



cuDNN Developer Guide

Table of Contents

Chapter 1. Overview.....	1
Chapter 2. Programming Model.....	2
Chapter 3. Convolution Formulas.....	3
Chapter 4. Notation.....	5
Chapter 5. Tensor Descriptor.....	6
5.1. WXYZ Tensor Descriptor.....	6
5.2. 4-D Tensor Descriptor.....	6
5.3. 5-D Tensor Description.....	7
5.4. Fully-packed Tensors.....	7
5.5. Partially-packed Tensors.....	7
5.6. Spatially Packed Tensors.....	8
5.7. Overlapping Tensors.....	8
Chapter 6. Data Layout Formats.....	9
6.1. Example.....	9
6.2. NCHW Memory Layout.....	10
6.3. NHWC Memory Layout.....	11
6.4. NC/32HW32 Memory Layout.....	12
Chapter 7. Thread Safety.....	15
Chapter 8. Reproducibility (determinism).....	16
Chapter 9. Scaling Parameters.....	17
Chapter 10. Tensor Core Operations.....	18
10.1. Basics.....	18
10.2. Convolution Functions.....	18
10.2.1. Prerequisites.....	19
10.2.2. Supported Algorithms.....	19
10.2.3. Data And Filter Formats.....	19
10.3. RNN Functions.....	20
10.3.1. Prerequisites.....	20
10.3.2. Supported Algorithms.....	20
10.3.3. Data And Filter Formats.....	20
10.4. Tensor Transformations.....	21
10.4.1. FP16 Data.....	21
10.4.2. FP32-to-FP16 Conversion.....	21

10.4.3. Padding.....	22
10.4.4. Folding.....	22
10.4.5. Conversion Between NCHW And NHWC.....	22
10.5. Guidelines For A Deep Learning Compiler.....	23
Chapter 11. GPU And Driver Requirements.....	24
Chapter 12. Backward Compatibility And Deprecation Policy.....	25
Chapter 13. Grouped Convolutions.....	27
Chapter 14. API Logging.....	29
Chapter 15. Features Of RNN Functions.....	31
Chapter 16. Mixed Precision Numerical Accuracy.....	34
Chapter 17. Operation Fusion Via The Backend API.....	35
Chapter 18. Troubleshooting.....	38
18.1. FAQs.....	38
18.2. How Do I Report A Bug?.....	41
18.3. Support.....	41
Chapter 19. Acknowledgments.....	42
19.1. University of Tennessee.....	42
19.2. University of California, Berkeley.....	42
19.3. Facebook AI Research, New York.....	43

List of Figures

Figure 1. Example with N=1, C=64, H=5, W=4.	10
Figure 2. NCHW Memory Layout	11
Figure 3. NHWC Memory Layout	12
Figure 4. NC/32HW32 Memory Layout	13
Figure 5. Scaling Parameters for Convolution	17
Figure 6. INT8 for cudnnConvolutionBiasActivationForward	17
Figure 7. Tensor Operation with FP16 Inputs	21
Figure 8. Tensor Operation with FP32 Inputs	21
Figure 9. A set of operation descriptors the user passes to the operation graph	36
Figure 10. The operation graph after finalization	36
Figure 11. Software stack with cuDNN.	38

List of Tables

Table 1. Convolution terms	3
Table 2. Two-step, deprecation policy	25
Table 3. API Logging Using Environment Variables	30

Chapter 1. Overview

NVIDIA® CUDA® Deep Neural Network library™ (cuDNN) is a GPU-accelerated library of primitives for deep neural networks. It provides highly tuned implementations of routines arising frequently in DNN applications:

- ▶ Convolution forward and backward, including cross-correlation
- ▶ Pooling forward and backward
- ▶ Softmax forward and backward
- ▶ Neuron activations forward and backward:
 - ▶ Rectified linear (ReLU)
 - ▶ Sigmoid
 - ▶ Hyperbolic tangent (TANH)
- ▶ Tensor transformation functions
- ▶ LRN, LCN and batch normalization forward and backward

cuDNN convolution routines aim for a performance that is competitive with the fastest GEMM (matrix multiply)-based implementations of such routines while using significantly less memory.

cuDNN features include customizable data layouts, supporting flexible dimension ordering, striding, and subregions for the 4D tensors used as inputs and outputs to all of its routines. This flexibility allows easy integration into any neural network implementation and avoids the input/output transposition steps sometimes necessary with GEMM-based convolutions.

cuDNN offers a context-based API that allows for easy multithreading and (optional) interoperability with NVIDIA® CUDA® streams.

Chapter 2. Programming Model

The cuDNN library exposes a Host API but assumes that for operations using the GPU, the necessary data is directly accessible from the device.

An application using cuDNN must initialize a handle to the library context by calling [`cudaDnnCreate\(\)`](#). This handle is explicitly passed to every subsequent library function that operates on GPU data. Once the application finishes using cuDNN, it can release the resources associated with the library handle using [`cudaDnnDestroy\(\)`](#). This approach allows the user to explicitly control the library's functioning when using multiple host threads, GPUs and CUDA Streams.

For example, an application can use [`cudaSetDevice`](#) to associate different devices with different host threads, and in each of those host threads, use a unique cuDNN handle that directs the library calls to the device associated with it. Thus the cuDNN library calls made with different handles will automatically run on different devices.

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

cuDNN API Compatibility

Beginning in cuDNN 7, the binary compatibility of a patch and minor releases is maintained as follows:

- ▶ Any patch release `x.y.z` is forward or backward-compatible with applications built against another cuDNN patch release `x.y.w` (meaning, of the same major and minor version number, but having $w \neq z$).
- ▶ cuDNN minor releases beginning with cuDNN 7 are binary backward-compatible with applications built against the same or earlier patch release (meaning, an application built against cuDNN `7.x` is binary compatible with cuDNN library `7.y`, where $y \geq x$).
- ▶ Applications compiled with a cuDNN version `7.y` are not guaranteed to work with `7.x` release when $y > x$.

Chapter 3. Convolution Formulas

This section describes the various convolution formulas implemented in convolution functions.

The convolution terms described in the table below apply to all the convolution formulas that follow.

Table 1. Convolution terms

Term	Description
x	Input (image) Tensor
w	Weight Tensor
y	Output Tensor
n	Current Batch Size
c	Current Input Channel
C	Total Input Channels
H	Input Image Height
W	Input Image Width
k	Current Output Channel
K	Total Output Channels
p	Current Output Height Position
q	Current Output Width Position
G	Group Count
pad	Padding Value
u	Vertical Subsample Stride (along Height)
v	Horizontal Subsample Stride (along Width)
dil_h	Vertical Dilation (along Height)
dil_w	Horizontal Dilation (along Width)
r	Current Filter Height
R	Total Filter Height
s	Current Filter Width
S	Total Filter Width

Term	Description
C_g	$\frac{C}{G}$
K_g	$\frac{K}{G}$

Normal Convolution (using cross-correlation mode)

$$y_{n, k, p, q} = \sum_c^C \sum_r^R \sum_s^S x_{n, c, p+r, q+s} \times w_{k, c, r, s}$$

Convolution with Padding

$$x_{<0, <0} = 0$$

$$x_{>H, >W} = 0$$

$$y_{n, k, p, q} = \sum_c^C \sum_r^R \sum_s^S x_{n, c, p+r-pad, q+s-pad} \times w_{k, c, r, s}$$

Convolution with Subsample-Striding

$$y_{n, k, p, q} = \sum_c^C \sum_r^R \sum_s^S x_{n, c, (p*u)+r, (q*v)+s} \times w_{k, c, r, s}$$

Convolution with Dilation

$$y_{n, k, p, q} = \sum_c^C \sum_r^R \sum_s^S x_{n, c, p+(r*dilh), q+(s*dilw)} \times w_{k, c, r, s}$$

Convolution using Convolution Mode

$$y_{n, k, p, q} = \sum_c^C \sum_r^R \sum_s^S x_{n, c, p+r, q+s} \times w_{k, c, R-r-1, S-s-1}$$

Convolution using Grouped Convolution

$$C_g = \frac{C}{G}$$

$$K_g = \frac{K}{G}$$

$$y_{n, k, p, q} = \sum_c^{C_g} \sum_r^R \sum_s^S x_{n, C_g*\text{floor}(k/K_g)+c, p+r, q+s} \times w_{k, c, r, s}$$

Chapter 4. Notation

As of cuDNN version 4, we have adopted a mathematically-inspired notation for layer inputs and outputs using x, y, dx, dy, b, w for common layer parameters. This was done to improve readability and ease of understanding of the meaning of the parameters. All layers now follow a uniform convention as below:

During inference

```
y = layerFunction(x, otherParams)
```

During backpropagation

```
(dx, dOtherParams) = layerFunctionGradient(x, y, dy, otherParams)
```

During convolution

For **convolution**, the notation is:

```
y = x*w+b
```

where:

- ▶ w is the matrix of filter weights
- ▶ x is the previous layer's data (during inference)
- ▶ y is the next layer's data
- ▶ b is the bias and $*$ is the convolution operator

In backpropagation routines, the parameters keep their meanings.

The parameters dx, dy, dw, db always refer to the gradient of the final network error function with respect to a given parameter. So dy in all backpropagation routines always refers to error gradient backpropagated through the network computation graph so far. Similarly, other parameters in more specialized layers, such as, for instance, $dMeans$ or $dBnBias$ refer to gradients of the loss function with regard to those parameters.



Note: w is used in the API for both the width of the x tensor and convolution filter matrix. To resolve this ambiguity we use w and `filter` notation interchangeably for the convolution filter weight matrix. The meaning is clear from the context since the layer width is always referenced near its height.

Chapter 5. Tensor Descriptor

The cuDNN library describes data holding images, videos and any other data with contents with a generic n-D tensor defined with the following parameters:

- ▶ a dimension `nbDims` from 3 to 8
- ▶ a data type (32-bit floating-point, 64 bit-floating point, 16-bit floating-point...)
- ▶ `dimA` integer array defining the size of each dimension
- ▶ `strideA` integer array defining the stride of each dimension (for example, the number of elements to add to reach the next element from the same dimension)

The first dimension of the tensor defines the batch size `n`, and the second dimension defines the number of feature maps `c`. This tensor definition allows, for example, to have some dimensions overlapping each other within the same tensor by having the stride of one dimension smaller than the product of the dimension and the stride of the next dimension. In cuDNN, unless specified otherwise, all routines will support tensors with overlapping dimensions for forward-pass input tensors, however, dimensions of the output tensors cannot overlap. Even though this tensor format supports negative strides (which can be useful for data mirroring), cuDNN routines do not support tensors with negative strides unless specified otherwise.

5.1. WXYZ Tensor Descriptor

Tensor descriptor formats are identified using acronyms, with each letter referencing a corresponding dimension. In this document, the usage of this terminology implies:

- ▶ all the strides are strictly positive
- ▶ the dimensions referenced by the letters are sorted in decreasing order of their respective strides

5.2. 4-D Tensor Descriptor

A 4-D Tensor descriptor is used to define the format for batches of 2D images with 4 letters: `n, c, h, w` for respectively the batch size, the number of feature maps, the height and the width. The letters are sorted in decreasing order of the strides. The commonly used 4-D tensor formats are:

- ▶ NCHW
- ▶ NHWC
- ▶ CHWN

5.3. 5-D Tensor Description

A 5-D Tensor descriptor is used to define the format of the batch of 3D images with 5 letters: N, C, D, H, W for respectively the batch size, the number of feature maps, the depth, the height, and the width. The letters are sorted in decreasing order of the strides. The commonly used 5-D tensor formats are called:

- ▶ NCDHW
- ▶ NDHWC
- ▶ CDHWN

5.4. Fully-packed Tensors

A tensor is defined as `XYZ-fully-packed` if and only if:

- ▶ the number of tensor dimensions is equal to the number of letters preceding the `fully-packed` suffix.
- ▶ the stride of the i -th dimension is equal to the product of the $(i+1)$ -th dimension by the $(i+1)$ -th stride.
- ▶ the stride of the last dimension is 1.

5.5. Partially-packed Tensors

The partially `XYZ-packed` terminology only applies in the context of a tensor format described with a superset of the letters used to define a partially-packed tensor. A `wXYZ` tensor is defined as `XYZ-packed` if and only if:

- ▶ The strides of all dimensions NOT referenced in the `-packed` suffix are greater or equal to the product of the next dimension by the next stride.
- ▶ The stride of each dimension referenced in the `-packed` suffix in position i is equal to the product of the $(i+1)$ -st dimension by the $(i+1)$ -st stride.
- ▶ If the last tensor's dimension is present in the `-packed` suffix, its stride is 1.

For example, an `NHWC` tensor `WC-packed` means that the `c_stride` is equal to 1 and `w_stride` is equal to `c_dim × c_stride`. In practice, the `-packed` suffix is usually applied to the minor dimensions of a tensor but can be applied to only the major dimensions; for example, an `NCHW` tensor that is only `N-packed`.

5.6. Spatially Packed Tensors

Spatially-packed tensors are defined as partially-packed in spatial dimensions.

For example, a spatially-packed 4D tensor would mean that the tensor is either NCHW HW-packed or CNHW HW-packed.

5.7. Overlapping Tensors

A tensor is defined to be overlapping if iterating over a full range of dimensions produces the same address more than once.

In practice an overlapped tensor will have $\text{stride}[i-1] < \text{stride}[i] * \text{dim}[i]$ for some of the i from $[1, \text{nbDims}]$ interval.

Chapter 6. Data Layout Formats

This section describes how cuDNN tensors are arranged in memory. See [cudnnTensorFormat_t](#) for enumerated tensor format types.

6.1. Example

Consider a batch of images in 4D with the following dimensions:

- ▶ **N** is the batch size; 1.
- ▶ **C** is the number of feature maps (i.e., number of channels); 64.
- ▶ **H** is the image height; 5.
- ▶ **W** is the image width; 4.

To keep the example simple, the image pixel elements are expressed as a sequence of integers, 0, 1, 2, 3, and so on. See [Figure 1](#).

Figure 1. Example with N=1, C=64, H=5, W=4.

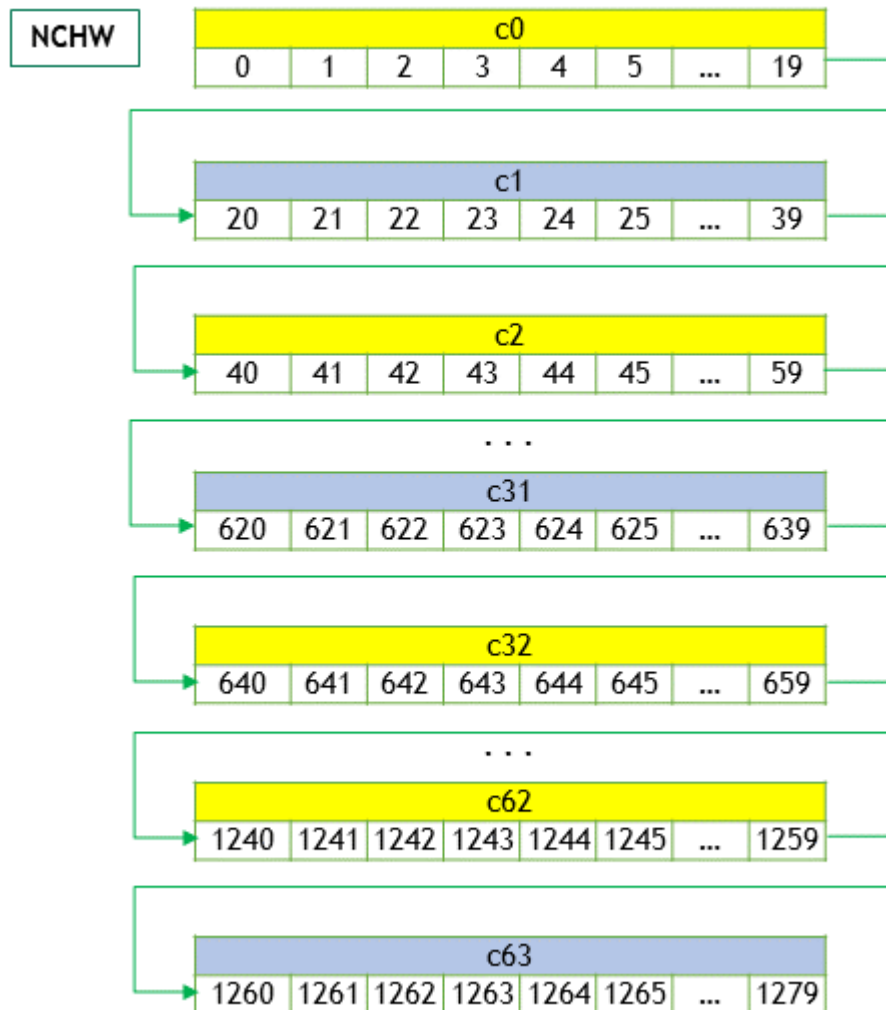
EXAMPLE N = 1 C = 64 H = 5 W = 4	c = 0	c = 1	c = 2																																																												
	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>0</td><td>1</td><td>2</td><td>3</td></tr><tr><td>4</td><td>5</td><td>6</td><td>7</td></tr><tr><td>8</td><td>9</td><td>10</td><td>11</td></tr><tr><td>12</td><td>13</td><td>14</td><td>15</td></tr><tr><td>16</td><td>17</td><td>18</td><td>19</td></tr></table>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>20</td><td>21</td><td>22</td><td>23</td></tr><tr><td>24</td><td>25</td><td>26</td><td>27</td></tr><tr><td>28</td><td>29</td><td>30</td><td>31</td></tr><tr><td>32</td><td>33</td><td>34</td><td>35</td></tr><tr><td>36</td><td>37</td><td>38</td><td>39</td></tr></table>	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>40</td><td>41</td><td>42</td><td>43</td></tr><tr><td>44</td><td>45</td><td>46</td><td>47</td></tr><tr><td>48</td><td>49</td><td>50</td><td>51</td></tr><tr><td>52</td><td>53</td><td>54</td><td>55</td></tr><tr><td>56</td><td>57</td><td>58</td><td>59</td></tr></table>	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59
	0	1	2	3																																																											
	4	5	6	7																																																											
	8	9	10	11																																																											
	12	13	14	15																																																											
16	17	18	19																																																												
20	21	22	23																																																												
24	25	26	27																																																												
28	29	30	31																																																												
32	33	34	35																																																												
36	37	38	39																																																												
40	41	42	43																																																												
44	45	46	47																																																												
48	49	50	51																																																												
52	53	54	55																																																												
56	57	58	59																																																												
...																																																															
c = 30	c = 31	c = 32																																																													
<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>600</td><td>601</td><td>602</td><td>603</td></tr><tr><td>604</td><td>605</td><td>606</td><td>607</td></tr><tr><td>608</td><td>609</td><td>610</td><td>611</td></tr><tr><td>612</td><td>613</td><td>614</td><td>615</td></tr><tr><td>616</td><td>617</td><td>618</td><td>619</td></tr></table>	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>620</td><td>621</td><td>622</td><td>623</td></tr><tr><td>624</td><td>625</td><td>626</td><td>627</td></tr><tr><td>628</td><td>629</td><td>630</td><td>631</td></tr><tr><td>632</td><td>633</td><td>634</td><td>635</td></tr><tr><td>636</td><td>637</td><td>638</td><td>639</td></tr></table>	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>640</td><td>641</td><td>642</td><td>643</td></tr><tr><td>644</td><td>645</td><td>646</td><td>647</td></tr><tr><td>648</td><td>649</td><td>650</td><td>651</td></tr><tr><td>652</td><td>653</td><td>654</td><td>655</td></tr><tr><td>656</td><td>657</td><td>658</td><td>659</td></tr></table>	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	
600	601	602	603																																																												
604	605	606	607																																																												
608	609	610	611																																																												
612	613	614	615																																																												
616	617	618	619																																																												
620	621	622	623																																																												
624	625	626	627																																																												
628	629	630	631																																																												
632	633	634	635																																																												
636	637	638	639																																																												
640	641	642	643																																																												
644	645	646	647																																																												
648	649	650	651																																																												
652	653	654	655																																																												
656	657	658	659																																																												
...																																																															
c = 62		c = 63																																																													
<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>1240</td><td>1241</td><td>1242</td><td>1243</td></tr><tr><td>1244</td><td>1245</td><td>1246</td><td>1247</td></tr><tr><td>1248</td><td>1249</td><td>1250</td><td>1251</td></tr><tr><td>1252</td><td>1253</td><td>1254</td><td>1255</td></tr><tr><td>1256</td><td>1257</td><td>1258</td><td>1259</td></tr></table>		1240	1241	1242	1243	1244	1245	1246	1247	1248	1249	1250	1251	1252	1253	1254	1255	1256	1257	1258	1259	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>1260</td><td>1261</td><td>1262</td><td>1263</td></tr><tr><td>1264</td><td>1265</td><td>1266</td><td>1267</td></tr><tr><td>1268</td><td>1269</td><td>1270</td><td>1271</td></tr><tr><td>1272</td><td>1273</td><td>1274</td><td>1275</td></tr><tr><td>1276</td><td>1277</td><td>1278</td><td>1279</td></tr></table>	1260	1261	1262	1263	1264	1265	1266	1267	1268	1269	1270	1271	1272	1273	1274	1275	1276	1277	1278	1279																					
1240	1241	1242	1243																																																												
1244	1245	1246	1247																																																												
1248	1249	1250	1251																																																												
1252	1253	1254	1255																																																												
1256	1257	1258	1259																																																												
1260	1261	1262	1263																																																												
1264	1265	1266	1267																																																												
1268	1269	1270	1271																																																												
1272	1273	1274	1275																																																												
1276	1277	1278	1279																																																												
...																																																															

6.2. NCHW Memory Layout

The above 4D Tensor is laid out in the memory in the NCHW format as below:

1. Beginning with the first channel (c=0), the elements are arranged contiguously in row-major order.
2. Continue with second and subsequent channels until the elements of all the channels are laid out. See [Figure 2](#).
3. Proceed to the next batch (if **N** is > 1).

Figure 2. NCHW Memory Layout

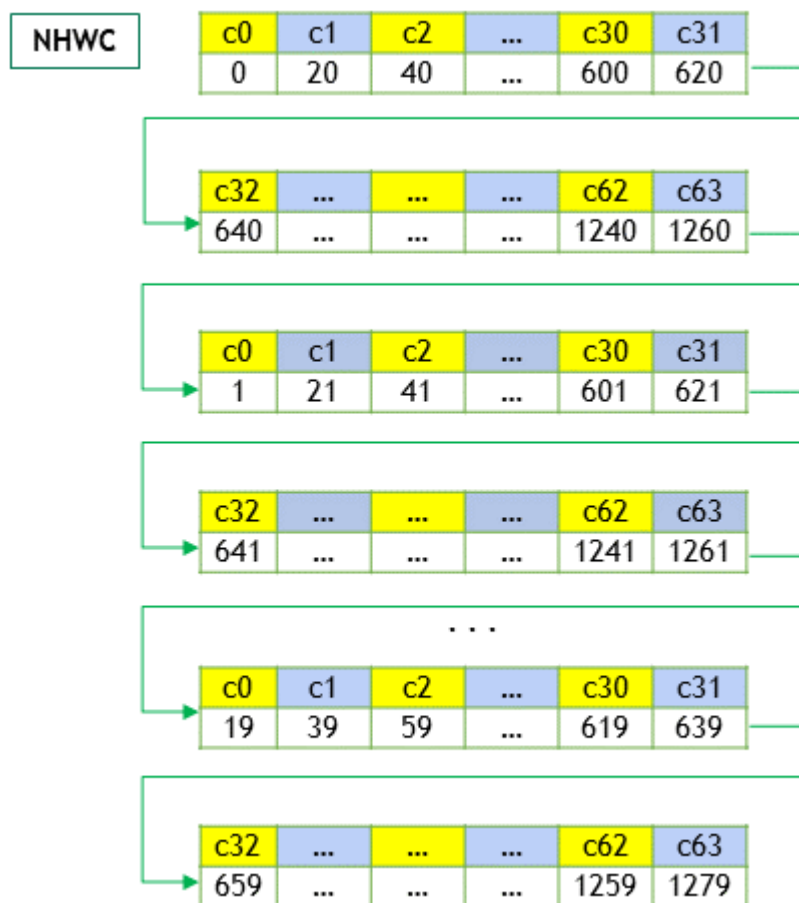


6.3. NHWC Memory Layout

For the NHWC memory layout, the corresponding elements in all the **C** channels are laid out first, as below:

1. Begin with the first element of channel 0, then proceed to the first element of channel 1, and so on, until the first elements of all the **C** channels are laid out.
2. Next, select the second element of channel 0, then proceed to the second element of channel 1, and so on, until the second element of all the channels are laid out.
3. Follow the row-major order of channel 0 and complete all the elements. See [Figure 3](#).
4. Proceed to the next batch (if **N** is > 1).

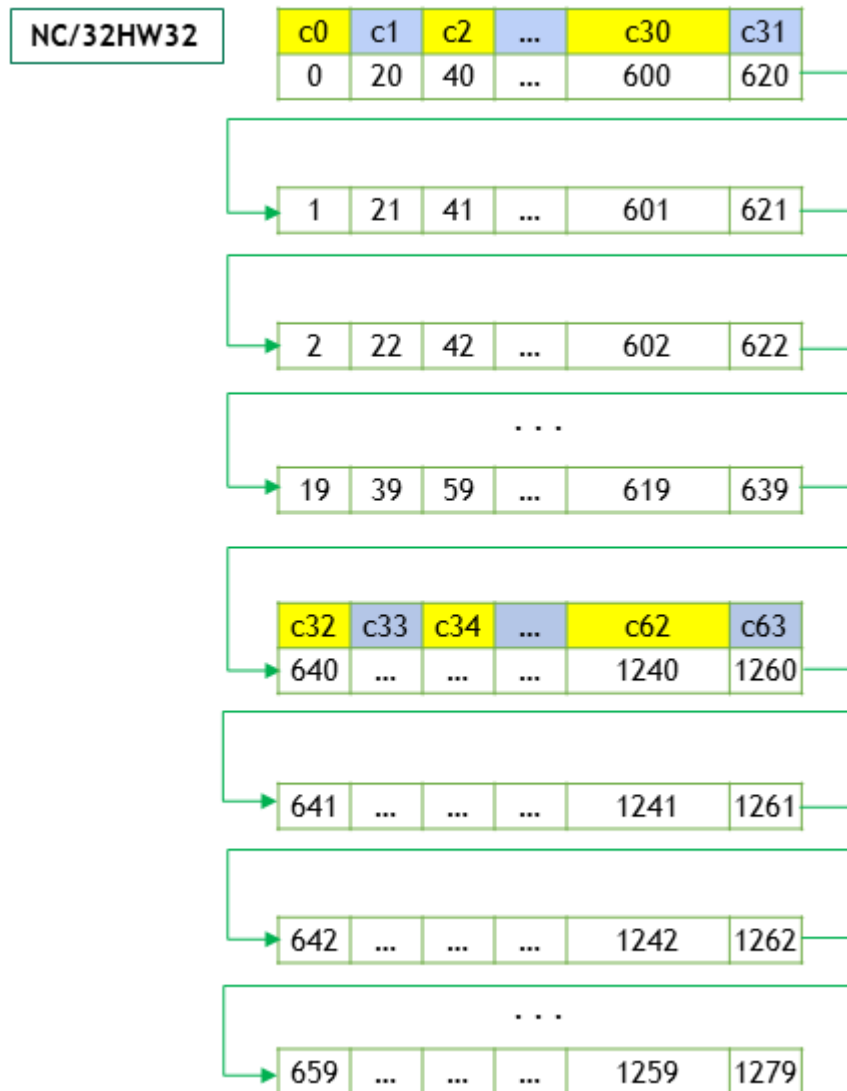
Figure 3. NHWC Memory Layout



6.4. NC/32HW32 Memory Layout

The NC/32HW32 is similar to NHWC, with a key difference. For the NC/32HW32 memory layout, the 64 channels are grouped into two groups of 32 channels each—first group consisting of channels c0 through c31, and the second group consisting of channels c32 through c63. Then each group is laid out using the NHWC format. See [Figure 4](#).

Figure 4. NC/32HW32 Memory Layout



For the generalized NC/xHWx layout format, the following observations apply:

- ▶ Only the channel dimension, **C**, is grouped into x channels each.
- ▶ When $x = 1$, each group has only one channel. Hence, the elements of one channel (i.e., one group) are arranged contiguously (in the row-major order), before proceeding to the next group (i.e., next channel). This is the same as NCHW format.
- ▶ When $x = C$, then NC/xHWx is identical to NHWC, i.e., the entire channel depth C is considered as a single group. The case $x = C$ can be thought of as vectorizing the entire C dimension as one big vector, laying out all the Cs, followed by the remaining dimensions, just like NHWC.
- ▶ The tensor format `CUDNN_TENSOR_NCHW_VECT_C` can also be interpreted in the following way: The NCHW INT8x32 format is really $N \times (C/32) \times H \times W \times 32$ (32 Cs for every

W), just as the NCHW INT8x4 format is $N \times (C/4) \times H \times W \times 4$ (4 Cs for every W). Hence the "VECT_C" name - each W is a vector (4 or 32) of Cs.

Chapter 7. Thread Safety

The library is thread-safe and its functions can be called from multiple host threads, as long as threads do not share the same cuDNN handle simultaneously.

Chapter 8. Reproducibility (determinism)

By design, most of cuDNN routines from a given version generate the same bit-wise results across runs when executed on GPUs with the same architecture and the same number of SMs. However, bit-wise reproducibility is not guaranteed across versions, as the implementation of a given routine may change. With the current release, the following routines do not guarantee reproducibility because they use atomic operations:

- ▶ `cudaConvolutionBackwardFilter` when `CUDNN_CONVOLUTION_BWD_FILTER_ALGO_0` or `CUDNN_CONVOLUTION_BWD_FILTER_ALGO_3` is used
- ▶ `cudaConvolutionBackwardData` when `CUDNN_CONVOLUTION_BWD_DATA_ALGO_0` is used
- ▶ `cudaPoolingBackward` when `CUDNN_POOLING_MAX` is used
- ▶ `cudaSpatialTfSamplerBackward`
- ▶ `cudaCTCLoss` and `cudaCTCLoss_v8` when `CUDNN CTC_LOSS_ALGO_NON_DETERMINISTIC` is used

Chapter 9. Scaling Parameters

Many cuDNN routines like [cudnnConvolutionForward\(\)](#) accept pointers in host memory to scaling factors `alpha` and `beta`. These scaling factors are used to blend the computed values with the prior values in the destination tensor as follows (see [Figure 5](#)):

```
dstValue = alpha*computedValue + beta*priorDstValue
```



Note: The `dstValue` is written to after being read.

Figure 5. Scaling Parameters for Convolution

When `beta` is zero, the output is not read and may contain uninitialized data (including NaN). These parameters are passed using a host memory pointer. The storage data types for `alpha` and `beta` are:

- ▶ `float` for HALF and FLOAT tensors, and
- ▶ `double` for DOUBLE tensors.



Note: For improved performance use `beta = 0.0`. Use a non-zero value for `beta` only when you need to blend the current output tensor values with the prior values of the output tensor.

Type Conversion

When the data input x , the filter input w and the output y are all in INT8 data type, the function [cudnnConvolutionBiasActivationForward\(\)](#) will perform the type conversion as shown in [Figure 6](#):



Note: Accumulators are 32-bit integers that wrap on overflow.

Figure 6. INT8 for cudnnConvolutionBiasActivationForward

Chapter 10. Tensor Core Operations

The cuDNN v7 library introduced the acceleration of compute-intensive routines using Tensor Core hardware on supported GPU SM versions. Tensor Core operations are supported beginning with the Volta GPU.

10.1. Basics

Tensor Core operations perform parallel floating-point accumulation of multiple floating-point product terms. Setting the math mode to `CUDNN_TENSOR_OP_MATH` via the [`cudaMathType_t`](#) enumerator indicates that the library will use Tensor Core operations. This enumerator specifies the available options to enable the Tensor Core and should be applied on a per-routine basis.

The default math mode is `CUDNN_DEFAULT_MATH`, which indicates that the Tensor Core operations will be avoided by the library. Because the `CUDNN_TENSOR_OP_MATH` mode uses the Tensor Cores, it is possible that these two modes generate slightly different numerical results due to different sequencing of the floating-point operations.

For example, the result of multiplying two matrices using Tensor Core operations is very close, but not always identical, to the result achieved using a sequence of scalar floating-point operations. For this reason, the cuDNN library requires an explicit user opt-in before enabling the use of Tensor Core operations.

However, experiments with training common deep learning models show negligible differences between using Tensor Core operations and scalar floating point paths, as measured by both the final network accuracy and the iteration count to convergence. Consequently, the cuDNN library treats both modes of operation as functionally indistinguishable and allows for the scalar paths to serve as legitimate fallbacks for cases in which the use of Tensor Core operations is unsuitable.

Kernels using Tensor Core operations are available for both convolutions and RNNs.

See also [Training with Mixed Precision](#).

10.2. Convolution Functions

10.2.1. Prerequisites

For the supported GPUs, the Tensor Core operations will be triggered for convolution functions only when `cudaSetConvolutionMathType()` is called on the appropriate convolution descriptor by setting the `mathType` to `CUDNN_TENSOR_OP_MATH` or `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION`.

10.2.2. Supported Algorithms

When the prerequisite is met, the below convolution functions can be run as Tensor Core operations:

- ▶ [cudaConvolutionForward\(\)](#)
- ▶ [cudaConvolutionBackwardData\(\)](#)
- ▶ [cudaConvolutionBackwardFilter\(\)](#)

See the table below for supported algorithms:

Supported Convolution Function	Supported Algos
<code>cudaConvolutionForward</code>	<ul style="list-style-type: none"> ▶ <code>CUDNN_CONVOLUTION_FWD_ALGO_IMPLICIT_PRECOMP_GEMM</code> ▶ <code>CUDNN_CONVOLUTION_FWD_ALGO_WINOGRAD_NONFUSED</code>
<code>cudaConvolutionBackwardData</code>	<ul style="list-style-type: none"> ▶ <code>CUDNN_CONVOLUTION_BWD_DATA_ALGO_1</code> ▶ <code>CUDNN_CONVOLUTION_BWD_DATA_ALGO_WINOGRAD_NONFUSED</code>
<code>cudaConvolutionBackwardFilter</code>	<ul style="list-style-type: none"> ▶ <code>CUDNN_CONVOLUTION_BWD_FILTER_ALGO_1</code> ▶ <code>CUDNN_CONVOLUTION_BWD_FILTER_ALGO_WINOGRAD_NONFUSED</code>

10.2.3. Data And Filter Formats

The cuDNN library may use padding, folding, and NCHW-to-NHWC transformations to call the Tensor Core operations. See [Tensor Transformations](#).

For algorithms other than `*_ALGO_WINOGRAD_NONFUSED`, when the following requirements are met, the cuDNN library will trigger the Tensor Core operations:

- ▶ Input, filter, and output descriptors (`xDesc`, `yDesc`, `wDesc`, `dxDesc`, `dyDesc` and `dwDesc` as applicable) are of the `dataType = CUDNN_DATA_HALF` (i.e., FP16). For FP32 `dataType` see [FP32-to-FP16 Conversion](#).
- ▶ The number of input and output feature maps (i.e., channel dimension `c`) is a multiple of 8. When the channel dimension is not a multiple of 8, see [Padding](#).
- ▶ The filter is of type `CUDNN_TENSOR_NCHW` or `CUDNN_TENSOR_NHWC`.
- ▶ If using a filter of type `CUDNN_TENSOR_NHWC`, then the input, filter, and output data pointers (`x`, `y`, `w`, `dx`, `dy`, and `dw` as applicable) are aligned to 128-bit boundaries.

10.3. RNN Functions

10.3.1. Prerequisites

Tensor core operations will be triggered for these RNN functions only when [cudnnSetRNNMatrixMathType\(\)](#) is called on the appropriate RNN descriptor setting `mathType` to `CUDNN_TENSOR_OP_MATH` or `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION`.

10.3.2. Supported Algorithms

When the above prerequisite is met, the RNN functions below can be run as Tensor Core operations:

- ▶ [cudnnRNNForwardInference\(\)](#)
- ▶ [cudnnRNNForwardTraining\(\)](#)
- ▶ [cudnnRNNBackwardData\(\)](#)
- ▶ [cudnnRNNBackwardWeights\(\)](#)
- ▶ [cudnnRNNForwardInferenceEx\(\)](#)
- ▶ [cudnnRNNForwardTrainingEx\(\)](#)
- ▶ [cudnnRNNBackwardDataEx\(\)](#)
- ▶ [cudnnRNNBackwardWeightsEx\(\)](#)

See the table below for the supported algorithms:

RNN Function	Support Algos
All RNN functions that support Tensor Core operations.	<ul style="list-style-type: none"> ▶ <code>CUDNN_RNN_ALGO_STANDARD</code> ▶ <code>CUDNN_RNN_ALGO_PERSIST_STATIC</code>

10.3.3. Data And Filter Formats

When the following requirements are met, then the cuDNN library will trigger the Tensor Core operations:

- ▶ For `algo = CUDNN_RNN_ALGO_STANDARD`:
 - ▶ The hidden state size, input size, and the batch size is a multiple of 8.
 - ▶ All user-provided tensors, workspace, and reserve space are aligned to 128-bit boundaries.
 - ▶ For FP16 input/output, the `CUDNN_TENSOR_OP_MATH` or `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` is selected.
 - ▶ For FP32 input/output, `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` is selected.

- ▶ For `algo = CUDNN_RNN_ALGO_PERSIST_STATIC`:
 - ▶ The hidden state size and the input size is a multiple of 32.
 - ▶ The batch size is a multiple of 8.
 - ▶ If the batch size exceeds 96 (for forward training or inference) or 32 (for backward data), then the batch size constraints may be stricter, and large power-of-two batch sizes may be needed.
 - ▶ All user-provided tensors, workspace, and reserve space are aligned to 128-bit boundaries.
 - ▶ For FP16 input/output, `CUDNN_TENSOR_OP_MATH` or `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` is selected.
 - ▶ For FP32 input/output, `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` is selected.

See also [Features Of RNN Functions](#).

10.4. Tensor Transformations

A few functions in the cuDNN library will perform transformations such as folding, padding, and NCHW-to-NHWC conversion while performing the actual function operation. See below.

10.4.1. FP16 Data

Tensor Cores operate on FP16 input data with FP32 accumulation. The FP16 multiply leads to a full-precision result that is accumulated in FP32 operations with the other products in a given dot product for a matrix with $m \times n \times k$ dimensions. See [Figure 7](#).

Figure 7. Tensor Operation with FP16 Inputs

10.4.2. FP32-to-FP16 Conversion

The cuDNN API allows the user to specify that FP32 input data may be copied and converted to FP16 data internally to use Tensor Core operations for potentially improved performance. This can be achieved by selecting `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` enum for [`cudnnMathType_t`](#) . In this mode, the FP32 tensors are internally down-converted to FP16, the Tensor Op math is performed, and finally up-converted to FP32 as outputs. See [Figure 8](#).

Figure 8. Tensor Operation with FP32 Inputs

For Convolutions

For convolutions, the FP32-to-FP16 conversion can be achieved by passing the `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` enum value to the [`cudnnSetConvolutionMathType\(\)`](#) call.

```
// Set the math type to allow cuDNN to use Tensor Cores:
```

```
checkCudnnErr (cudnnSetConvolutionMathType (cudnnConvDesc,
CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION));
```

For RNNs

For RNNs, the FP32-to-FP16 conversion can be achieved by passing the CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION enum value to the [cudnnSetRNNMatrixMathType\(\)](#) call to allow FP32 data to be converted for use in RNNs.

```
// Set the math type to allow cuDNN to use Tensor Cores:
checkCudnnErr (cudnnSetRNNMatrixMathType (cudnnRnnDesc,
CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION));
```

10.4.3. Padding

For packed NCHW data, when the channel dimension is not a multiple of 8, then the cuDNN library will pad the tensors as needed to enable Tensor Core operations. This padding is automatic for packed NCHW data in both the CUDNN_TENSOR_OP_MATH and the CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION cases.

The padding occurs with a negligible loss of performance. Hence, the NCHW tensor dimensions such as below are allowed:

```
// Set NCHW Tensor dimensions, not necessarily as multiples of eight (only the input
tensor is shown here):
int dimA[] = {1, 7, 32, 32};
int strideA[] = {7168, 1024, 32, 1};
```

10.4.4. Folding

In the folding operation, the cuDNN library implicitly performs the formatting of input tensors and saves the input tensors in an internal workspace. This can lead to an acceleration of the call to Tensor Cores.

Folding enables the input tensors to be transformed into a format that the Tensor Cores support (i.e., no strides).

10.4.5. Conversion Between NCHW And NHWC

Tensor Cores require that the tensors be in the NHWC data layout. Conversion between NCHW and NHWC is performed when the user requests Tensor Op math. However, as stated in [Basics](#), a request to use Tensor Cores is just that, a request and Tensor Cores may not be used in some cases. The cuDNN library converts between NCHW and NHWC if and only if Tensor Cores are requested and are actually used.

If your input (and output) are NCHW, then expect a layout change.

Non-Tensor Op convolutions will not perform conversions between NCHW and NHWC.

In very rare and difficult-to-qualify cases that are a complex function of padding and filter sizes, it is possible that Tensor Ops is not enabled. In such cases, users should pre-pad.

10.5. Guidelines For A Deep Learning Compiler

For a deep learning compiler, the following are the key guidelines:

- ▶ Make sure that the convolution operation is eligible for Tensor Cores by avoiding any combinations of large padding and large filters.
- ▶ Transform the inputs and filters to NHWC, pre-pad channel and batch size to be a multiple of 8.
- ▶ Make sure that all user-provided tensors, workspace, and reserve space are aligned to 128-bit boundaries.

Chapter 11. GPU And Driver Requirements

For the latest compatibility software versions of the OS, CUDA, the CUDA driver, and the NVIDIA hardware, see the [cuDNN Support Matrix](#).

Chapter 12. Backward Compatibility And Deprecation Policy

cuDNN version 8 introduces a new API deprecation policy to enable a faster pace of innovation.

The old deprecation policy required three major library releases to complete an API update. During this process, the original function name was first assigned to the legacy API, and then to the revised API, depending on the library version. The user wishing to migrate to the new API version had to update his or her code twice. In the first update, the original call `foo()` had to be changed to `foo_vN()`, where `N` is the new major cuDNN version. After the next major cuDNN release, the `foo_vN()` function had to be renamed back as `foo()`. Clearly, the above process could be difficult for code maintenance, especially when many functions are upgraded.

A streamlined, two-step, deprecation policy will be used for all API changes starting with cuDNN version 8. Let us explain the process using two subsequent, major cuDNN releases, version 8 and 9:

Table 2. Two-step, deprecation policy

cuDNN version	Explanation
Major release 8	The updated API is introduced as <code>foo_v8()</code> . The deprecated API <code>foo()</code> is kept unchanged to maintain backward compatibility until the next major release.
Major release 9	The deprecated API <code>foo()</code> is permanently removed and its name is not reused. The <code>foo_v8()</code> function supersedes the retired call <code>foo()</code> .

If the existing API needs to be updated, a new function flavor is introduced with the `_v` tag followed by the current, major cuDNN version. In the next major release, the deprecated function is removed, and its name is never reused. A brand-new API is first introduced without the `_v` tag.

The revised deprecation scheme allows us to retire the legacy API in just one major release. Similarly to the previous API deprecation policy, the user is able to compile the legacy code without any changes using the next major release of the cuDNN library. The backward compatibility ends when another major cuDNN release is introduced.

The updated function name embeds the information in which the cuDNN version of the API call was modified. As a result, the API changes will be easier to track and document.

The new deprecation policy is applied also to pending API changes from previous cuDNN releases. For example, according to the old deprecation policy, `cudaSetRNNDescriptor_v6()` should be removed in cuDNN version 8 and the upgraded call `cudaSetRNNDescriptor()` with the same arguments and behavior should be kept. Instead, the new deprecation policy is applied to this case and the tagged function is kept.

Prototypes of deprecated functions will be prepended in cuDNN version 8 headers using the `CUDNN_DEPRECATED` macro. When the `-DCUDNN_WARN_DEPRECATED` switch is passed to the compiler, any deprecated function call in the user's code will emit a compiler warning, for example:

```
warning: 'cudaStatus_t cudaSetRNNMatrixMathType(cudaRNNDescriptor_t, cudaMathType_t)' is deprecated [-Wdeprecated-declarations]
```

Or

```
warning C4996: 'cudaSetRNNMatrixMathType': was declared deprecated
```

The above warnings are disabled by default to avoid potential build breaks in software setups where compiler warnings are treated as errors.

Note that the simple swapping of older cuDNN version 7 shared library files will not work with the cuDNN version 8 release. The user source code needs to be recompiled from scratch with the cuDNN version 8 headers and linked with the version 8 libraries.

Chapter 13. Grouped Convolutions

cuDNN supports grouped convolutions by setting `groupCount > 1` for the convolution descriptor `convDesc`, using `cudaSetConvolutionGroupCount()`.



Note: By default, the convolution descriptor `convDesc` is set to `groupCount` of 1.

Basic Idea

Conceptually, in grouped convolutions, the input channels and the filter channels are split into `groupCount` number of independent groups, with each group having a reduced number of channels. The convolution operation is then performed separately on these input and filter groups.

For example, consider the following: if the number of input channels is 4, and the number of filter channels of 12. For a normal, ungrouped convolution, the number of computation operations performed are 12×4 .

If the `groupCount` is set to 2, then there are now two input channel groups of two input channels each, and two filter channel groups of six filter channels each.

As a result, each grouped convolution will now perform 2×6 computation operations, and two such grouped convolutions are performed. Hence the computation savings are 2x: $(12 \times 4) / (2 \times (2 \times 6))$

cuDNN Grouped Convolution

- ▶ When using `groupCount` for grouped convolutions, you must still define all tensor descriptors so that they describe the size of the entire convolution, instead of specifying the sizes per group.
- ▶ Grouped convolutions are supported for all formats that are currently supported by the functions `cudaConvolutionForward()`, `cudaConvolutionBackwardData()` and `cudaConvolutionBackwardFilter()`.
- ▶ The tensor stridings that are set for `groupCount` of 1 are also valid for any group count.
- ▶ By default, the convolution descriptor `convDesc` is set to `groupCount` of 1.



Note: See [Convolution Formulas](#) for the math behind the cuDNN Grouped Convolution.

Example

Below is an example showing the dimensions and strides for grouped convolutions for NCHW format, for 2D convolution.



Note: The symbols "*" and "/" are used to indicate multiplication and division.

xDesc or dxDesc:

- ▶ **Dimensions:** [batch_size, input_channel, x_height, x_width]
- ▶ **Strides:** [input_channels*x_height*x_width, x_height*x_width, x_width, 1]

wDesc or dwDesc:

- ▶ **Dimensions:** [output_channels, input_channels/groupCount, w_height, w_width]
- ▶ **Format:** NCHW

convDesc:

- ▶ **Group Count:** groupCount

yDesc or dyDesc:

- ▶ **Dimensions:** [batch_size, output_channels, y_height, y_width]
- ▶ **Strides:** [output_channels*y_height*y_width, y_height*y_width, y_width, 1]

Chapter 14. API Logging

cuDNN API logging is a tool that records all input parameters passed into every cuDNN API function call. This functionality is disabled by default, and can be enabled through methods described in this section.

The log output contains variable names, data types, parameter values, device pointers, process ID, thread ID, cuDNN handle, CUDA stream ID, and metadata such as time of the function call in microseconds.

When logging is enabled, the log output will be handled by the built-in default callback function. The user may also write their own callback function, and use the [`cudaSetCallback\(\)`](#) to pass in the function pointer of their own callback function. The following is a sample output of the API log.

```
Function cudnnSetActivationDescriptor() called:
mode: type=cudnnActivationMode_t; val=CUDNN_ACTIVATION_RELU (1);
reluNanOpt: type=cudnnNanPropagation_t; val=CUDNN_NOT_PROPAGATE_NAN (0);
coef: type=double; val=1000.000000;
Time: 2017-11-21T14:14:21.366171 (0d+0h+1m+5s since start)
Process: 21264, Thread: 21264, cudnn_handle: NULL, cudnn_stream: NULL.
```

There are two methods to enable API logging.

Method 1: Using Environment Variables

To enable API logging using environment variables, follow these steps:

- ▶ Set the environment variable `CUDNN_LOGINFO_DBG` to "1", and
- ▶ Set the environment variable `CUDNN_LOGDEST_DBG` to one of the following:
 - ▶ `stdout`, `stderr`, or a user-desired file path, for example, `/home/userName1/log.txt`.
- ▶ Include the conversion specifiers in the file name. For example:
 - ▶ To include date and time in the file name, use the date and time conversion specifiers: `log_%Y_%m_%d_%H_%M_%S.txt`. The conversion specifiers will be automatically replaced with the date and time when the program is initiated, resulting in `log_2017_11_21_09_41_00.txt`.
 - ▶ To include the process id in the file name, use the `%i` conversion specifier: `log_%Y_%m_%d_%H_%M_%S_%i.txt` for the result: `log_2017_11_21_09_41_00_21264.txt` when the process id is 21264. When you have several processes running, using the process

id conversion specifier will prevent these processes from writing to the same file at the same time.



Note: The supported conversion specifiers are similar to the `strftime` function.

If the file already exists, the log will overwrite the existing file.



Note: These environmental variables are only checked once at the initialization. Any subsequent changes in these environmental variables will not be effective in the current run. Also note that these environment settings can be overridden by Method 2 below.

See also [Table 3](#) for the impact on the performance of API logging using environment variables.

Table 3. API Logging Using Environment Variables

Environment variables	CUDNN_LOGINFO_DBG=0	CUDNN_LOGINFO_DBG=1
CUDNN_LOGDEST_DBG not set	No logging output No performance loss	No logging output No performance loss
CUDNN_LOGDEST_DBG=NULL	No logging output No performance loss	No logging output No performance loss
CUDNN_LOGDEST_DBG=stdout or stderr	No logging output No performance loss	Logging to stdout or stderr Some performance loss
CUDNN_LOGDEST_DBG=filename.txt	No logging output No performance loss	Logging to filename.txt Some performance loss

Method 2


Method 2: To use API function calls to enable API logging, refer to the API description of [`cudaSetCallback\(\)`](#) and [`cudaGetCallback\(\)`](#).

Chapter 15. Features Of RNN Functions

The RNN functions are:

- ▶ [cudnnRNNForwardInference\(\)](#)
- ▶ [cudnnRNNForwardTraining\(\)](#)
- ▶ [cudnnRNNBackwardData\(\)](#)
- ▶ [cudnnRNNBackwardWeights\(\)](#)
- ▶ [cudnnRNNForwardInferenceEx\(\)](#)
- ▶ [cudnnRNNForwardTrainingEx\(\)](#)
- ▶ [cudnnRNNBackwardDataEx\(\)](#)
- ▶ [cudnnRNNBackwardWeightsEx\(\)](#)

See the table below for a list of features supported by each RNN function:

 **Note:**

For each of these terms, the short-form versions shown in the parenthesis are used in the tables below for brevity: `CUDNN_RNN_ALGO_STANDARD` (`_ALGO_STANDARD`), `CUDNN_RNN_ALGO_PERSIST_STATIC` (`_ALGO_PERSIST_STATIC`), `CUDNN_RNN_ALGO_PERSIST_DYNAMIC` (`_ALGO_PERSIST_DYNAMIC`), and `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` (`_ALLOW_CONVERSION`).

Functions	Input/output layout supported	Supports variable sequence length in batch	Commonly supported
cudnnRNNForwardInference() cudnnRNNForwardTraining() cudnnRNNBackwardData() cudnnRNNBackwardWeights()	Only Sequence major, packed (non-padded)	Only with <code>_ALGO_STANDARD</code> Require input sequences descending sorted according to length	Mode (cell type) supported: <code>CUDNN_RNN_RELU</code> , <code>CUDNN_RNN_TANH</code> , <code>CUDNN_LSTM</code> , <code>CUDNN_GRU</code>
cudnnRNNForwardInferenceEx() cudnnRNNForwardTrainingEx()	▶ Sequence major, unpacked	Only with <code>_ALGO_STANDARD</code>	Algo supported ¹ (see the table below for an

¹ Do not mix different algos for different steps of training. It's also not recommended to mix non-extended and extended API for different steps of training.

Functions	Input/output layout supported	Supports variable sequence length in batch	Commonly supported
cudnnRNNBackwardDataEx cudnnRNNBackwardWeightsEx	▶ Batch major unpacked ² ▶ Sequence major packed ²	For unpacked layout ² , no input sorting required. For packed layout, require input sequences descending sorted according to length	elaboration on these algorithms): _ALGO_STANDARD, _ALGO_PERSIST_STATIC, _ALGO_PERSIST_DYNAMIC Math mode supported: CUDNN_DEFAULT_MATH, CUDNN_TENSOR_OP_MATH (will automatically fall back if run on pre-Volta or if algo doesn't support Tensor Cores) _ALLOW_CONVERSION (may do down conversion to utilize Tensor Cores) Direction mode supported: CUDNN_UNIDIRECTIONAL, CUDNN_BIDIRECTIONAL RNN input mode: CUDNN_LINEAR_INPUT, CUDNN_SKIP_INPUT

The following table provides the features supported by the algorithms referred in the above table: CUDNN_RNN_ALGO_STANDARD, CUDNN_RNN_ALGO_PERSIST_STATIC, and CUDNN_RNN_ALGO_PERSIST_DYNAMIC.

Features	_ALGO_STANDARD	_ALGO_PERSIST_STATIC	_ALGO_PERSIST_DYNAMIC
Half input	Supported		
Single accumulation	Half intermediate storage		
Half output	Single accumulation		
Single input	Supported		
Single accumulation	If running on Volta, with CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION ¹ , will down-convert and use half intermediate storage.		
Single output	Otherwise: Single intermediate storage Single accumulation		
Double input	Supported	Not Supported	Supported
Double accumulation			

² To use unpacked layout, user need to set CUDNN_RNN_PADDED_IO_ENABLED through cudnnSetRNNPaddingMode().

Features	<code>_ALGO_STANDARD</code>	<code>_ALGO_PERSIST_STAT</code>	<code>_ALGO_PERSIST_DYNAMIC</code>
Double output	Double intermediate storage Double accumulation		Double intermediate storage Double accumulation
LSTM recurrent projection	Supported	Not Supported	Not Supported
LSTM cell clipping	Supported		
Variable sequence length in batch	Supported	Not Supported	Not Supported
Tensor Cores on Volta/Xavier	Supported For half input/output, acceleration requires setting <code>CUDNN_TENSOR_OP_MATH</code> ³ or <code>CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION</code> ³ Acceleration requires <code>inputSize</code> and <code>hiddenSize</code> to be a multiple of 8 For single input/output, acceleration requires setting <code>CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION</code> ³ Acceleration requires <code>inputSize</code> and <code>hiddenSize</code> to be a multiple of 8		Not Supported, will execute normally ignoring <code>CUDNN_TENSOR_OP_MATH</code> ³ or <code>_ALLOW_CONVERSION</code> ³
Other limitations		Max problem size is limited by GPU specifications.	Requires real time compilation through NVRTC

³ `CUDNN_TENSOR_OP_MATH` or `CUDNN_TENSOR_OP_MATH_ALLOW_CONVERSION` can be set through `cudaSetRNNMatrixMathType()`.

Chapter 16. Mixed Precision Numerical Accuracy

When the computation precision and the output precision are not the same, it is possible that the numerical accuracy will vary from one algorithm to the other.

For example, when the computation is performed in FP32 and the output is in FP16, the `CUDNN_CONVOLUTION_BWD_FILTER_ALGO_0` (`ALGO_0`) has lower accuracy compared to the `CUDNN_CONVOLUTION_BWD_FILTER_ALGO_1` (`ALGO_1`). This is because `ALGO_0` does not use extra workspace, and is forced to accumulate the intermediate results in FP16, i.e., half precision float, and this reduces the accuracy. The `ALGO_1`, on the other hand, uses additional workspace to accumulate the intermediate values in FP32, i.e., full precision float.

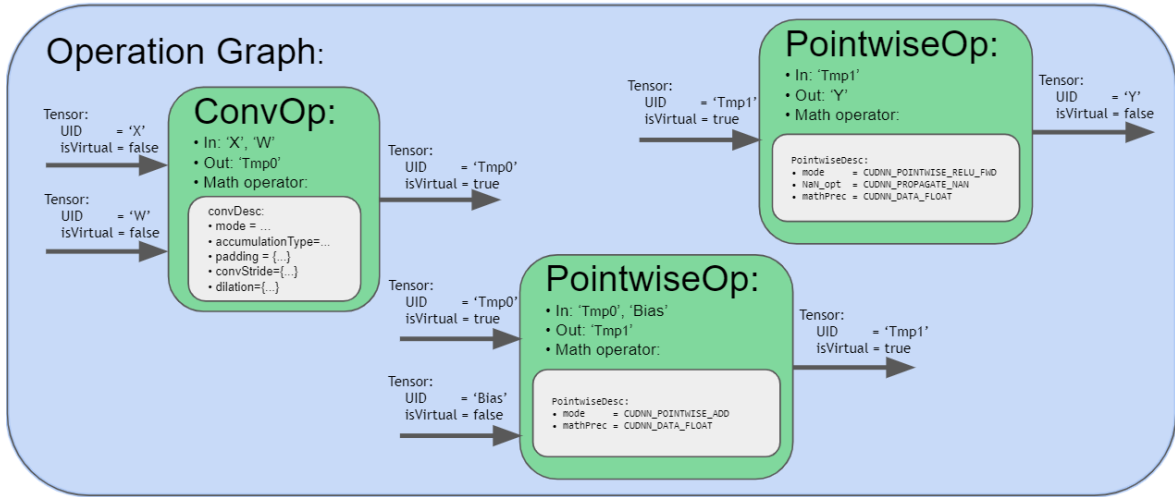
Chapter 17. Operation Fusion Via The Backend API

Introduced in cuDNN 8.0, operation fusion can be achieved via the backend API. The general workflow is similar to running unfused operations, except that instead of creating a single operation Operation Graph, the user may specify a multi-operation Operation Graph. Here we illustrate the flow via an example.

In the following example, the user would like to implement a fusion operation of convolution, bias, and activation.

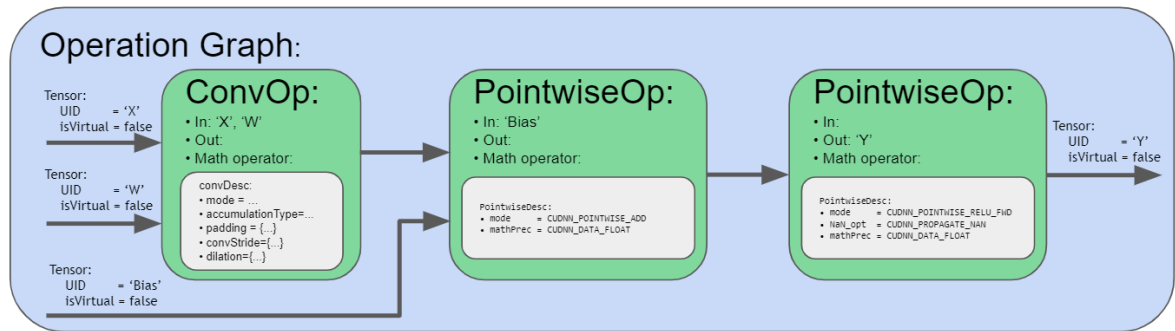
- ▶ First, the user should create three cuDNN backend operation descriptors - one convolution operation descriptor and two pointwise operation descriptors. Depending on the pointwise mode in the pointwise descriptor, a pointwise operation descriptor can be set up to describe an activation operation or a bias operation. By specifying the backend tensor *Tmp0* as both the output of the convolution operation and the input of the bias operation, this allows cuDNN to infer the dataflow between the operations. The same applies to tensor *Tmp1*. Here assume the user doesn't need the intermediate results *Tmp0* and *Tmp1* for any other use, then the user can specify them to be virtual tensors, so the memory I/Os can later be optimized out.
 - ▶ Note for the purpose of fusion, users should not construct in-place operations where any of the input UIDs matches any of its own output UIDs. Such in-place operations will be considered cyclic in later graph analysis and deemed unsupported.
 - ▶ Also note that the operation descriptors can be passed into cuDNN in any order, as the tensor UIDs are enough to determine the dependencies in the graph.

Figure 9. A set of operation descriptors the user passes to the operation graph



- ▶ Second, upon finalizing the operation graph, cuDNN will perform the dataflow analysis to establish the dependency relationship between operations and connect the edges, as illustrated in the figure below. In this step, cuDNN will also perform various checks to confirm the validity of the graph.

Figure 10. The operation graph after finalization



- ▶ Third, with the finalized operation graph, there are two options:
 - ▶ For most users that prefer cuDNN to recommend the best engine and knob choices, they can query cuDNN's heuristics to get a list of engine configs and choose from them. After that, the user can construct the execution plan using the chosen engine config. Note the heuristics support for fusion use cases are not yet available. This will be available in the coming releases.
 - ▶ For expert users, they can query the engines that can support this operation graph. For each engine, the user can further query the numerical notes and adjustable knobs. Numerical notes would inform the user about the numerical behavior of the engine such as whether it does datatype down conversion at the input or during output

reduction. The adjustable knobs allow fine grained control of the engine’s behavior and performance. With the engine choice and the knob choice determined, the user can construct the backend engine, backend engine config, and further the execution plan.

Note for operation fusion use cases, there are two different mechanisms in cuDNN to support them. First, there are engines containing offline compiled kernels that can support certain fusion patterns. These engines try to match the user provided operation graph with their supported fusion pattern. If there is a match, then that particular engine is deemed suitable for this use case. In addition, there are also runtime fusion engines to be made available in the upcoming releases. Instead of passively matching the user graph, such engines actively walk the graph and assemble code blocks to form a CUDA kernel and compile on the fly. Such runtime fusion engines are much more flexible in its range of support. However, because the construction of the execution plans requires runtime compilation, the one-time CPU overhead is higher than the other engines.

- ▶ Finally, with the execution plan constructed and when it comes time to run it, the user should construct the backend variant pack by providing the workspace pointer, an array of UIDs, and an array of device pointers. The UIDs and the pointers should be in the corresponding order. With the handle, the execution plan and variant pack, the execution API can be called and the computation is carried out on the GPU.

The table below briefly summarizes the current fusion support in cuDNN. We will be adding additional support in the upcoming releases. We welcome feature suggestions. For feedback, email cuda@nvidia.com.

Fusion Graph Pattern	Supported Device Compute Capabilities	Supported Data Config and Layout	Supported Engine Types
Conv_Bias_Add_activation	All that cuDNN supports	Same as cudnnConvolutionBiasAdd	Pattern matching engines, runtime fusion engines ⁴
Scale_Bias_Activation	Compute capability 70 or above	PSEUDO_HALF_CONFIG, NHWC layout	Pattern matching engines, runtime fusion engines ¹
Convolution_Pointwise	Compute capability 75 or above	Flexible	Runtime fusion engines ¹
Gemm_Pointwise ¹	Compute capability 75 or above	Flexible	Runtime fusion engines ¹

⁴ As of cuDNN 8.0.0, the runtime fusion engines are not yet available. They will be made available in the upcoming releases.

Chapter 18. Troubleshooting

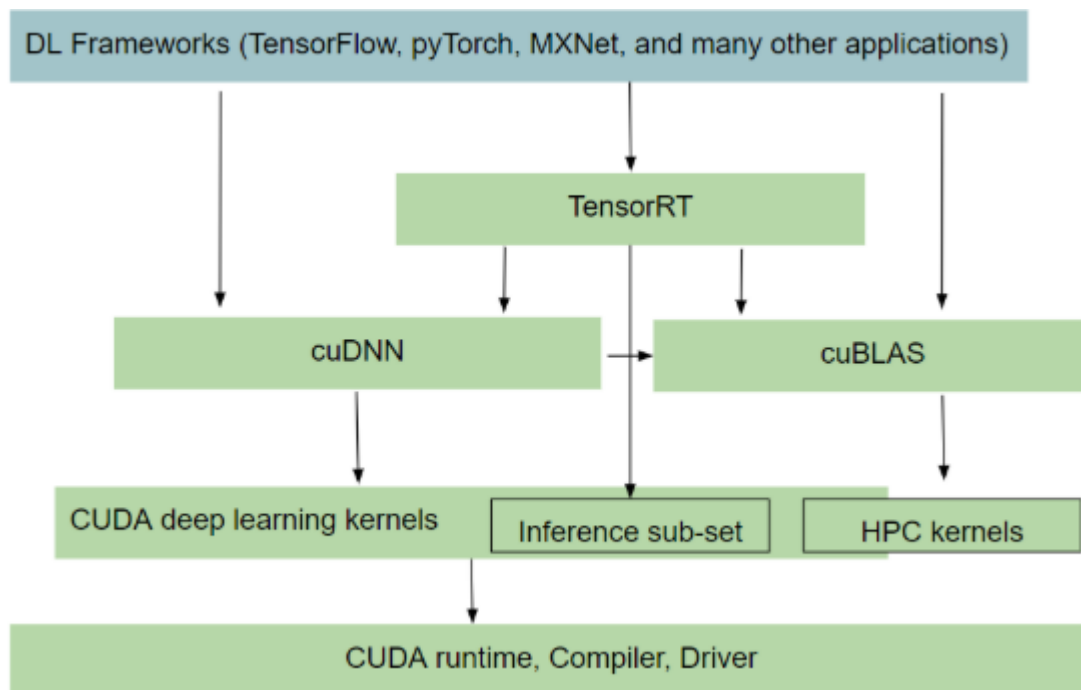
The following sections help answer the most commonly asked questions regarding typical use cases.

18.1. FAQs

Q: Where in the software stack does cuDNN sit? What is the interaction between CUDA, cuDNN, and TensorRT?

A: The following graphic shows how cuDNN relates to other software in the stack.

Figure 11. Software stack with cuDNN.



Q: I'm not sure if I should use cuDNN for inference or training. How does it compare with TensorRT?

A: cuDNN provides the building blocks for common routines such as convolution, pooling, activation and RNN/LSTMs. You can use cuDNN for both training and inference. However, where it differs from TensorRT is that the latter (TensorRT) is a programmable inference accelerator; just like a framework. TensorRT sees the whole graph and optimizes the network by fusing/combining layers and optimizing kernel selection for improved latency, throughput, power efficiency and for reducing memory requirements.

A rule of thumb you can apply is to check out TensorRT, see if it meets your inference needs, if it doesn't, then look at cuDNN for a closer, more in-depth perspective.

Q: How does heuristics in cuDNN work? How does it know what is the optimal solution for a given problem?

A: NVIDIA actively monitors the Deep Learning space for important problem specifications such as commonly used models. The heuristics are produced by sampling a portion of these problem specifications with available computational choices. Over time, more models are discovered and incorporated into the heuristics.

Q: Is cuDNN going to support running arbitrary graphs?

A: No, we don't plan to become a framework and execute the whole graph one op at a time. At this time, we are focused on a subgraph given by the user, where we try to produce an optimized fusion kernel. We will document what the rules regarding what can be fused and what cannot. The goal is to support general and flexible fusion, however, it will take time and there will be limits in what it can do in the cuDNN version 8.0.0 launch.

Q: What's the difference between TensorRT, TensorFlow/XLA's fusion, and cuDNN's fusion?

A: TensorRT and TensorFlow are frameworks; they see the whole graph and can do global optimization, however, they generally only fuse pointwise ops together. On the other hand, cuDNN targets a subgraph, but can fuse convolutions with pointwise ops, thus providing potentially better performance. CuDNN fusion kernels can be utilized by TensorRT and TensorFlow/XLA as part of their global graph optimization.

Q: Can I write an application calling cuDNN directly?

A: Yes, you can call the C/C++ API directly. Usually, data scientists would wait for framework integration and use the Python API which is more convenient. However, if your use case requires better performance, you can target the cuDNN API directly.

Q: How does mixed precision training work?

A: Several components need to work together to make mixed precision training possible. CuDNN needs to support the layers with the required datatype config and have optimized kernels that run very fast. In addition, there is a module called automatic mixed precision (AMP) in frameworks which intelligently decides which op can run in a lower precision without affecting convergence and minimize the number of type conversions/transposes in the entire graph. These work together to give you speed up. For more information, see [Mixed Precision Numerical Accuracy](#).

Q: How can I pick the fastest convolution kernels with cuDNN version 8.0.0?

A: In the API introduced in cuDNN v8, convolution kernels are grouped by similar computation and numerical properties into engines. Every engine has a queryable set of performance tuning knobs. A computation case such as a convolution operation graph can be computed using different valid combinations of engines and their knobs, known as an engine configuration. Users can query an array of engine configurations for any given computation case ordered by performance, from fastest to slowest according to cuDNN's own heuristics. Alternately, users can generate all possible engine configurations by querying the engine count and available knobs for each engine. This generated list could be used for auto-tuning or the user could create their own heuristics.

Q: Why is cuDNN version 8.0 convolution API call much slower on the first call than subsequent calls?

A: Due to the library split, cuDNN version 8.0 API will only load the necessary kernels on the first API call that requires it. In previous versions, this load would have been observed in the first cuDNN API call that triggers CUDA context initialization, typically `cudaCreate()`. In version 8.0, this is delayed until the first sub-library call that triggers CUDA context initialization. Users who desire to have CUDA context preloaded can call the new `cudaCnnInferVersionCheck()` API (or its related cousins), which has the side effect of initializing a CUDA context. This will reduce the run time for all subsequent API calls.

Q: How do I build the cuDNN version 8.0.0 split library?

A: cuDNN v8.0 library is split into multiple sub-libraries. Each library contains a subset of the API. Users can link directly against the individual libraries or link with a `dlopen` layer which follows a plugin architecture.

To link against an individual library, users can directly specify it and its dependencies on the linker command line. For example, for infer libraries: `-lcudnn_adv_infer, -lcudnn_cnn_infer, or -lcudnn_ops_infer`.

For all libraries, `-lcudnn_adv_train, -lcudnn_cnn_train, -lcudnn_ops_train, -lcudnn_adv_infer, -lcudnn_cnn_infer, and -lcudnn_ops_infer`.

The dependency order is documented in the cuDNN [8.0.0 Preview Release Notes](#) and the [cuDNN API Reference](#).

Alternatively, the user can continue to link against a shim layer (`-libcuda`) which can `dlopen` the correct library that provides the implementation of the function. When the function is called for the first time, the dynamic loading of the library takes place.

Linker argument:

```
-lcudnn
```

Q: What are the new APIs in cuDNN version 8.0.0?

A: The new cuDNN APIs are listed in the cuDNN 8.0.0 Release Notes as well as in the [API Changes For cuDNN 8.0.0](#).

18.2. How Do I Report A Bug?

We appreciate all types of feedback. If you encounter any issues, please report them by following these steps.

1. Register for the [NVIDIA Developer website](#).
2. Log in to the developer site.
3. Click on your name in the upper right corner.
4. Click **My account > My Bugs** and select **Submit a New Bug**.
5. Fill out the bug reporting page. Be descriptive and if possible, provide the steps that you are following to help reproduce the problem.
6. Click **Submit a bug**.

18.3. Support

Support, resources, and information about cuDNN can be found online at <https://developer.nvidia.com/cudnn>. This includes downloads, webinars, [NVIDIA Developer Forums](#), and more.

For questions or to provide feedback, please contact cuDNN@nvidia.com.

Chapter 19. Acknowledgments

Some of the cuDNN library routines were derived from code developed by others and are subject to the following:

19.1. University of Tennessee

```
Copyright (c) 2010 The University of Tennessee.
```

```
All rights reserved.
```

```
Redistribution and use in source and binary forms, with or without  
modification, are permitted provided that the following conditions are  
met:
```

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer listed in this license in the documentation and/or other materials provided with the distribution.
- * Neither the name of the copyright holders nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

```
THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS  
"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT  
LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR  
A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT  
OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL,  
SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT  
LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE,  
DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY  
THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT  
(INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE  
OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
```

19.2. University of California, Berkeley

```
COPYRIGHT
```

```
All contributions by the University of California:  
Copyright (c) 2014, The Regents of the University of California (Regents)  
All rights reserved.
```

```
All other contributions:  
Copyright (c) 2014, the respective contributors
```


All rights reserved.

Caffe uses a shared copyright model: each contributor holds copyright over their contributions to Caffe. The project versioning records all such contribution and copyright details. If a contributor wants to further mark their specific copyright on a particular contribution, they should indicate their copyright solely in the commit message of the change when it is committed.

LICENSE

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

1. Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

CONTRIBUTION AGREEMENT

By contributing to the BVLC/caffe repository through pull-request, comment, or otherwise, the contributor releases their content to the license and copyright terms herein.

19.3. Facebook AI Research, New York

Copyright (c) 2014, Facebook, Inc. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- * Neither the name Facebook nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS

SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Additional Grant of Patent Rights

"Software" means fbcunn software distributed by Facebook, Inc.

Facebook hereby grants you a perpetual, worldwide, royalty-free, non-exclusive, irrevocable (subject to the termination provision below) license under any rights in any patent claims owned by Facebook, to make, have made, use, sell, offer to sell, import, and otherwise transfer the Software. For avoidance of doubt, no license is granted under Facebook's rights in any patent claims that are infringed by (i) modifications to the Software made by you or a third party, or (ii) the Software in combination with any software or other technology provided by you or a third party.

The license granted hereunder will terminate, automatically and without notice, for anyone that makes any claim (including by filing any lawsuit, assertion or other action) alleging (a) direct, indirect, or contributory infringement or inducement to infringe any patent: (i) by Facebook or any of its subsidiaries or affiliates, whether or not such claim is related to the Software, (ii) by any party if such claim arises in whole or in part from any software, product or service of Facebook or any of its subsidiaries or affiliates, whether or not such claim is related to the Software, or (iii) by any party relating to the Software; or (b) that any right in any patent claim of Facebook is invalid or unenforceable.

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.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk.

NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer's product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.

THIS DOCUMENT AND ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

VESA DisplayPort

DisplayPort and DisplayPort Compliance Logo, DisplayPort Compliance Logo for Dual-mode Sources, and DisplayPort Compliance Logo for Active Cables are trademarks owned by the Video Electronics Standards Association in the United States and other countries.

HDMI

HDMI, the HDMI logo, and High-Definition Multimedia Interface are trademarks or registered trademarks of HDMI Licensing LLC.

OpenCL

OpenCL is a trademark of Apple Inc. used under license to the Khronos Group Inc.

Trademarks

NVIDIA, the NVIDIA logo, and cuBLAS, CUDA, CUDA Toolkit, cuDNN, DALI, DIGITS, DGX, DGX-1, DGX-2, DGX Station, DLProf, GPU, JetPack, Jetson, Kepler, Maxwell, NCCL, Nsight Compute, Nsight Systems, NVCAffe, NVIDIA Ampere GPU architecture, NVIDIA Deep Learning SDK, NVIDIA Developer Program, NVIDIA GPU Cloud, NVLink, NVSHMEM, PerfWorks, Pascal, SDK Manager, T4, Tegra, TensorRT, TensorRT Inference Server, Tesla, TF-TRT, Triton Inference Server, Turing, and Volta are trademarks and/or registered trademarks of NVIDIA Corporation in the United States and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2017-2020 NVIDIA Corporation. All rights reserved.

