

streaming processor

Related terms:

[Soft Error](#), [Graphics](#), [Kernel](#), [Memory Access](#), [Multiprocessors](#), [Vulnerabilities](#)

[View all Topics](#)

Learn more about streaming processor

Adaptive sparse matrix representation for efficient matrix-vector multiplication

P. Zardoshti, ... H. Sarbazi-Azad, in [Advances in GPU Research and Practice](#), 2017

3.1 Hardware Architectures

The GPU architecture consists of several hundred simple cores called [streaming processors](#) (SPs), which are grouped in a set of streaming [multiprocessors](#) (SMs). Each SM [executes instructions](#) in a [single-instruction multiple-threads](#) (SIMT) mode and supports a [multithreading](#) execution mechanism. GPUs are capable of utilizing large amounts of [memory bandwidth](#). To further increase usable bandwidth, modern GPUs also contain a series of on-chip caches. In modern GPU architecture, each thread has its own private register. For memory, a thread has its own memory space, which is called the local memory. A block also has its own 64 KB on-chip memory that can be configured as shared memory and [L1 cache](#). All threads in a block can access the same shared memory and L1 cache, while threads in other blocks cannot. The [L2 cache](#) is the primary point of data [unification](#) between the SM units, servicing all load, store, and texture requests and providing efficient, high-speed [data sharing](#) across the GPU. The entire [kernel](#) also has a global memory. All threads in a block can access the global memory space. Also, all threads can access read-only constant memory and [texture memory](#) spaces [35].

NVIDIA GeForce GTX 480

NVIDIA GeForce GTX 480

The NVIDIA GPU is based on the GF100 core with 480 cores. It works with a clock frequency of 1.4 GHz, and has 1.5 GB GDDR5 with a 177 GB/s bandwidth. This GPU consists of 15 SMs. Each SM has a memory space that can be configured as 48K of shared memory with 16K of L1 cache, or as 16K of shared memory with 48K of L1 cache [35].

NVIDIA Tesla K20X

There are 2688 cores with a clock frequency of 0.73 GHz in the Kepler GK110 architecture. It has 6 GB GDDR5 with 250 GB/s bandwidth. In addition, Kepler allows additional flexibility in configuring the memory space of SMs by splitting a 32K/32K split between allocation of shared memory and L1 cache [36]. Table 2 summarizes the specifications of these architectures.

Table 2. Architectural Analysis of Evaluated Platforms

Architecture	Architecture	GTX 480 GeForce	Tesla K20X	GTX 480 GeForce
Model	Model	GF100	GK100	GF100
Core	Core	480	2688	480
Clock (GHz)	Clock (GHz)	1.4	0.73	1.4
DP Peak (GFlop/s)	DP Peak (GFlop/s)	168	1310	168
SP Peak (GFlop/s)	SP Peak (GFlop/s)	1345	3950	1345
Memory bandwidth (GB/s)	Memory bandwidth (GB/s)	177.4	250	177.4

[Read full chapter](#)

GPU Acceleration of Rarefied Gas Dynamic Simulations

Aldo Frezzotti, ... Aldo Frezzotti, [GPU Computing: GPU Code Editing](#), 2012, Jade Edition, 2012

15.4.2 Performance Evaluation

The performance analysis has been performed on a commercially available GPU GeForce GTX 260 using the GTX 260 model. The GTX 260 GPU model consists of 24 streaming multiprocessors with 8 cores (SM) clocked at 1.242 GHz each for a peak performance of 115.4 GFlop/s in single precision. The GPU has 896 MB of device memory with a memory bandwidth of 111.9 GB/s. The graphic processing has been hosted by a personal computer equipped with 4

GB of [main memory](#) with a memory bandwidth of 12.8 GB/s and an Intel Core Duo Quad Q9300 CPU, running at 2.5 GHz for a peak theoretical performance of a single core of 20 GFLOP/s. The host code has been compiled using the gcc/g++ [compiler](#) with optimization option “-O3” whereas no SSE instructions have been used.

GB of [main memory](#) with a memory bandwidth of 12.8 GB/s and an Intel Core Duo Quad Q9300 CPU, running at 2.5 GHz for a peak theoretical performance of a single core of 20 GFLOP/s. The host code has been compiled using the gcc/g++ [compiler](#) with optimization option “-O3” whereas no SSE instructions have been used.

The efficiency of the BHS and BGK codes has been assessed by computing the [speedup factor](#) $S = T_{CPU}/T_{GPU}$, where T_{CPU}/T_{GPU} are the times used by the CPU and GPU, respectively. It is measured after initial setup and do not include the initial and final iterations as executed to transfer data between the disjoint GPU and CPU by [spdc65](#). Figure 15.2 shows the time in seconds which is spent in the streaming step and only the [collision](#) step by the BHS code and the BGK code. For the BGK code, the collision step is more time consuming than the streaming step, while taking at most 31% of overall computing time. This is in contrast to the BGK code where the evaluation of the collision integral is [computationally demanding](#).

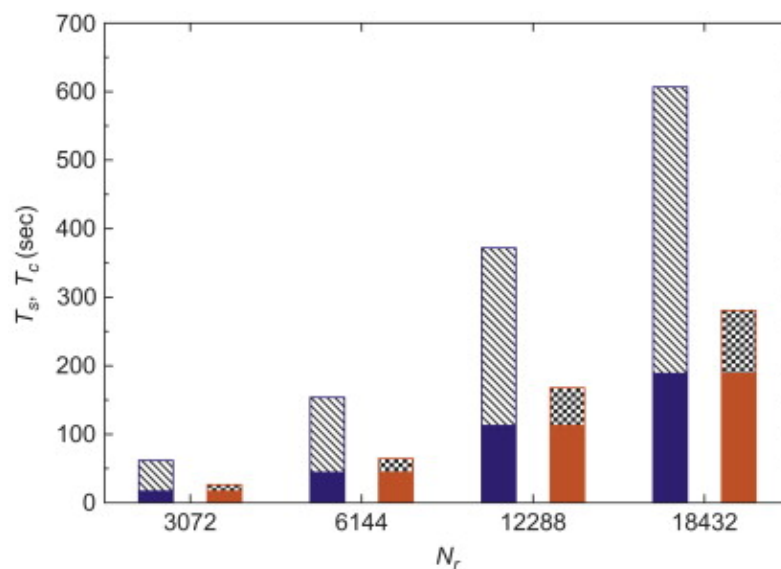
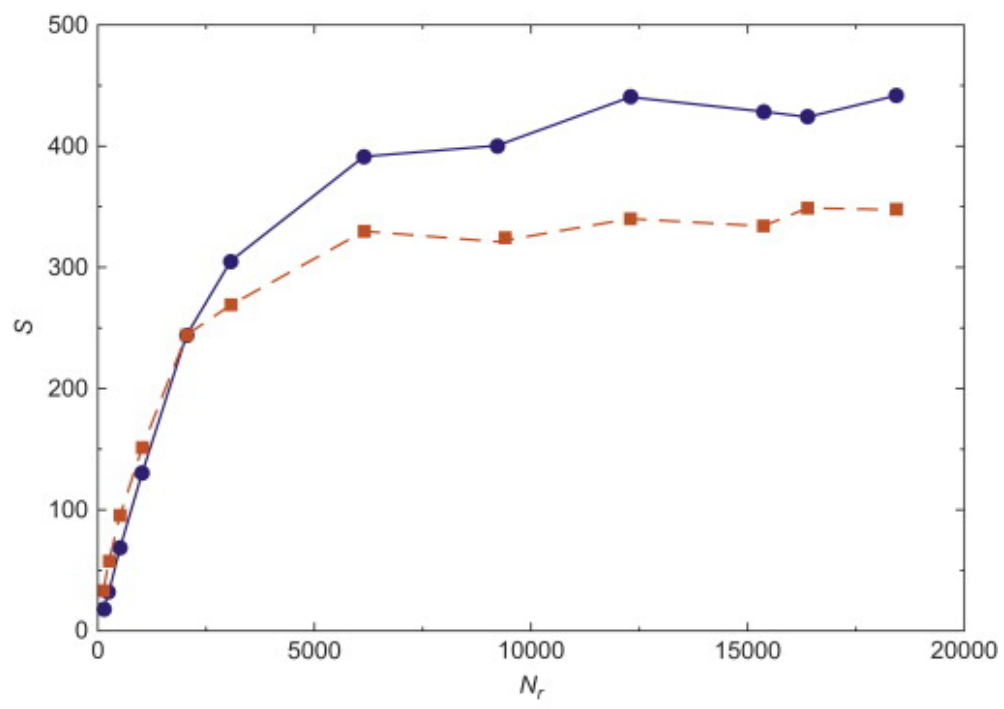


Figure 15.2. Time in seconds spent on streaming step (dark bar) and on collision step for the BHS code (solid bar) and BGK code (striped bar). $N_v = 8000$.

Figure 15.3 reports the obtained speedup factors as a function of the number of spatial grid points, N_r . With the number of cells, till N_r approximately exceeds 10,000, the speedup factor speeds up about 440 and 350 for the BHS and the BGK code, respectively. As shown by the speed-up curves, the GPU power is not fully exploited till the number of cores or threads reaches a threshold. Beyond this point, the speedup computation time approximately behaves as a [linear function](#) of N_r .



power. Altogether, the codes which solve the BHS and BGKW equations show a performance of 100 and 75 GFLOP/s, respectively.

power. Altogether, the codes which solve the BHS and BGKW equations show a performance of 100 and 75 GFLOP/s, respectively.

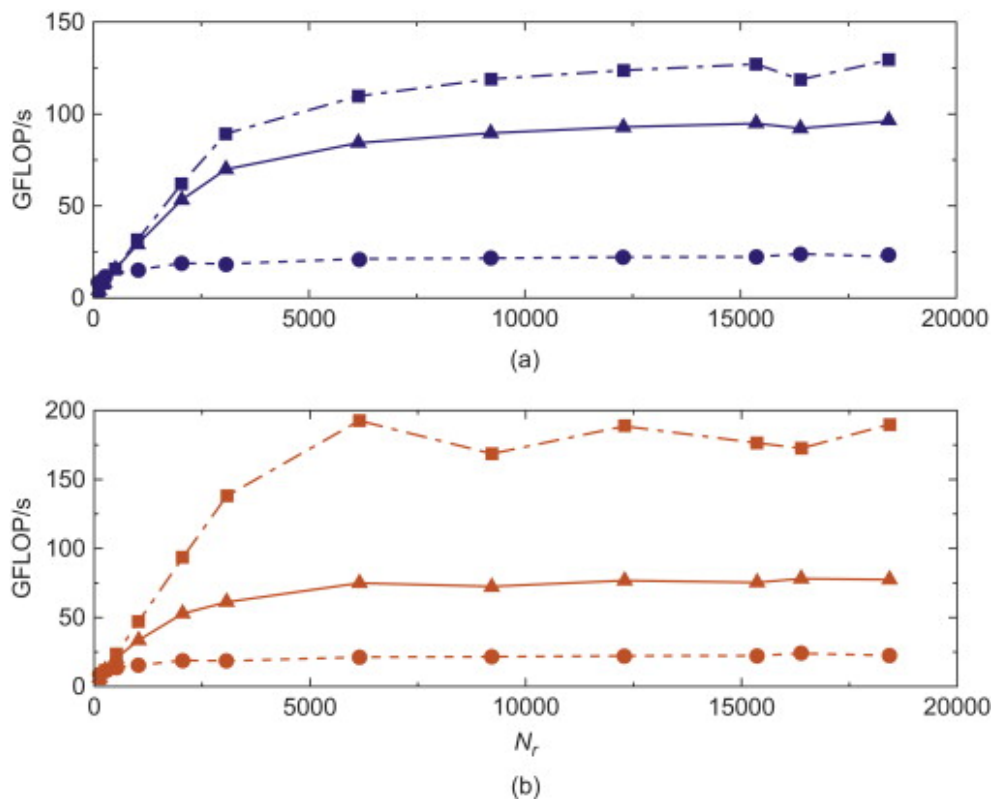


Figure 15.4. Real GFLOP/s vs. N_r for the BHS and BGKW codes. (a) BHS code with $N_t = 64$. (b) BGKW code with $N_t = 64$. Dashed lines with circles: streaming step. Dashed lines with squares: collision step. Solid lines with triangles: overall code.

A similar analysis can be performed for the GPU-based codes. The time required by the CPU to execute the floating-point operations is however not negligible with respect to the time required for the transfer of data to/from the memory and, hence, it has to be taken into account in the estimation of the GFLOP/s. The number of accesses to the global memory is supposed to be the same as above whereas the CPU bandwidth and the computing performance per single core are about 12.8 GB/s and 20 GFLOP/s, respectively. Therefore the sequential codes which solve the BHS and BGKW equations deliver 9.9 GFLOP/s and 11.1 GFLOP/s, respectively, with the effective performance which are about 0.28 GFLOP/s and 0.29 GFLOP/s, respectively. This suggests that the high speed-up factors shown in Figure 15.3 are mainly due to the fact that the sequential codes are not fully optimized. However, it is worth noticing that the ideal analysis discussed above is also valid for the GPU. For instance, as indicated by benchmarks, the GPU effective bandwidth is approximately one-half of the main memory bandwidth whereas for the CPU case for the GPU. Accordingly, the absolute GFLOP/s of the GPU-based codes should be halved. Moreover, the performance of the sequential codes developed here are comparable

to the performance of similar codes described in literature. For instance, in Ref. [9] the solution of the cavity flow problem is obtained by integrating the linearized BGKW equation in about 3 hours. Using the execution time of this sequential code to compute the speed-up factor yields a result as high as 265. It is worth noticing that in Ref. [9] the number of velocity variables has been reduced by a standard projection method [1]. The value of the speed-up factor given above has been obtained by modifying the parallel code to take advantage of the same reduction technique. To the best of the authors' knowledge, no direct solutions of the cavity flow problem based on the BHS equation have been provided until now. However, a similar reasoning can be applied to the one-dimensional unsteady Couette flow. In Ref. [11] its solution is achieved in approximately 9 hours, which gives a speed-up factor of the BHS code of about 54. This results is still lower than the speed-up factor reported in Figure 15.3, but this is due to the fact that the discretizations used to provide solutions as accurate as the ones reported in Ref. [11] does not allow to fully exploit the computational power of the GPU.

to the performance of similar codes described in literature. For instance, in Ref. [9] the solution of the cavity flow problem is obtained by integrating the linearized BGKW equation in about 3 hours. Using the execution time of this sequential code to compute the speed-up factor yields a result as high as 265. It is worth noticing that in Ref. [9] the number of velocity variables has been reduced by a standard projection method [1]. The value of the speed-up factor given above has been obtained by modifying the parallel code to take advantage of the same reduction technique. To the best of the authors' knowledge, no direct solutions of the cavity flow problem based on the BHS equation have been provided until now. However, a similar reasoning can be applied to the one-dimensional unsteady Couette flow. In Ref. [11] its solution is achieved in approximately 9 hours, which gives a speed-up factor of the BHS code of about 54. This results is still lower than the speed-up factor reported in Figure 15.3, but this is due to the fact that the discretizations used to provide solutions as accurate as the ones reported in Ref. [11] does not allow to fully exploit the computational power of the GPU.

We conclude by observing that the best problem examined here has shown that the size of physical meshes is toward the application and the application to complex two- or three-dimensional flows, direct numerical flows, for the trending capability.

> Read full chapter

Advances in COMPUTERS

Jack Dongarra, ... Mustafa Tili, Computers, 2008

3.3.3.1 Molecular Dynamics.

As with the Cell, in the GPU the algorithm that calculates new accelerations from the locations of the atoms and several constants. For our constants, then, the choice is to have one input array comprising the positions and comprising the new accelerations. The accelerations were compiled into the source using the provided JIT compilation program initialization.

We set up the GPU code to compute the acceleration for each location in the output array. That is, each shader program has to compute an acceleration for one atom by checking for all other atoms and by accumulating forces into a single acceleration for the target atom. The target length limits prevent us from searching for the target atom in a single pass, and so with 4096 atoms on the GPU, the algorithm must pass through the input array. After the GPU is finished, the resulting accelerations are read back

into [main memory](#), where the host CPU proceeds with the current time step. At the next time step, the updated positions are re-sent to the GPU and new accelerations computed again.

into [main memory](#), where the host CPU proceeds with the current time step. At the next time step, the updated positions are re-sent to the GPU and new accelerations computed again.

Figure 13 shows performance of NVIDIA GeForce 7900GTX GPU. This figure includes results from the host CPU (a 2.4 GHz Opteron) as a reference for scaling comparison. The graph shows the relative performance of the GPU below 1024 atoms shows the point at which the overhead of offloading this acceleration computation to the GPU becomes a significant fraction of the total runtime than the $O(N^2)$ calculation itself. The cost of $O(N)$ costs for each time step includes sending the position and acceleration array across the PCIe bus every time step, and every time step, and that the results show that there is a [lower bound](#) on problem size where a GPU will be faster than a CPU. However, the massive [parallelism](#) of the GPU helps it to maintain a speedup above 1000 atoms.

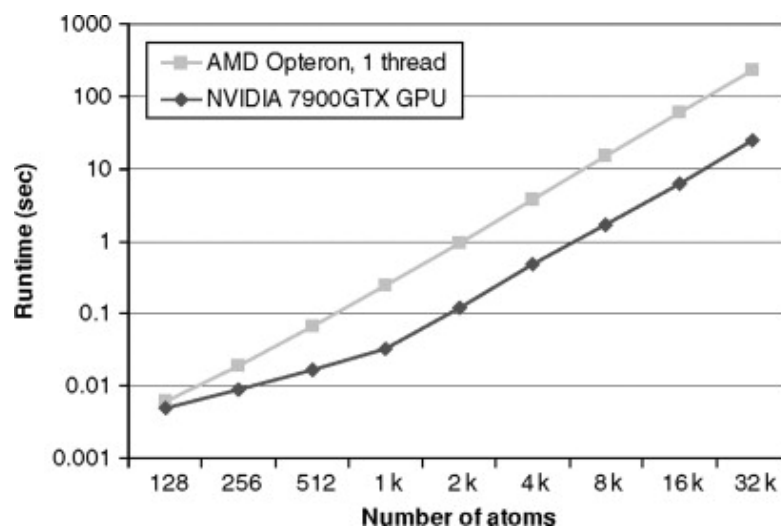


Fig. 13. Performance scaling of GPU with CPU for comparison.

[Read full chapter](#)

GPU alignment of two and three sequences

J. Li, ... S. Sahni, in [Algorithms and GPU Research and Practice](#), 2017

2 GPU architecture

The GPU [algorithm](#) in this chapter targets the NVIDIA Tesla C2050 GPU, which is also known as the C2050 Fermi. The C2050 has 44 multiprocessors (SMs) and each SM has 32 cores (SPs) giving the C2050 a total of 1408 cores. The GPU architecture is shown in Figure 14.

448 SPs or cores. Fig. 1 shows the architecture of a C2050 SM. Although each SP of a C2050 has its own integer, [single- and double-precision](#) units, the 32 SPs of an SM share 4 single-precision [transcendental function](#) units. An SM has 64 KB of on-chip memory that can be “configured as 48 KB of shared memory with 16 KB of [L1 cache](#) (default setting) or as 16 KB of shared memory with 48 KB of L1 cache” [25]. Additionally, there are 32K 32-bit registers per SM and 3 GB of off-chip [device/global memory](#) that is shared by all 14 SMs. The [peak performance](#) of a C2050 is 1288 GFlops (or 1.288 TFlops) of single-precision operations, and 515 GFlops of double-precision operations and the [power consumption](#) is 238 W [26]. Once again, the peak of 1288 GFlops requires that MADDs and SF instructions be dual-issued. When there are MADDs alone, the peak single-precision rate is 1.03 GFlops.

448 SPs or cores. Fig. 1 shows the architecture of a C2050 SM. Although each SP of a C2050 has its own integer, [single- and double-precision](#) units, the 32 SPs of an SM share 4 single-precision [transcendental function](#) units. An SM has 64 KB of on-chip memory that can be “configured as 48 KB of shared memory with 16 KB of [L1 cache](#) (default setting) or as 16 KB of shared memory with 48 KB of L1 cache” [25]. Additionally, there are 32K 32-bit registers per SM and 3 GB of off-chip [device/global memory](#) that is shared by all 14 SMs. The [peak performance](#) of a C2050 is 1288 GFlops (or 1.288 TFlops) of single-precision operations, and 515 GFlops of double-precision operations and the [power consumption](#) is 238 W [26]. Once again, the peak of 1288 GFlops requires that MADDs and SF instructions be dual-issued. When there are MADDs alone, the peak single-precision rate is 1.03 GFlops.

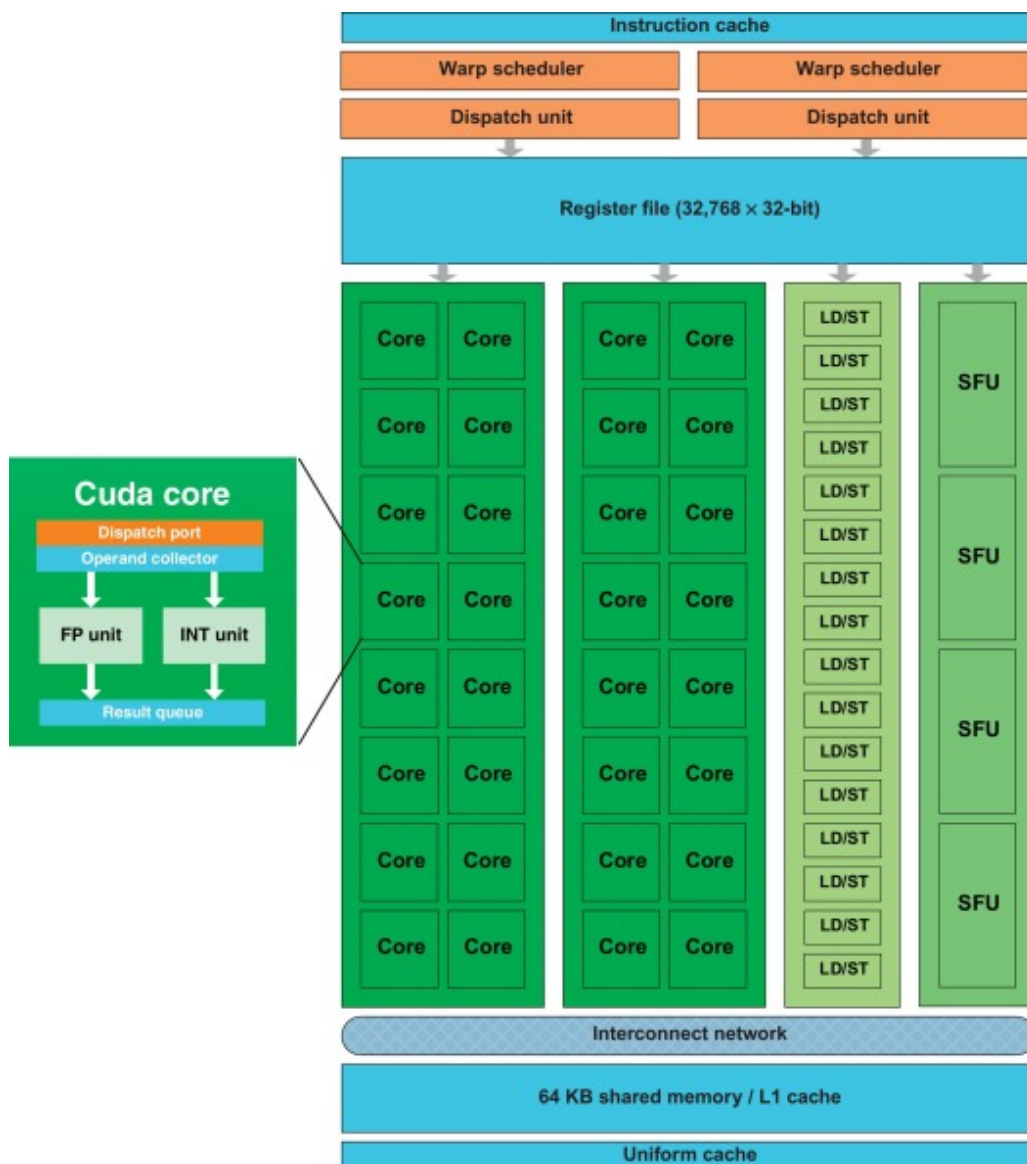


Fig. 1. Architecture of the NVIDIA SM of [27] NVIDIA Fermi [27].

In NVIDIA parlance, the **compute capability** of the C2050 is 2.0. The key challenge in deriving high performance for GPUs is to **effectively minimize the memory traffic** between the SMs and between the SMs and the GPU. Data that

is used repeatedly should go to [registers or shared memory](#), while data that is used less frequently but of larger size should go to device memory.

is used repeatedly should go to [registers or shared memory](#), while data that is used less frequently but of larger size should go to device memory.

> [Read full chapter](#)

GPGPU: General-Purpose Computing on the GPU

Ying Tan, in [Gpu-Based Parallel Systems](#), 2016
Ying Tan, in [Gpu-Based Parallel Systems](#), 2016

2.3.4 Single-Instruction, Multiple-Thread (SIMT) Thread (SIMT)

A CUDA-enabled GPU can have multiple streaming multiprocessors (SMs), which are equivalent to GPU cores. Each SM can have certain number of [separate processors \(cores, SPs\)](#) with respect to the [specific architecture](#).

When a CUDA program invokes a host API to launch a kernel, all blocks are distributed equally to the SMs with available SMs with capacity. The threads of a thread block execute concurrently on multiple processors in the multiprocessor as a unit, and [multiple thread blocks](#) can execute on multiple processors. As running blocks finish, the blocks of inactive threads are created on the vacated SMs, as illustrated in Fig. 2.9.

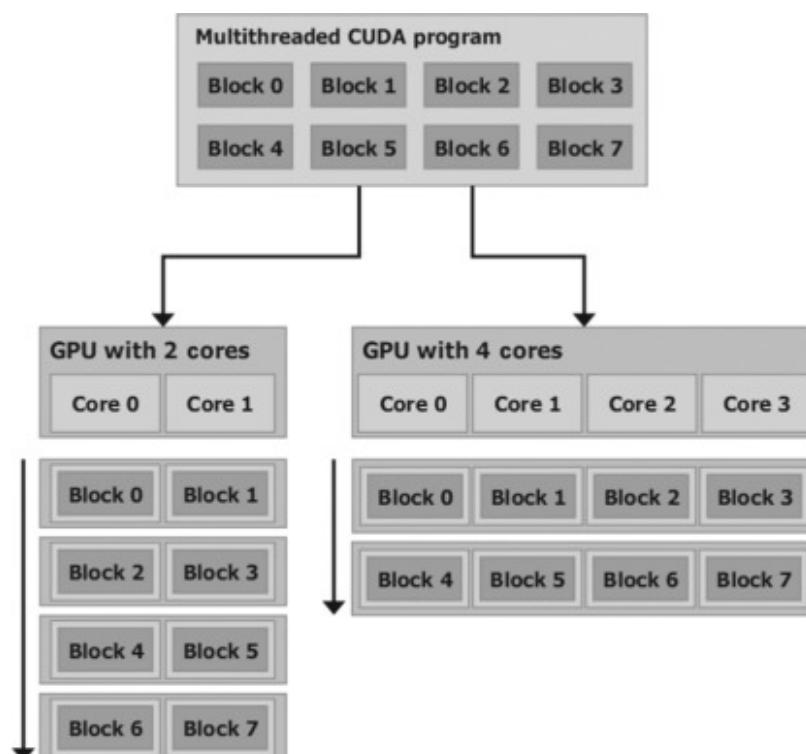


Fig. 2.9. Automatic Scalability

Fig. 2.9. Automatic Scalability

To manage such a large amount of threads, the employment of thread architecture called **single-instruction, multiple-thread (SIMT)** or **multiple-thread (SIMT)**.

The multiprocessor threads, processes, schedules, and executes threads in groups of 32 parallel threads called **warps**. In a single warp, all threads start together at the same program address, but they have different **instruction address** counter and register state and are therefore executed independently.

When a multiprocessor is given a block of code to execute, it partitions them into warps and each warp gets scheduled by a scheduler for execution. The way a block is partitioned into warps is in the same way as the way a warp contains threads of consecutive threads, with the first warp containing thread 0.

A warp executes on a processor core. A warp can be efficient or inefficient. Efficiency is realized when all 32 threads in a warp agree on their threads of path. If threads of a warp diverge via a data-dependent branch, the warp serially executes each branch path taken, disabling the threads that path, and when all paths complete, the threads are re-enabled to re-execute back to the point of divergence. **Branch divergence** occurs only within a warp, and different warps can execute independently regardless of whether they are executing the same or disjoint code paths.

> [Read full chapter](#)

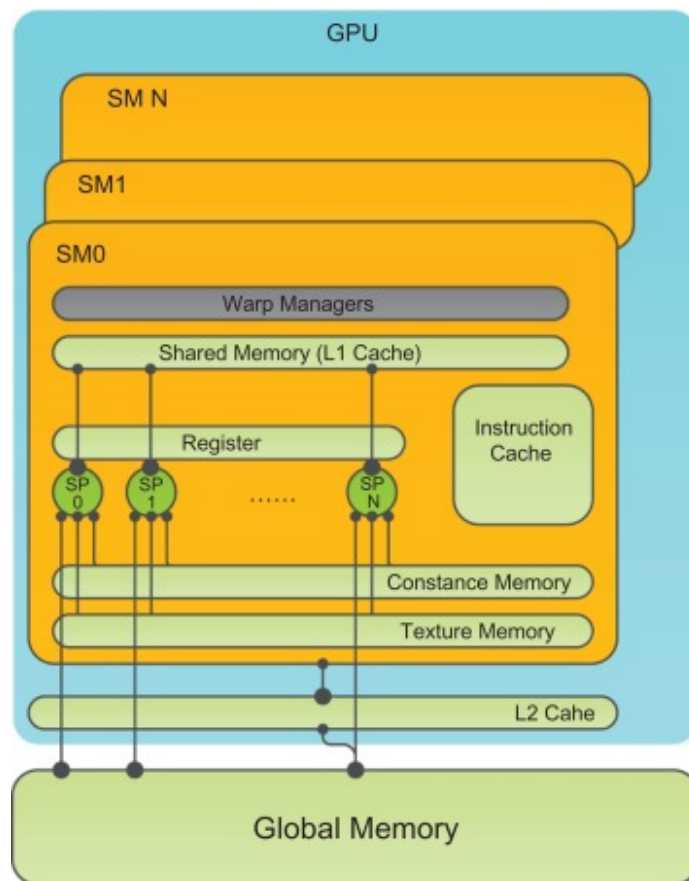
Deep Learning and Its Parallelization

X. Li, ... W. Zheng, X. Li, W. Zheng, in [Big Data](#), 2016

The architecture of GPU

GPU is connected to the host PC Express (PCIe) and it has its own device memory with capacity up to several gigabytes in GPU devices. GPU hardware mainly consists of **memory, streaming (SMs), and processors (SPs)**.

GPU is an array of SMs that consist of SPs (Fig. 10). Each SM can execute in parallel with the other SMs. Each SM consists of 32 SPs and 2880 cores and each core can execute a sequential thread (single-instruction, multiple thread). Each SM is associated with a private local memory (L1 cache), read-only texture memory, registers, and each L2 cache. Each L2 cache has a fully pipelined integer ALU, a floating point unit, a compare unit, and a branch unit [17,18].



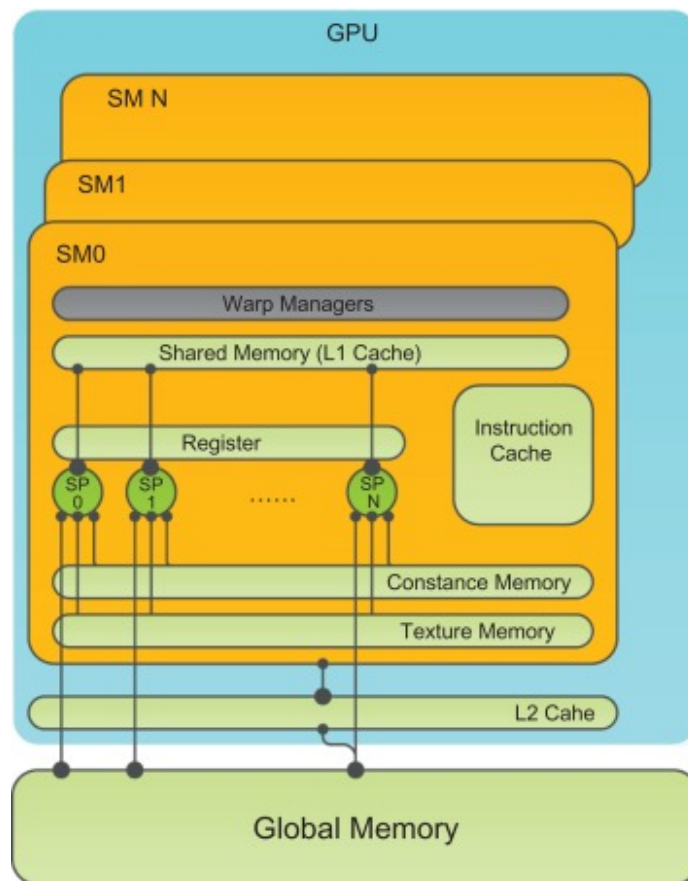


Fig. 10. A simplified diagram of GPU hardware.

Each thread can access registers and shared memory at a high speed in parallel, and we also refer to shared memory as Register memory. Registers are private to individual threads, and each thread can only access its own registers. Thirty-two threads are grouped into a warp, which is scheduled by a warp manager, and the threads in a warp are executed in parallel. The GPU architecture can be enhanced by extending the number of SMs and memory resources.

> [Read full chapter](#)

GPU programming

Gerassimos Barlas, Gerassimos Barlas, GPU Programming, 2015

6.3 CUDA'S execution model: multiple processors and warps

GPU cores are essentially vector processors, capable of applying the same instruction on a large collection of data in parallel. So, when a GPU is run on a GPU core, the same instruction is repeatedly executed by a large collection of processing units called SIMD. A group of SIMDs is called a warp, and the execution of

a single control unit is called a streaming [multiprocessor](#), or SM. A GPU can contain multiple SMs, each running each own kernel. Since each thread runs on its own SP, we will refer to SPs as cores (Nvidia documentation calls them CUDA cores), although a more purist approach would be to treat SMs as cores. Nvidia calls this [execution model Single-Instruction, Multiple Threads](#) (SIMT).

[cution context](#) that is maintained on-chip. This contradicts the arrangement used by [multithreading](#) on CPUs, where a very expensive context switch (involving the saving of CPU registers) accompanies thread switching.

Concurrent kernels/device

of a conditional operation leads them to different paths? The answer is that all the divergent paths are evaluated (if threads branch into them) in sequence until the paths merge again. The threads that do not follow the path currently being executed are stalled. So, given the kernel in Listing 6.4, the execution of a warp would be as shown in Figure 6.4.

of a conditional operation leads them to different paths? The answer is that all the divergent paths are evaluated (if threads branch into them) in sequence until the paths merge again. The threads that do not follow the path currently being executed are stalled. So, given the kernel in Listing 6.4, the execution of a warp would be as shown in Figure 6.4.

```
doSmtElse();
}
doFinal();
}
```

Listing 6.4. An example of a kernel that would stall half of the SMs in half of the threads in a warp stalled.

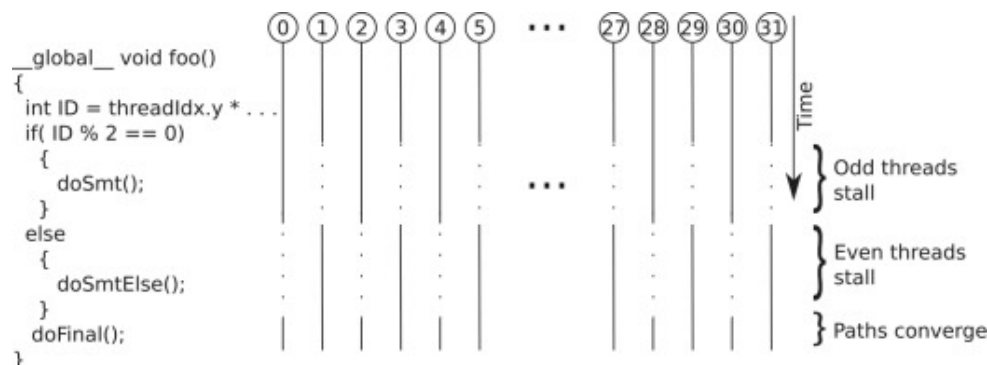


Figure 6.4. An illustration of the execution of the kernel in Listing 6.4 by a warp of threads. The dotted lines in the dotted lines indicate a stall.

```
__global__ void foo()
{
    int ID = threadIdx.y * blockDim.x + threadIdx.x;
    if( ID % 2 == 0)
    {
        doSmt();
    }
    else
    {
        doSmtElse();
    }
    doFinal();
}
```

A solution to this problem is discussed in Section 6.7.2.

[> Read full chapter](#)

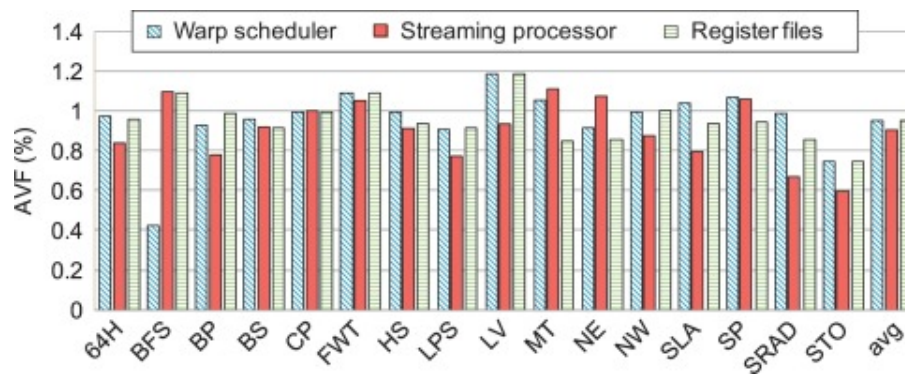
Addressing hardwareability challenges in general-purpose GPUs

J. Tan, X. Fu, in [Advances in GPU Research and Practice](#), 2017

Dynamic warp formation

Branch divergence is a major cause for performance degradation in GPGPUs. As we discussed earlier, the immediate (PDOM) or late (LDOM) lacks the capability to reconverge threads at the beginning of the next iteration. To further improve

the performance. DWF is proposed in Ref. [24] to efficiently handle the [threads' divergence](#). It groups threads from multiple warps but branching to the same target into a new and complete warp, and issues it into the SIMD pipeline. Therefore the parallel [streaming processors](#) in the SM are fully utilized, and their performance is enhanced.



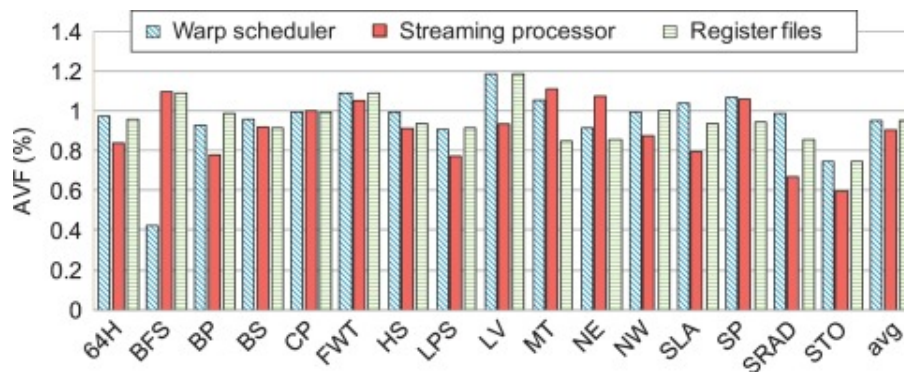


Fig. 6. The normalized AVF of GPU processing microarchitecture structures under the impact of DWF.

The streaming processor's AVF is highly related to the kernel execution time. Since the instruction completion time is determined by the instruction type, a shorter kernel execution time implies a higher quantity of ACE instructions per cycle in the streaming processor. And the streaming processors are more susceptible to soft errors when the kernel executes intuitively, DWF would hurt the streaming processor's soft error susceptibility because it targets on optimizing performance and reduces the kernel execution time. Interestingly, DWF decreases streaming processor's AVF by 10% on average in Fig. 6. This is because DWF decreases the benchmarks (as of benchmarks) (as shown in Fig. 7). When randomly reordering instructions, DWF may lose opportunities to combine memory accesses originally from the same warp, and introduce extra off-chip memory accesses, which negatively affect the performance. When the benefit of DWF in reducing the kernel execution time is outweighed by the increased memory access requests, performance starts to degrade.

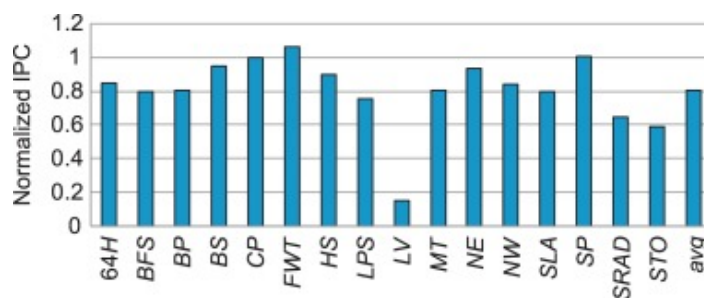


Fig. 7. The normalized IPC under DWF.

> [Read full chapter](#)

An Introduction to OpenCL™

David B. Kirk, Wen-mo W. Hwu, [Read full chapter](#)
 (Second Edition), 2013

14.3 Device Architecture

14.3 Device Architecture

Like CUDA, OpenCL is a heterogeneous computing system as a host and one or more OpenCL devices. The host is a general-purpose CPU that executes the host program. Figure 14.2 shows the structure of an OpenCL device. Each device consists of one or more units (CUs) that correspond to CUDA streaming multiprocessors (SMs); however, a CU also corresponds to CPU cores or other types of execution units such as DSPs and FPGAs.

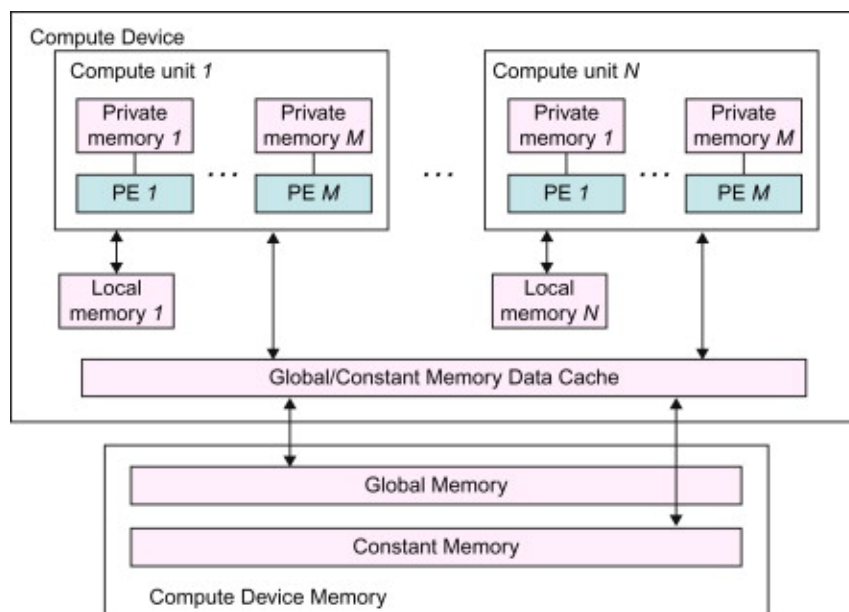


Figure 14.2. Conceptual OpenCL device architecture.

Each compute unit, in turn, consists of one or more processing elements (PEs), which corresponds to the corresponding stream processors (SPs) in CUDA. A computation on a device ultimately happens in individual PEs.

Like CUDA, OpenCL also exposes a variety of memory types that can be used by programmers. Figure 14.2 illustrates the types: global, constant, local, and private. Table 14.3 summarizes the supported memory types and the mapping of these OpenCL memory types to CUDA memory types. The OpenCL global memory corresponds to the CUDA global memory. Unlike CUDA, the global memory can be dynamically allocated by the host program and supports read/write access by both host and devices.

Table 14.3. Mapping of OpenCL Memory Types to CUDA Memory Types

Memory Type	Memory Access	Device Access	CUDA Equivalent
Global memory	Global location; read/write access	No allocation; read/write access by all work items in all work groups, large and slow,	Global allocation; read/write access

		but may be cached in some devices	
Constant memory	Constant allocation; read/write access	Static allocation; read-only access by all work items	Dynamic allocation; read/write access
Local memory	Dynamic allocation; no access	Static allocation; shared read/write access by all work items in a work group	Dynamic allocation; no access
Private memory	Static allocation; no access	Static allocation; read/write access by a single work item	Register and local memory

Unlike CUDA, the OpenCL local memory is not dynamically allocated by the host. Like CUDA, the constant memory supports read-only access by the host and read-only access by devices. To support this, OpenCL provides a device query that returns the constant memory size supported by the device.

The mapping of OpenCL local memory to CUDA memory types is more interesting. The OpenCL local memory corresponds to CUDA shared memory. The OpenCL local memory is dynamically allocated by the host or statically allocated on the device. In the latter case, the OpenCL local memory is accessed by the host and it supports shared read/write access by all work items. The private memory of OpenCL corresponds to the CUDA automatic variables.

> [Read full chapter](#)