             const cuDoubleComplex *beta,
                             cuDoubleComplex *C,
                             size_t ldc );

```

This function performs the symmetric packed matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side == CUBLAS\_SIDE\_LEFT} \\ \alpha BA + \beta C & \text{if side == CUBLAS\_SIDE\_RIGHT} \end{cases}$$


where  $A$  is a  $n \times n$  symmetric matrix stored in packed format,  $B$  and  $C$  are  $m \times n$  matrices, and  $\alpha$  and  $\beta$  are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the element  $A(i, j)$  is stored in the memory location `AP[i + ((2*n-j+1)*j)/2]` for  $j = 1, \dots, n$  and  $i \geq j$ .

Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix  $A$  are packed together column by column without gaps, so that the

element  $A(i, j)$  is stored in the memory location  $AP[i + (j * (j + 1)) / 2]$  for  $j = 1, \dots, n$  and  $i \leq j$ . Consequently, the packed format requires only  $\frac{n(n+1)}{2}$  elements for storage.

 **Note:** The packed matrix AP must be located on the host whereas the other matrices can be located on the host or any GPU device.

Param.	Memory	In/out	Meaning
handle		input	Handle to the cuBLASx API context.
side		input	Indicates if matrix A is on the left or right of B.
uplo		input	Indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
m		input	Number of rows of matrix A and B, with matrix A sized accordingly.
n		input	Number of columns of matrix C and A, with matrix A sized accordingly.
alpha	host	input	<type> scalar used for multiplication.
AP	host	input	<type> array with A stored in packed format.
B	host or device	input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1, m)$ .
ldb		input	Leading dimension of two-dimensional array used to store matrix B.
beta	host	input	<type> scalar used for multiplication, if $\beta == 0$ then C does not have to be a valid input.
C	host or device	in/out	<type> array of dimension $ldc \times n$ with $ldc \geq \max(1, m)$ .
ldc		input	Leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Value	Meaning
CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	The library was not initialized.
CUBLAS_STATUS_INVALID_VALUE	The parameters $m, n < 0$ .
CUBLAS_STATUS_NOT_SUPPORTED	The matrix AP is located on a GPU device.
CUBLAS_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

For references please refer to:

[ssymm](#), [dsymm](#), [csymm](#), [zsymm](#)

---

# Appendix A. Using the cuBLAS Legacy API

This section does not provide a full reference of each Legacy API datatype and entry point. Instead, it describes how to use the API, especially where this is different from the regular cuBLAS API.

Note that in this section, all references to the “cuBLAS Library” refer to the Legacy cuBLAS API only.



**WARNING: The legacy cuBLAS API is deprecated and will be removed in future release.**

## A.1. Error Status

The `cublasStatus` type is used for function status returns. The cuBLAS Library helper functions return status directly, while the status of core functions can be retrieved using `cublasGetError()`. Note that reading the error status via `cublasGetError()` resets the internal error state to `CUBLAS_STATUS_SUCCESS`. Currently, the following values for are defined as:

Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	The library was not initialized.
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	The resource allocation failed.
<code>CUBLAS_STATUS_INVALID_VALUE</code>	An invalid numerical value was used as an argument.
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	An absent device architectural feature is required.
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	An access to GPU memory space failed.
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	The GPU program failed to execute.
<code>CUBLAS_STATUS_INTERNAL_ERROR</code>	An internal operation failed.
<code>CUBLAS_STATUS_NOT_SUPPORTED</code>	The feature required is not supported.

This legacy type corresponds to type `cublasStatus_t` in the cuBLAS library API.

## A.2. Initialization and Shutdown

The functions `cublasInit()` and `cublasShutdown()` are used to initialize and shut down the cuBLAS library. It is recommended for `cublasInit()` to be called before any other function is invoked. It allocates hardware resources on the GPU device that is currently bound to the host thread from which it was invoked.

The legacy initialization and shutdown functions are similar to the cuBLAS library API routines `cublasCreate()` and `cublasDestroy()`.

## A.3. Thread Safety

The legacy API is not thread-safe when used with multiple host threads and devices. It is recommended to be used only when utmost compatibility with Fortran is required and when a single host thread is used to setup the library and make all the functions calls.

## A.4. Memory Management

The memory used by the legacy cuBLAS library API is allocated and released using functions `cublasAlloc()` and `cublasFree()`, respectively. These functions create and destroy an object in the GPU memory space capable of holding an array of `n` elements, where each element requires `elemSize` bytes of storage. Please see the legacy cuBLAS API header file `cublas.h` for the prototypes of these functions.

The function `cublasAlloc()` is a wrapper around the function `cudaMalloc()`, therefore device pointers returned by `cublasAlloc()` can be passed to any CUDA device kernel functions. However, these device pointers cannot be dereferenced in the host code. The function `cublasFree()` is a wrapper around the function `cudaFree()`.

## A.5. Scalar Parameters

In the legacy cuBLAS API, scalar parameters are passed by value from the host. Also, the few functions that do return a scalar result, such as `dot()` and `norm2()`, return the resulting value on the host, and hence these routines will wait for kernel execution on the device to complete before returning, which makes parallelism with streams impractical. However, the majority of functions do not return any value, in order to be more compatible with Fortran and the existing BLAS libraries.

## A.6. Helper Functions

This section lists the helper functions provided by the legacy cuBLAS API and their functionality. For the exact prototypes of these functions, refer to the legacy cuBLAS API header file `cublas.h`.

Helper function	Meaning
<code>cublasInit()</code>	Initialize the library.
<code>cublasShutdown()</code>	Shuts down the library.
<code>cublasGetError()</code>	Retrieves the error status of the library.
<code>cublasSetKernelStream()</code>	Sets the stream to be used by the library.
<code>cublasAlloc()</code>	Allocates the device memory for the library.
<code>cublasFree()</code>	Releases the device memory allocated for the library.
<code>cublasSetVector()</code>	Copies a vector $x$ on the host to a vector on the GPU.
<code>cublasGetVector()</code>	Copies a vector $x$ on the GPU to a vector on the host.
<code>cublasSetMatrix()</code>	Copies a $m \times n$ tile from a matrix on the host to the GPU.
<code>cublasGetMatrix()</code>	Copies a $m \times n$ tile from a matrix on the GPU to the host.
<code>cublasSetVectorAsync()</code>	Similar to <code>cublasSetVector()</code> , but the copy is asynchronous.
<code>cublasGetVectorAsync()</code>	Similar to <code>cublasGetVector()</code> , but the copy is asynchronous.
<code>cublasSetMatrixAsync()</code>	Similar to <code>cublasSetMatrix()</code> , but the copy is asynchronous.
<code>cublasGetMatrixAsync()</code>	Similar to <code>cublasGetMatrix()</code> , but the copy is asynchronous.

## A.7. Level-1,2,3 Functions

The Level-1,2,3 cuBLAS functions (also called core functions) have the same name and behavior as the ones listed in the sections 3, 4 and 5 in this document. Please refer to the legacy cuBLAS API header file `cublas.h` for their exact prototype. Also, the next section talks



a bit more about the differences between the legacy and the cuBLAS API prototypes, more specifically how to convert the function calls from one API to another.

## A.8. Converting Legacy to the cuBLAS API

There are a few general rules that can be used to convert from legacy to the cuBLAS API:

- ▶ Exchange the header file `cublas.h` for `cublas_v2.h`.
- ▶ Exchange the type `cublasStatus` for `cublasStatus_t`.
- ▶ Exchange the function `cublasSetKernelStream()` for `cublasSetStream()`.
- ▶ Exchange the function `cublasAlloc()` and `cublasFree()` for `cudaMalloc()` and `cudaFree()`, respectively. Note that `cudaMalloc()` expects the size of the allocated memory to be provided in bytes (usually simply provide `n × elemSize` to allocate `n` elements, each of size `elemSize` bytes).
- ▶ Declare the `cublasHandle_t` cuBLAS library handle.
- ▶ Initialize the handle using `cublasCreate()`. Also, release the handle once finished using `cublasDestroy()`.
- ▶ Add the handle as the first parameter to all the cuBLAS library function calls.
- ▶ Change the scalar parameters to be passed by reference, instead of by value (usually, simply adding the “&” symbol in C/C++ is enough, because the parameters are passed by reference on the host by *default*). However, note that if the routine is running asynchronously, then the variable holding the scalar parameter cannot be changed until the kernels that the routine dispatches are completed. See the [CUDA C++ Programming Guide](#) for a detailed discussion of how to use streams.
- ▶ Change the parameter characters 'N' or 'n' (non-transpose operation), 'T' or 't' (transpose operation) and 'C' or 'c' (conjugate transpose operation) to `CUBLAS_OP_N`, `CUBLAS_OP_T` and `CUBLAS_OP_C`, respectively.
- ▶ Change the parameter characters 'L' or 'l' (lower part filled) and 'U' or 'u' (upper part filled) to `CUBLAS_FILL_MODE_LOWER` and `CUBLAS_FILL_MODE_UPPER`, respectively.
- ▶ Change the parameter characters 'N' or 'n' (non-unit diagonal) and 'U' or 'u' (unit diagonal) to `CUBLAS_DIAG_NON_UNIT` and `CUBLAS_DIAG_UNIT`, respectively.
- ▶ Change the parameter characters 'L' or 'l' (left side) and 'R' or 'r' (right side) to `CUBLAS_SIDE_LEFT` and `CUBLAS_SIDE_RIGHT`, respectively.
- ▶ If the legacy API function returns a scalar value, add an extra scalar parameter of the same type passed by reference, as the last parameter to the same function.
- ▶ Instead of using `cublasGetError`, use the return value of the function itself to check for errors.
- ▶ Finally, use the function prototypes in the header files `cublas.h` and `cublas_v2.h` to check the code for correctness.

## A.9. Examples

For sample code references that use the legacy cuBLAS API please see the two examples below. They show an application written in C using the legacy cuBLAS library API with two indexing styles (Example A.1. "Application Using C and cuBLAS: 1-based indexing" and Example A.2. "Application Using C and cuBLAS: 0-based Indexing"). This application is analogous to the one using the cuBLAS library API that is shown in the Introduction chapter.

### Example A.1. Application Using C and cuBLAS: 1-based indexing

```
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "cublas.h"
#define M 6
#define N 5
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1))

static __inline__ void modify (float *m, int ldm, int n, int p, int q, float
alpha, float beta){
    cublasSscal (n-q+1, alpha, &m[IDX2F(p,q,ldm)], ldm);
    cublasSscal (ldm-p+1, beta, &m[IDX2F(p,q,ldm)], 1);
}

int main (void){
    int i, j;
    cublasStatus stat;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            a[IDX2F(i,j,M)] = (float)((i-1) * M + j);
        }
    }
    cublasInit();
    stat = cublasAlloc (M*N, sizeof(*a), (void**)&devPtrA);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("device memory allocation failed");
        cublasShutdown();
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    modify (devPtrA, M, N, 2, 3, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
}
```

```

    }
    cublasFree (devPtrA);
    cublasShutdown();
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            printf ("%7.0f", a[IDX2F(i,j,M)]);
        }
        printf ("\n");
    }
    free(a);
    return EXIT_SUCCESS;
}

```

### Example A.2. Application Using C and cuBLAS: 0-based indexing

```

//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "cublas.h"
#define M 6
#define N 5
#define IDX2C(i,j,ld) (((j)*(ld))+(i))

static __inline__ void modify (float *m, int ldm, int n, int p, int q, float
alpha, float beta){
    cublasSscal (n-q, alpha, &m[IDX2C(p,q,ldm)], ldm);
    cublasSscal (ldm-p, beta, &m[IDX2C(p,q,ldm)], 1);
}

int main (void){
    int i, j;
    cublasStatus stat;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            a[IDX2C(i,j,M)] = (float)(i * M + j + 1);
        }
    }
    cublasInit();
    stat = cublasAlloc (M*N, sizeof(*a), (void**)&devPtrA);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("device memory allocation failed");
        cublasShutdown();
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    modify (devPtrA, M, N, 1, 2, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != cuBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cublasFree (devPtrA);
    }
}

```

```
        cublasShutdown();  
        return EXIT_FAILURE;  
    }  
    cublasFree (devPtrA);  
    cublasShutdown();  
    for (j = 0; j < N; j++) {  
        for (i = 0; i < M; i++) {  
            printf ("%7.0f", a[IDX2C(i,j,M)]);  
        }  
        printf ("\n");  
    }  
    free(a);  
    return EXIT_SUCCESS;  
}
```

---

## Appendix B. cuBLAS Fortran Bindings

The cuBLAS library is implemented using the C-based CUDA toolchain. Thus, it provides a C-style API. This makes interfacing to applications written in C and C++ trivial, but the library can also be used by applications written in Fortran. In particular, the cuBLAS library uses 1-based indexing and Fortran-style column-major storage for multidimensional data to simplify interfacing to Fortran applications. Unfortunately, Fortran-to-C calling conventions are not standardized and differ by platform and toolchain. In particular, differences may exist in the following areas:

- ▶ symbol names (capitalization, name decoration)
- ▶ argument passing (by value or reference)
- ▶ passing of string arguments (length information)
- ▶ passing of pointer arguments (size of the pointer)
- ▶ returning floating-point or compound data types (for example single-precision or complex data types)

To provide maximum flexibility in addressing those differences, the cuBLAS Fortran interface is provided in the form of wrapper functions and is part of the Toolkit delivery. The C source code of those wrapper functions is located in the `src` directory and provided in two different forms:

- ▶ the thunking wrapper interface located in the file `fortran_thunking.c`
- ▶ the direct wrapper interface located in the file `fortran.c`

The code of one of those two files needs to be compiled into an application for it to call the cuBLAS API functions. Providing source code allows users to make any changes necessary for a particular platform and toolchain.

The code in those two C files has been used to demonstrate interoperability with the compilers g77 3.2.3 and g95 0.91 on 32-bit Linux, g77 3.4.5 and g95 0.91 on 64-bit Linux, Intel Fortran 9.0 and Intel Fortran 10.0 on 32-bit and 64-bit Microsoft Windows XP, and g77 3.4.0 and g95 0.92 on Mac OS X.

Note that for g77, use of the compiler flag `-fno-second-underscore` is required to use these wrappers as provided. Also, the use of the default calling conventions with regard to argument and return value passing is expected. Using the flag `-fno-f2c` changes the default calling convention with respect to these two items.

The thunking wrappers allow interfacing to existing Fortran applications without any changes to the application. During each call, the wrappers allocate GPU memory, copy source data from CPU memory space to GPU memory space, call cuBLAS, and finally copy back the results to CPU memory space and deallocate the GPU memory. As this process causes very significant call overhead, these wrappers are intended for light testing, not for production code. To use the thunking wrappers, the application needs to be compiled with the file `fortran_thunking.c`.

The direct wrappers, intended for production code, substitute device pointers for vector and matrix arguments in all BLAS functions. To use these interfaces, existing applications need to be modified slightly to allocate and deallocate data structures in GPU memory space (using `cuBLAS_ALLOC` and `cuBLAS_FREE`) and to copy data between GPU and CPU memory spaces (using `cuBLAS_SET_VECTOR`, `cuBLAS_GET_VECTOR`, `cuBLAS_SET_MATRIX`, and `cuBLAS_GET_MATRIX`). The sample wrappers provided in `fortran.c` map device pointers to the OS-dependent type `size_t`, which is 32-bit wide on 32-bit platforms and 64-bit wide on a 64-bit platforms.

One approach to deal with index arithmetic on device pointers in Fortran code is to use C-style macros, and use the C preprocessor to expand these, as shown in the example below. On Linux and Mac OS X, one way of pre-processing is to use the option `-E -x f77-cpp-input` when using the `g77` compiler, or simply the option `-cpp` when using `g95` or `gfortran`. On Windows platforms with Microsoft Visual C/C++, using `cl -EP` achieves similar results.

```
! Example B.1. Fortran 77 Application Executing on the Host
! -----
subroutine modify ( m, ldm, n, p, q, alpha, beta )
  implicit none
  integer ldm, n, p, q
  real*4 m (ldm, *) , alpha , beta
  external cublas_sscal
  call cublas_sscal (n-p+1, alpha , m(p,q), ldm)
  call cublas_sscal (ldm-p+1, beta, m(p,q), 1)
  return
end

program matrixmod
  implicit none
  integer M,N
  parameter (M=6, N=5)
  real*4 a(M,N)
  integer i, j
  external cublas_init
  external cublas_shutdown

  do j = 1, N
    do i = 1, M
      a(i, j) = (i-1)*M + j
    enddo
  enddo
  call cublas_init
  call modify ( a, M, N, 2, 3, 16.0, 12.0 )
  call cublas_shutdown
  do j = 1 , N
    do i = 1 , M
      write(*,"(F7.0$)") a(i,j)
    enddo
    write (*,*) ""
  enddo
  stop
```

```
end
```

When traditional fixed-form Fortran 77 code is ported to use the cuBLAS library, line length often increases when the BLAS calls are exchanged for cuBLAS calls. Longer function names and possible macro expansion are contributing factors. Inadvertently exceeding the maximum line length can lead to runtime errors that are difficult to find, so care should be taken not to exceed the 72-column limit if fixed-form is retained.

The examples in this chapter show a small application implemented in Fortran 77 on the host and the same application with the non-thunking wrappers after it has been ported to use the cuBLAS library.

The second example should be compiled with ARCH\_64 defined as 1 on 64-bit OS system and as 0 on 32-bit OS system. For example for g95 or gfortran, this can be done directly on the command line by using the option `-cpp -DARCH_64=1`.

```
! Example B.2. Same Application Using Non-thunking cuBLAS Calls
!-----
#define IDX2F (i,j,ld) (((j)-1)*(ld))+((i)-1))
subroutine modify ( devPtrM, ldm, n, p, q, alpha, beta )
  implicit none
  integer sizeof_real
  parameter (sizeof_real=4)
  integer ldm, n, p, q
#if ARCH_64
  integer*8 devPtrM
#else
  integer*4 devPtrM
#endif
  real*4 alpha, beta
  call cublas_sscal ( n-p+1, alpha,
    1 devPtrM+IDX2F(p, q, ldm)*sizeof_real,
    2 ldm)
  call cublas_sscal(ldm-p+1, beta,
    1 devPtrM+IDX2F(p, q, ldm)*sizeof_real,
    2 1)
  return
end
program matrixmod
  implicit none
  integer M,N,sizeof_real
#if ARCH_64
  integer*8 devPtrA
#else
  integer*4 devPtrA
#endif
  parameter(M=6,N=5,sizeof_real=4)
  real*4 a(M,N)
  integer i,j,stat
  external cublas_init, cublas_set_matrix, cublas_get_matrix
  external cublas_shutdown, cublas_alloc
  integer cublas_alloc, cublas_set_matrix, cublas_get_matrix
  do j=1,N
    do i=1,M
      a(i,j)=(i-1)*M+j
    enddo
  enddo
  call cublas_init
  stat= cublas_alloc(M*N, sizeof_real, devPtrA)
  if (stat.NE.0) then
    write(*,*) "device memory allocation failed"
    call cublas_shutdown
    stop
  end if
end program
```

```

endif
stat = cublas_set_matrix(M,N,sizeof_real,a,M,devPtrA,M)
if (stat.NE.0) then
  call cublas_free( devPtrA )
  write(*,*) "data download failed"
  call cublas_shutdown
  stop
endif

```

---

--- Code block continues below. Space added for formatting purposes. ---

---

```

call modify(devPtrA, M, N, 2, 3, 16.0, 12.0)
stat = cublas_get_matrix(M, N, sizeof_real, devPtrA, M, a, M )
if (stat.NE.0) then
  call cublas_free ( devPtrA )
  write(*,*) "data upload failed"
  call cublas_shutdown
  stop
endif
call cublas_free ( devPtrA )
call cublas_shutdown
do j = 1 , N
  do i = 1 , M
    write (*,"(F7.0$)") a(i,j)
  enddo
  write (*,*) ""
enddo
stop
end

```



---

## Appendix C. Interaction with Other Libraries and Tools

This section describes important requirements and recommendations that ensure correct use of cuBLAS with other libraries and utilities.

### C.1. `nvprune`

`nvprune` enables pruning relocatable host objects and static libraries to only contain device code for the specific target architectures. In case of cuBLAS, particular care must be taken if using `nvprune` with compute capabilities, whose minor revision number is different than 0. To reduce binary size, cuBLAS may only store major revision equivalents of CUDA binary files for kernels reused between different minor revision versions. Therefore, to ensure that a pruned library does not fail for arbitrary problems, the user must keep binaries for a selected architecture and all prior minor architectures in its major architecture.

For example, the following call prunes `libcublas_static.a` to contain only `sm_75` (Volta) and `sm_70` (Turing) cubins:

```
nvprune --generate-code code=sm_70 --generate-code code=sm_75 libcublasLt_static.a -o libcublasLt_static_sm70_sm75.a
```

which should be used instead of:

```
nvprune -arch=sm_75 libcublasLt_static.a -o libcublasLt_static_sm75.a
```

---

## Appendix D. Acknowledgements

NVIDIA would like to thank the following individuals and institutions for their contributions:

- ▶ Portions of the SGEMM, DGEMM, CGEMM and ZGEMM library routines were written by Vasily Volkov of the University of California.
- ▶ Portions of the SGEMM, DGEMM and ZGEMM library routines were written by Davide Barbieri of the University of Rome Tor Vergata.
- ▶ Portions of the DGEMM and SGEMM library routines optimized for Fermi architecture were developed by the University of Tennessee. Subsequently, several other routines that are optimized for the Fermi architecture have been derived from these initial DGEMM and SGEMM implementations.
- ▶ The substantial optimizations of the STRSV, DTRSV, CTRSV and ZTRSV library routines were developed by Jonathan Hogg of The Science and Technology Facilities Council (STFC). Subsequently, some optimizations of the STRSM, DTRSM, CTRSM and ZTRSM have been derived from these TRSV implementations.
- ▶ Substantial optimizations of the SYMV and HEMV library routines were developed by Ahmad Abdelfattah, David Keyes and Hatem Ltaief of King Abdullah University of Science and Technology (KAUST).
- ▶ Substantial optimizations of the TRMM and TRSM library routines were developed by Ali Charara, David Keyes and Hatem Ltaief of King Abdullah University of Science and Technology (KAUST).
- ▶ This product includes {fmt} - A modern formatting library <https://fmt.dev> Copyright (c) 2012 - present, Victor Zverovich.
- ▶ This product includes spdlog - Fast C++ logging library. <https://github.com/gabime/spdlog> The MIT License (MIT).
- ▶ This product includes SIMD Library for Evaluating Elementary Functions, vectorized libm and DFT <https://sleef.org> Boost Software License - Version 1.0 - August 17th, 2003.
- ▶ This product includes Frozen - a header-only, constexpr alternative to gperf for C++14 users. <https://github.com/serge-sans-paille/frozen> Apache License - Version 2.0, January 2004.
- ▶ This product includes Boost C++ Libraries - free peer-reviewed portable C++ source libraries <https://www.boost.org/> Boost Software License - Version 1.0 - August 17th, 2003.

## Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation ("NVIDIA") makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer ("Terms of Sale"). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

## OpenCL

OpenCL is a trademark of Apple Inc. used under license to the Khronos Group Inc.

## Trademarks

NVIDIA and the NVIDIA logo are trademarks or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

## Copyright

© 2012-2022 NVIDIA Corporation & affiliates. All rights reserved.