

Tuning CUDA Applications for Kepler

Application Note

Table of Contents

Chapter 1. Kepler Tuning Guide	1
1.1. NVIDIA Kepler Compute Architecture	1
1.2. CUDA Best Practices	1
1.3. Application Compatibility	2
1.4. Kepler Tuning	2
1.4.1. Device Utilization and Occupancy	2
1.4.2. Managing Coarse-Grained Parallelism	3
1.4.2.1. Concurrent Kernels	3
1.4.2.2. Hyper-Q	4
1.4.2.3. Dynamic Parallelism	4
1.4.3. Shared Memory and Warp Shuffle	5
1.4.3.1. Shared Memory Bandwidth	5
1.4.3.2. Shared Memory Capacity	5
1.4.3.3. Warp Shuffle	5
1.4.4. Memory Throughput	6
1.4.4.1. Increased Addressable Registers Per Thread	6
1.4.4.2. L1 Cache	6
1.4.4.3. Read-Only Data Cache	6
1.4.4.4. Fast Global Memory Atomics	7
1.4.4.5. Global Memory Bandwidth and GPU Boost	7
1.4.4.6. 2D Memory Copies	7
1.4.5. Instruction Throughput	7
1.4.5.1. Single-precision vs. Double-precision	7
1.4.6. GPU Boost	8
1.4.7. Multi-GPU	8
1.4.8. PCIe 3.0	8
1.4.9. Warp-synchronous Programming	9
Appendix A. References	.10
Appendix B. Revision History	11

Chapter 1. Kepler Tuning Guide

NVIDIA Kepler Compute Architecture

Kepler is NVIDIA's 3rd-generation architecture for CUDA compute applications. Kepler retains and extends the same CUDA programming model as in earlier NVIDIA architectures such as Fermi, and applications that follow the best practices for the Fermi architecture should typically see speedups on the Kepler architecture without any code changes. This guide summarizes the ways that an application can be fine-tuned to gain additional speedups by leveraging Kepler architectural features. 1

The Kepler architecture comprises two major variants: GK104 and GK110. A detailed overview of the major improvements in $GK104^2$ and $GK110^3$ over the earlier Fermi architecture are described in a pair of whitepapers [1][2] entitled NVIDIA GeForce GTX 680: The fastest, most efficient GPU ever built for GK104 and NVIDIA's Next Generation CUDA Compute Architecture: Kepler GK110 for GK110.

For details on the programming features discussed in this guide, please refer to the CUDA C ++ Programming Guide. Details on the architectural features are covered in the architecture whitepapers referenced above. Some of the Kepler features described in this guide are specific to GK110, as noted; if not specified, Kepler features refer to both GK104 and GK110.

CUDA Best Practices

The performance guidelines and best practices described in the CUDA C++ Programming Guide [3] and the CUDA C++ Best Practices Guide [4] apply to all CUDA-capable GPU architectures. Programmers must primarily focus on following those recommendations to achieve the best performance.

The high-priority recommendations from those guides are as follows:

Find ways to parallelize sequential code,

Throughout this guide, Fermi refers to devices of compute capability 2.x and Kepler refers to devices of compute capability 3.x. GK104 has compute capability 3.0; GK110 has compute capability 3.5.

The features of GK107 are similar to those of GK104.

³ The features of GK20A (compute capability 3.2) and of GK208 and GK110B (compute capability 3.5) and of GK210 (compute capability 3.7) are similar to those of GK110 except where noted, though numbers of multiprocessors and various throughputs or bandwidths may differ.

- Minimize data transfers between the host and the device,
- Adjust kernel launch configuration to maximize device utilization,
- Ensure global memory accesses are coalesced,
- Minimize redundant accesses to global memory whenever possible,
- Avoid different execution paths within the same warp.

Application Compatibility

Before addressing the specific performance tuning issues covered in this guide, refer to the Kepler Compatibility Guide for CUDA Applications to ensure that your application is being compiled in a way that will be compatible with Kepler.

Note that many of the GK110 architectural features described in this document require the device code in the application to be compiled for its native compute capability 3.5 target architecture (sm 35).

1.4. Kepler Tuning

1.4.1. Device Utilization and Occupancy

Kepler's new Streaming Multiprocessor, called SMX, has significantly more CUDA Cores than the SM of Fermi GPUs, yielding a throughput improvement of 2-3x per clock. 4 Furthermore, GK110 has increased memory bandwidth over Fermi and GK104. To match these throughput increases, we need roughly twice as much parallelism per multiprocessor on Kepler GPUs, via either an increased number of active warps of threads or increased instruction-level parallelism (ILP) or some combination thereof.

Balancing this is the fact that GK104 ships with only 8 multiprocessors, half of the size of Fermi GF110, meaning that GK104 needs roughly the same total amount of parallelism as is needed by Fermi GF110, though it needs more parallelism per multiprocessor to achieve this. Since GK110 can have up to 15 multiprocessors, which is similar to the number of multiprocessors of Fermi GF110, then GK110 typically needs a larger amount of parallelism than Fermi or GK104.

To enable the increased per-multiprocessor warp occupancy beneficial to both GK104 and GK110, several important multiprocessor resources have been significantly increased in SMX:

Kepler increases the size of the register file over Fermi by 2x per multiprocessor; GK210 increases this by a further 2x. On Fermi, the number of registers available was the primary limiting factor of occupancy for many kernels. On Kepler, these kernels can automatically fit more thread blocks per multiprocessor. For example, a kernel using 63 registers per thread and 256 threads per block can fit at most 16 concurrent warps per multiprocessor on Fermi (out of a maximum of 48, i.e., 33% theoretical occupancy). The

⁴ Note, however, that Kepler clocks are generally lower than Fermi clocks for improved power efficiency.

same configuration can fit 32 warps per multiprocessor on GK104 and GK110 (out of a maximum of 64, i.e., 50% theoretical occupancy), or 64 warps (100% theoretical occupancy) on GK210.

- Kepler has increased the maximum number of simultaneous blocks per multiprocessor from 8 to 16. As a result, kernels having their occupancy limited due to reaching the maximum number of thread blocks per multiprocessor will see increased theoretical occupancy in Kepler.
- ▶ GK210 more than doubles the shared memory capacity versus Fermi and earlier Kepler GPUs.

Note that these automatic occupancy improvements require kernel launches with sufficient total thread blocks to fill Kepler. For this reason, it remains a best practice to launch kernels with significantly more thread blocks than necessary to fill current GPUs, allowing this kind of scaling to occur naturally without modifications to the application. The CUDA Occupancy Calculator [5] spreadsheet is a valuable tool in visualizing the achievable occupancy for various kernel launch configurations.

Also note that Kepler GPUs can utilize ILP in place of thread/warp-level parallelism (TLP) more readily than Fermi GPUs can. Furthermore, some degree of ILP in conjunction with TLP is required by Kepler GPUs in order to approach peak single-precision performance, since SMX's warp scheduler issues one or two independent instructions from each of four warps per clock. ILP can be increased by means of, for example, processing several data items concurrently per thread or unrolling loops in the device code, though note that either of these approaches may also increase register pressure.

1.4.2. Managing Coarse-Grained Parallelism

Since GK110 requires more concurrently active threads than either GK104 or Fermi, GK110 introduces several features that can assist applications having more limited parallelism, where the expanded multiprocessor resources described in Device Utilization and Occupancy are difficult to leverage from any single kernel launch. These improvements allow the application to more readily use several concurrent kernel grids to fill GK110:

1.4.2.1. Concurrent Kernels

Since the introduction of Fermi, applications have had the ability to launch several kernels concurrently. This provides a mechanism by which applications can fill the device with several smaller kernel launches simultaneously as opposed to a single larger one. On Fermi and on GK104, at most 16 kernels can execute concurrently; GK110 allows up to 32 concurrent kernels to execute, which can provide a speedup for applications with necessarily small (but independent) kernel launches.



Note: GK20A devices of compute capability 3.2 limit the number of concurrent kernels to 4.

1.4.2.2. Hyper-Q

GK110 further improves this with the addition of Hyper-Q, which removes the false dependencies that can be introduced among CUDA streams in cases of suboptimal kernel launch or memory copy ordering across streams in Fermi or GK104. Hyper-Q allows GK110 to handle the concurrent kernels and/or memory transfers in separate CUDA streams truly independently, rather than serializing the several streams into a single work queue at the hardware level. This allows applications to enqueue work into separate CUDA streams without considering the relative order of insertion of otherwise independent work, making concurrency of multiple kernels as well as overlapping of memory copies with computation much more readily achievable on GK110.

CUDA streams are automatically mapped onto Hyper-Q's multiple hardware work queues via connections to the hardware allocated by the CUDA Driver. While it is possible to allocate more CUDA streams than there are connections, this simply implies that the driver will alias several streams onto some or all of those connections. The CUDA DEVICE MAX CONNECTIONS environment variable can be used to specify the preferred number of connections to be allocated to the driver. The default is 8 (or fewer if CUDA Multi-Process Service is in use); the architectural maximum for GK110 is 32.

CUDA Multi-Process Service (MPS) presents another means by which applications can take advantage of Hyper-Q, wherein several host processes (typically MPI processes) share access to and submit work to the same GPU concurrently, each process receiving some subset of the available connections to that GPU. Using CUDA MPS, processes can achieve overlap of their respective memory transfers and computations with or without the use of CUDA streams, although at the cost of some added latency of work submission and a few other caveats. For more information see the CUDA MPS Overview[6].



Note: Hyper-Q is not supported on GK20A.

Dynamic Parallelism 1.4.2.3.

GK110 also introduces a new architectural feature called Dynamic Parallelism, which allows the GPU to create additional work for itself. A programming model enhancement that leverages this architectural feature was introduced in CUDA 5.0 to enable kernels running on GK110 to launch additional kernels onto the same GPU. Nested kernel launches are done via the same <<<>>> triple-angle bracket notation used from the host and can make use of the familiar CUDA streams interface to specify whether or not the kernels launched are independent of one another. More than one GPU thread can simultaneously launch kernel grids (of the same or different kernels), further increasing the application's flexibility in keeping the GPU filled with parallel work.



Note: Dynamic Parallelism is not supported on GK20A.

1.4.3. Shared Memory and Warp Shuffle

1.4.3.1. Shared Memory Bandwidth

In balance with the increased computational throughput in Kepler's SMX described in Device Utilization and Occupancy, shared memory bandwidth in SMX is twice that of Fermi's SM. This bandwidth increase is exposed to the application through a configurable new 8-byte shared memory bank mode. When this mode is enabled, 64-bit (8-byte) shared memory accesses (such as loading a double-precision floating point number from shared memory) achieve twice the effective bandwidth of 32-bit (4-byte) accesses. Applications that are sensitive to shared memory bandwidth can benefit from enabling this mode as long as their kernels' accesses to shared memory are for 8-byte entities wherever possible.

1.4.3.2. Shared Memory Capacity

Fermi introduced an L1 cache in addition to the shared memory available since the earliest CUDA-capable GPUs. In Fermi, the shared memory and the L1 cache share the same physical on-chip storage, and a split of 48 KB shared memory / 16 KB L1 cache or vice versa can be selected per application or per kernel launch. Kepler continues this pattern and introduces an additional setting of 32 KB shared memory / 32 KB L1 cache, the use of which may benefit L1 hit rate in kernels that need more than 16 KB but less than 48 KB of shared memory per multiprocessor.

Since the maximum shared memory capacity per multiprocessor remains 48 KB on GK104 and GK110, however, applications that depend on shared memory capacity either at a perblock level for data exchange or at a per-thread level for additional thread-private storage may require some rebalancing on these GPUs to improve achievable occupancy. The Warp Shuffle operation for data-exchange uses of shared memory and the <u>Increased Addressable Registers</u> Per Thread as an alternative to thread-private uses of shared memory present two possible alternatives to achieve this rebalancing.

GK210 improves on this by increasing the shared memory capacity per multiprocessor for each of the configurations described above by a further 64 KB (i.e., the application can select 112 KB, 96 KB, or 80 KB of shared memory), which provides automatic occupancy improvements for kernels limited by shared memory capacity.



Note: The maximum shared memory per thread block for all Kepler GPUs, including GK210, remains 48 KB.

1.4.3.3. Warp Shuffle

Kepler introduces a new warp-level intrinsic called the *shuffle* operation. This feature allows the threads of a warp to exchange data with each other directly without going through shared (or global) memory. The shuffle instruction also has lower latency than shared memory access and does not consume shared memory space for data exchange, so this can present an attractive way for applications to rapidly interchange data among threads.

1.4.4. Memory Throughput

1.4.4.1. Increased Addressable Registers Per Thread

GK110 increases the maximum number of registers addressable per thread from 63 to 255. This can improve performance of bandwidth-limited kernels that have significant register spilling on Fermi or GK104. Experimentation should be used to determine the optimum balance of spilling vs. occupancy, however, as significant increases in the number of registers used per thread naturally decreases the warp occupancy that can be achieved, which trades off latency due to memory traffic for arithmetic latency due to fewer concurrent warps.

1.4.4.2. L1 Cache

L1 caching in Kepler GPUs is reserved only for local memory accesses, such as register spills and stack data. Global loads are cached in L2 only (or in the Read-Only Data Cache).

GK110B-based products such as the Tesla K40 GPU Accelerator, GK20A, and GK210 retain this behavior by default but also allow applications to opt-in to the Fermi-style behavior of caching both global and local loads in L1. To select this mode, pass the -xptxas -dlcm=ca flag to nvcc at compile time.

Read-Only Data Cache 1.4.4.3.

GK110 adds the ability for read-only data in global memory to be loaded through the same cache used by the texture pipeline via a standard pointer without the need to bind a texture beforehand and without the sizing limitations of standard textures. Since this is a separate cache with a separate memory pipe and with relaxed memory coalescing rules, use of this feature can benefit the performance of bandwidth-limited kernels. This feature will be automatically enabled and utilized where possible by the compiler when compiling for GK110 as long as certain conditions are met. Foremost among these requirements is that the data must be guaranteed read-only for the duration of the kernel, as the read-only data cache is incoherent with respect to writes. In order to allow the compiler to detect that this condition is satisfied, a necessary (but not always sufficient) condition is that pointers used for loading such data should be marked with both the const and restrict qualifiers. Note that adding these qualifiers where applicable can improve code generation quality via other mechanisms on earlier GPUs as well.

In cases where more explicit control over the read-only data cache mechanism is desired than the const restrict qualifiers provide, or where the code is sufficiently complex that the compiler is unable to detect that the read-only data cache is safe to use, the ldg() intrinsic can be used in place of a normal pointer dereference to force the load to go through the readonly data cache.

Note that the read-only data cache accessed via const restrict is separate and distinct from the constant cache acessed via the constant qualifier. Data loaded through the constant cache must be relatively small and must be accessed uniformly for good performance (i.e., all threads of a warp should access the same location at any given time), whereas data loaded through the read-only data cache can be much larger and can

be accessed in a non-uniform pattern. These two data paths can be used simultaneously for different data if desired

1.4.4.4. Fast Global Memory Atomics

Global memory atomic operations have dramatically higher throughput on Kepler than on Fermi. Algorithms requiring multiple threads to update the same location in memory concurrently have at times on earlier GPUs resorted to complex data rearrangements in order to minimize the number of atomics required. Given the improvements in global memory atomic performance, many atomics can be performed on Kepler nearly as quickly as memory loads. This may simplify implementations requiring atomicity or enable algorithms previously deemed impractical.

Global Memory Bandwidth and GPU Boost 1.4.4.5.

GK110B provides higher memory clocks (and, by extension, higher peak global memory bandwidth) than GK110. For the GK110B-based Tesla K40 GPU Accelerator, while all of the GPU Boost clock settings use the same 3GHz memory clock, the effective memory bandwidth utilization can typically be increased by using the highest boost setting for SM core clocks as well.

1.4.4.6. 2D Memory Copies

The effective bandwidth of cudaMemcpy2D() operations is best when avoiding the use of small device pitches together with large host pitches (>64 KB).

Instruction Throughput

While the maximum instructions per clock (IPC) of both floating-point and integer operations has been either increased or maintained in Kepler as compared to Fermi, the relative ratios of maximum IPC for various specific instructions has changed somewhat. Refer to the CUDA C++ Programming Guide for details.

Single-precision vs. Double-precision

As one example of these instruction throughput ratios, an important difference between GK104 and GK110 is the ratio of peak single-precision to peak double-precision floating point performance. Whereas GK104 focuses primarily on high single-precision throughput, GK110 significantly improves the peak double-precision throughput over Fermi as well. Applications that depend heavily on high double-precision performance will generally perform best with GK110.



Note: GK208 and GK20A devices, like GK104, are single-precision focused and therefore have throughput ratios more similar to GK104 than to GK110.

1.4.6. GPU Boost

NVIDIA GPU Boost is a feature available on the GK110B-based Tesla K40 GPU Accelerator that makes use of power headroom to run the SM core clock to a higher frequency. While the default clock for K40 is set to the base clock, which is necessary for some applications that are demanding on power (e.g., DGEMM), many application workloads are less demanding on power and can take advantage of a higher boost clock setting for added performance.

GK210 further improves GPU Boost functionality. For example, the Tesla K80 GPU Accelerator provides many more possible boost clock settings than Tesla K40, and Tesla K80 defaults to a dynamic boost mode, wherein power headroom is leveraged by increasing clock speeds automatically.

GPU Boost clocks (and optional disabling of dynamic boost for GK210) can be selected through nvidia-smi or NVML.

1 4 7 Multi-GPU

NVIDIA's Tesla K10 GPU Accelerator is a dual-GK104 board; the Tesla K80 GPU Accelerator is a dual-GK210 board. As with other dual-GPU NVIDIA boards, the two GPUs on these board will appear as two separate CUDA devices; they have separate memories and operate independently. As such, applications that will target the Tesla K10 or K80 GPU Accelerator but that are not yet multi-GPU aware should begin preparing for the multi-GPU paradigm. Since dual-GPU boards appear to the host application exactly the same as two separate single-GPU boards, enabling applications for multi-GPU can benefit application performance on a wide range of systems where more than one CUDA-capable GPU can be installed.

1.4.8. PCle 3.0

Kepler's interconnection to the host system has been enhanced to support PCIe 3.0. For applications where host-to-device, device-to-host, or device-to-device transfer time is significant and not easily overlapped with computation, the additional bandwidth provided by PCIe 3.0, given the requisite host system support, over the earlier PCIe 2.0 specification supported by Fermi GPUs should boost application performance without modifications to the application. Note that best PCIe transfer speeds to or from system memory with either PCIe generation are achieved when using pinned system memory.



Note: In the Tesla K10 and K80 GPU Accelerator products, the two GPUs sharing a board are connected via an on-board PCIe 3.0 switch. Since these GPUs are also capable of GPUDirect



Peer-to-Peer transfers, the inter-device memory transfers between GPUs on the same board can run at PCIe 3.0 speeds even if the host system supports only PCIe 2.0 or earlier.



Note: While the Kepler architecture is compliant with the PCIe 3.0 specification, not all Keplerbased products support PCIe 3.0 speeds. For example, while Tesla K10, K40, and K80 support PCIe 3.0. Tesla K20 and K20X do not.



Note: PCIe 3.0 throughputs may be improved in some circumstances by using the highestavailable GPU Boost clock.

Warp-synchronous Programming 1.4.9.

As a means of mitigating the cost of repeated block-level synchronizations, particularly in parallel primitives such as reduction and prefix sum, some programmers exploit the knowledge that threads in a warp execute in lock-step with each other to omit syncthreads () in some places where it is semantically necessary for correctness in the CUDA programming model.

The absence of an explicit synchronization in a program where different threads communicate via memory constitutes a data race condition or synchronization error. Warp-synchronous programs are unsafe and easily broken by evolutionary improvements to the optimization strategies used by the CUDA compiler toolchain, which generally has no visibility into crossthread interactions of this variety in the absence of barriers, or by changes to the hardware memory subsystem's behavior. Such programs also tend to assume that the warp size is 32 threads, which may not necessarily be the case for all future CUDA-capable architectures.

Therefore, programmers should avoid warp-synchronous programming to ensure futureproof correctness in CUDA applications. When threads in a block must communicate or synchronize with each other, regardless of whether those threads are expected to be in the same warp or not, the appropriate barrier primitives should be used. Note that the Warp Shuffle primitive presents a future-proof, supported mechanism for intra-warp communication that can safely be used as an alternative in many cases.

Appendix A. References

[1] NVIDIA GeForce GTX 680: The fastest, most efficient GPU ever built.

http://www.geforce.com/Active/en US/en US/pdf/GeForce-GTX-680-Whitepaper-FINAL.pdf

[2] NVIDIA's Next Generation CUDA Compute Architecture: Kepler GK110.

http://www.nvidia.com/content/PDF/kepler/NVIDIA-Kepler-GK110-Architecture-Whitepaper.pdf

[3] CUDA C Programming Guide.

http://docs.nvidia.com/cuda/cuda-c-programming-guide/

[4] CUDA C Best Practices Guide.

http://docs.nvidia.com/cuda/cuda-c-best-practices-quide/

[5] CUDA Occupancy Calculator spreadsheet.

http://developer.download.nvidia.com/compute/cuda/CUDA Occupancy calculator.xls

[6] Sharing A GPU Between MPI Processes: Multi-Process Service (MPS) Overview.

http://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf

Appendix B. Revision History

Version 0.9

CUDA 5.0 Preview Release

Version 1.0

- Added discussion of ILP vs TLP (see Device Utilization and Occupancy).
- Expanded discussion of cache behaviors (see Memory Throughput).
- Added section regarding Warp-synchronous Programming.
- Added section regarding <u>2D Memory Copies</u>.
- Minor corrections and clarifications.

Version 1.1

Clarified const restrict discussion and mentioned ldg() intrinsic in Read-Only Data Cache.

Version 1.2

- \triangleright Add references to GK110B, which allows an opt-in to the caching of global loads in the <u>L1</u> Cache and enables higher clock speeds via GPU Boost.
- Expand discussion of ILP in <u>Device Utilization and Occupancy</u>.
- Expand discussion of Hyper-Q, adding mention of CUDA DEVICE MAX CONNECTIONS and CUDA Multi-Process Service (MPS).
- Clarification of PCIe 3.0 support.
- Add hyperlinks to all endnote references.

Version 1.3

Add references to GK210, which increases register file and shared memory capacities and enables additional GPU Boost modes versus GK110B.

Version 1.4

▶ Updated references to the CUDA C++ Programming Guide and CUDA C++ Best Practices

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

© -2021 NVIDIA Corporation & affiliates. All rights reserved.

