CUDA设备与架构

CUDA设备的特点

CPU和GPU之间的协作

Explicit Resource Allocation. 与共享内存模型不同,CUDA需要程序员显式地在GPU上分配内存和计算资源。这种显式控制允许通过利用GPU在并行计算中的优势来微调性能。

Program Entry and Control Flow. CPU(通常称为主机)负责程序的整体控制流。它初始化数据、分配内存并在GPU上启动kernels。GPU完成任务后,将状态报告给CPU,CPU根据结果继续控制流。

Synchronization and Status Reporting. 在启动kernel后,CPU可以使用 cudaDeviceSynchronize 与GPU同步,以确保kernel完成执行,然后继续进行。

GPU内部执行

Kernel Execution. Kernel是用CUDA C/C++编写的函数,在GPU上执行。Kernel由CPU启动,指定thread blocks和每块的threads数。

Thread Hierarchy. CUDA编程模型使用层次结构组织threads。Threads被分组到blocks中,blocks被分组到grid中。这种层次结构允许可扩展的并行性。

Memory Hierarchy. CUDA提供不同类型的内存,每种内存都有其特性和使用场景,包括global memory、shared memory、constant memory和registers。

Automatic Distribution of Data and Program:

- NVIDIA编译器 (nvcc) 在编译时将数据和程序指令分配到不同的thread blocks中。每个thread block包含多个threads,这些threads进一步分为多个warps进行并行执行。
- 例如:上述kernel代码由nvcc编译,nvcc确定threads和blocks的分配。

Warp执行与调度

Warp Formation:

- block内的threads被分组为warps,每个warp包含32个threads。warp是GPU上的基本执行单元。 硬件管理warp的调度和执行,以确保高效的并行处理。
- 例如:如果一个block有64个threads,它将被分为2个warps,每个warp有32个threads。

Warp Scheduling:

- warp调度器根据warp的准备执行情况动态调度warp。这包括监控每个warp的状态(例如,等待内存访问,准备执行)并在warp之间切换以隐藏延迟并保持GPU核心忙碌。
- 例如: 当warp 0等待内存加载时,调度器可以切换到warp 1执行下一条指令,从而最大化吞吐量。

Latency Hiding and Context Switching:

调度器通过切换到其他准备执行的warp来隐藏内存延迟。这类似于CPU使用上下文切换来处理中断和保持响应性。

自有DRAM (Device Memory)

- Dedicated Device Memory
 - CUDA设备具有自己的专用DRAM,称为device memory,与主机系统的内存(RAM)分离。
 这种device memory对于GPU存储大型数据集、纹理和处理中所需的中间结果至关重要。内存空间的隔离确保了GPU可以快速访问其内存,而不会受到CPU内存带宽的瓶颈影响。
- Memory Hierarchy and Types. CUDA设备中的内存层次结构包括几种内存类型,每种内存都有不同的特性:
 - 。 **Global Memory**: 这是GPU的主内存空间,所有threads都可以访问。它容量大但延迟相对较高。
 - o **Shared Memory**: Shared memory是thread block内的线程共享的小型低延迟内存空间。它比global memory快得多,用于线程间通信和缓存频繁访问的数据。
 - Constant Memory: 这是一种只读内存空间,具有缓存优化,用于将相同的值广播给多个threads。
 - **Texture Memory**: 这种内存优化用于典型图形应用中的只读访问模式,并提供硬件加速的过滤和寻址模式。
 - o Registers: Registers是最快的内存类型,私有给每个线程,用于存储临时变量。
- Memory Management and Data Transfer
 - o 为了利用GPU的内存,CPU必须管理主机内存和设备内存之间的数据传输。可以使用CUDA API函数来实现,如 cudaMalloc() 分配device memory, cudaMemcpy() 在主机和设备之间 复制数据,以及 cudaFree() 释放device memory。
 - 高效的数据传输和内存管理对于实现CUDA应用的高性能至关重要。开发者必须最小化数据传输并仔细管理内存分配以避免瓶颈。使用CUDA streams和异步内存操作的计算与数据传输重叠等技术有助于优化性能。
- Hardware Implementation
 - GPU的内存控制器管理对device memory的访问,处理来自SP的读写请求。内存控制器使用 诸如内存合并等技术来优化访问模式并减少延迟。此外,GPU架构支持高带宽内存接口,允许 在GPU和其内存之间快速数据传输。
 - Device memory与host memory的分离意味着数据必须在它们之间显式传输。GPU驱动程序和CUDA runtime管理这些数据传输,通过PCIe总线协调数据传输,确保数据完整性和带宽的高效使用。

线程的并行执行

- Overview of Parallel Execution
 - CUDA设备的核心优势之一是能够并行执行成干上万个threads。这种能力源于GPU的架构, 设计用于处理大规模并行性。每个thread执行较大计算的一小部分,从而实现大规模并行性。 这对于矩阵乘法、图像处理和科学模拟等任务特别有利,因为同样的操作需要独立应用于许多 数据元素。
- Thread Hierarchy and Management. Threads在CUDA中组织成网格和块的层次结构。Grid由多个blocks组成,每个block包含多个threads。这种层次结构允许灵活的组织和高效的并行任务执行。
 - 。 Grids: 层次结构的顶层,表示执行kernel所需的所有线程集合。
 - 。 **Blocks**: Grid的子部分,包含固定数量的threads。 Blocks独立调度,可以以任何顺序执行。
 - o Threads: Block内的单个执行单元。每个thread执行相同的kernel代码,但操作不同的数据。

- Warp and Scheduler: Warp是32个线程的组,在任何给定时间执行相同的指令。GPU的调度器管理这些warps,并动态调度它们以最大化利用率和隐藏内存延迟。
 - Warp Scheduling: GPU调度器可以在warp之间切换,以确保计算单元保持忙碌,即使某些线程在等待内存操作完成。
- Allocation Mechanisms
 - Static Allocation: 在kernel 启动时定义线程、blocks和grids的数量。这种静态分配简化了执行模型,但需要仔细计划以优化资源使用。
 - Dynamic Allocation: 某些CUDA功能允许在kernel执行期间动态分配资源,例如动态并行性,其中一个kernel可以启动其他kernels。
- Comparison with Sequential Programs
 - 在顺序程序中,单线程按顺序处理指令。这种方法简单但对于大型数据集或计算密集型任务可能很慢。
 - 相比之下,CUDA的并行执行模型将任务分解为许多小部分,允许同时执行。这种并行性显著加快了适合任务的处理时间。
- Comparison with Distributed Memory Programs
 - 使用MPI (Message Passing Interface) 的分布式内存程序涉及在不同节点上运行的多个独立 进程,每个进程都有自己的内存。这些进程通过消息传递进行通信。
 - CUDA的并行执行模型类似于将工作分配给许多线程,但不同之处在于其GPU内的共享内存架构。Block内的线程可以使用shared memory进行通信和同步,相比于分布式系统中的显式消息传递提供了更集成的方法。

Kernel在数据并行应用中的定义和作用

- **Definition of Kernels**: Kernel是用CUDA C/C++编写的函数,在GPU上执行。它们定义每个thread将执行的计算。Kernel启动时,由指定数量的threads并行执行。
- 在数据并行应用中,kernel应用于大数据集,每个thread处理一部分数据。并行执行模型通过将工作分配给许多threads来显著加速计算任务。
 - 例如:矩阵乘法。在矩阵乘法中,每个thread可以被分配计算输出矩阵的单个元素。Kernel将 执行其分配的元素的必要计算,访问输入矩阵的相应元素。
 - 例如: 图像处理。在图像处理过程中,每个thread可能处理单个像素或图像的小区域。这种并 行方法使实时图像增强和转换成为可能。
- Memory Allocation for Kernels
 - o **Global Memory**: Kernels通常访问global memory以进行输入和输出数据。每个thread根据 其唯一的thread ID计算其内存地址。
 - Shared Memory: 在一个block内, threads可以使用shared memory共享数据并协调其计算。这减少了对较慢的global memory的访问需求,提高了性能。
 - Registers and Local Memory: 每个thread都有自己的一组寄存器,用于存储临时变量。如果变量过多,则使用local memory (一种较慢的外部芯片内存)。
- Kernel Launch Configuration: Kernel 启动配置指定了grid和block的维度:
 - 1 KernelFunction<<<numBlocks, numThreadsPerBlock>>>(arguments);
 - o numBlocks: 指定grid包含的blocks数量。
 - o numThreadsPerBlock: 指定每个block内的threads数量。

- **Synchronization within Kernels**: CUDA提供同步机制,如__syncthreads(),确保block内的 所有threads在代码的某一点到达后再继续。这对于协调shared memory访问和避免竞态条件至关 重要。
- Comparison with Sequential Programs
 - 顺序程序一次执行一条指令。单线程的矩阵乘法将按顺序计算输出矩阵的每个元素。
 - o 并行CUDA kernel的矩阵乘法同时计算多个元素,显著减少计算时间。
- Comparison with Distributed Memory Programs
 - o 在分布式内存程序中,每个进程独立处理数据集的一部分,进程之间通过消息传递进行通信。
 - CUDA kernels使用shared memory进行快速的block内通信,并使用global memory进行较大数据交换,提供了更紧密的并行计算内存模型。

Array of Streaming Processors (SPs)

Array of Streaming Processors (SPs)

每个CUDA设备中的Streaming Multiprocessor (SM)包含多个Streaming Processors (SPs)。这些SPs是执行kernel指令的基本核心。一个SM的架构通常包含几十个SPs,使其能够同时处理许多线程。

- Arithmetic Logic Units (ALUs). 每个SP包含ALUs,负责执行整数算术和逻辑操作。这些单元处理加法、减法、位运算和比较等任务。ALUs对于执行控制流操作和整数密集型计算至关重要。
- Floating-Point Units (FPUs). 除了ALUs, SPs还配备了FPUs, 专门用于执行浮点算术运算。FPUs 处理加法、减法、乘法、除法和平方根等操作。这些单元对于科学计算、模拟和其他需要高精度的任务至关重要。
- Instruction Set
 - SPs执行丰富的指令集,包括整数和浮点运算。指令集设计用于支持通用计算、图形渲染和科学计算所需的广泛操作。
 - Common Instructions: 包括算术运算(add, sub, mul, div)、逻辑运算(and, or, xor)、数据移动指令(load, store)和控制流指令(branch, call)。
 - Specialized Instructions: 包括超越函数 (sin, cos, exp, log) 和面向并行的指令如同步和原子操作。

Pipeline Stages. SPs采用流水线架构,使其能够同时处理多条指令。流水线分为几个阶段,每个阶段处理指令处理的特定部分:

- 1. Fetch: 从指令缓存中取指令。
- 2. Decode: 解释指令并确定所需的操作和资源。
- 3. Execute: 使用ALUs或FPUs执行算术或逻辑操作。
- 4. Memory Access: 如果指令涉及数据移动,则读取或写入内存。
- 5. Write Back: 将计算结果写回寄存器文件或内存。

SPs的流水线架构通过允许多条指令在不同的执行阶段并行处理,实现高指令吞吐量。这种并行性确保计算单元高效利用,最大化GPU性能。

Warp Execution

在CUDA中,一个warp由32个线程组成,这些线程同步执行指令。这种SIMD(单指令多数据)模型利用数据并行性实现高性能。warp中的所有线程同时执行相同的指令,但操作不同的数据元素。

Example: Vector Addition. 这里,每个线程计算结果向量 C 的一个元素,该元素由 A 和 B 的对应元素相加得到。

```
1 __global__ void vectorAdd(float *A, float *B, float *C, int N) {
2    int idx = threadIdx.x + blockDim.x * blockIdx.x;
3    if (idx < N) {
4        C[idx] = A[idx] + B[idx];
5    }
6 }</pre>
```

Branch Divergence. 类似于CPU中的分支预测,GPU也会遭遇分支发散问题。如果warp中的不同线程由于条件语句而采取不同的执行路径,warp会顺序执行每个路径,从而降低效率。如果warp中的线程遵循if语句的不同分支,warp将顺序执行两个分支,从而降低并行效率。

```
__global__ void vectorAddWithCondition(float *A, float *B, float *C, int N) {
1
 2
        int idx = threadIdx.x + blockDim.x * blockIdx.x;
 3
        if (idx < N) {
            if (A[idx] > 0) {
 4
 5
                C[idx] = A[idx] + B[idx];
 6
            } else {
 7
                C[idx] = B[idx];
8
            }
9
        }
10 }
```

Warp Scheduling

GPU的warp调度器根据warp的准备情况动态选择下一个执行的warp。这种动态调度有助于隐藏由于内存访问或其他停滞造成的延迟,类似于现代CPU如何处理指令级并行和乱序执行。

Example: 假设warp 0正在等待从全局内存加载数据。warp调度器可以切换到准备执行下一条指令的warp 1,从而保持SM忙碌。这里,当一个warp等待全局内存加载时,另一个warp可以执行以保持流水线忙碌,类似于CPU使用超线程技术保持执行单元忙碌。

```
__global___ void processVectors(float *A, float *B, float *C, int N) {
    int idx = threadIdx.x + blockDim.x * blockIdx.x;
    if (idx < N) {
        float a = A[idx]; // 从全局内存加载
        float b = B[idx]; // 从全局内存加载
        C[idx] = a + b; // 执行计算
    }
}
```

Synchronization

warp内的线程可以使用特殊指令如__syncthreads()进行同步,确保所有线程在继续之前到达某一点。这类似于CPU并行编程中的屏障。

```
cudaCopy code__global__ void matrixTranspose(float *in, float *out, int
 1
    width) {
 2
        __shared__ float tile[32][32];
 3
        int x = blockIdx.x * 32 + threadIdx.x;
        int y = blockIdx.y * 32 + threadIdx.y;
 4
 5
        tile[threadIdx.y][threadIdx.x] = in[y * width + x];
 6
        __syncthreads();
 7
        x = blockIdx.y * 32 + threadIdx.x;
        y = blockIdx.x * 32 + threadIdx.y;
8
9
        out[y * width + x] = tile[threadIdx.x][threadIdx.y];
10
   }
```

__syncthreads()确保block中的所有线程在读取数据之前已经将数据写入shared memory,这类似于多线程CPU程序中的屏障同步。

Inter-Warp Communication

不同warp之间的通信通常依赖于全局内存。每个warp将数据写入全局内存,其他warp可以读取这些数据。这类似于操作系统中的进程间通信(IPC)机制,进程通过共享内存或消息传递进行通信。这里,每个block计算部分和并将其写入全局内存。最终的和使用原子操作累加,防止竞态条件,类似于多线程CPU程序中使用互斥锁或原子操作。

```
1
    __global__ void accumulateResults(float *results, float *finalSum, int N) {
 2
        int idx = threadIdx.x + blockDim.x * blockIdx.x;
 3
        __shared__ float sharedSum[32];
 4
        sharedSum[threadIdx.x] = (idx < N) ? results[idx] : 0;</pre>
 5
         __syncthreads();
 6
        if (threadIdx.x == 0) {
 7
            float blockSum = 0;
            for (int i = 0; i < 32; ++i) {
 8
 9
                 blockSum += sharedSum[i];
10
11
            atomicAdd(finalSum, blockSum);
        }
12
13
    }
```

Cooperative Groups and Warp-Level Primitives. 现代GPU架构通过协作组和warp级原语提供更高效的warp间通信机制。这类似于CPU并行编程中的高级同步技术,如线程池或工作窃取队列。在这个例子中,__shf1_down_sync 是一个warp级原语,允许warp内的线程直接通信,减少了对全局内存访问的需求,提高了效率。这类似于在CPU上使用SIMD指令如SSE或AVX进行高效的线程内通信。

```
1  __global__ void reduceSum(float *data, float *result, int N) {
2   unsigned int tid = threadIdx.x;
3   float sum = 0.0f;
4   for (int i = tid; i < N; i += blockDim.x) {
5     sum += data[i];
6   }
7   sum = warpReduceSum(sum);</pre>
```

```
8 if (tid % warpSize == 0) {
 9
            atomicAdd(result, sum);
10
        }
11
   }
12
13
    __inline__ __device__ float warpReduceSum(float val) {
         for (int offset = warpSize / 2; offset > 0; offset /= 2) {
14
            val += __shfl_down_sync(FULL_MASK, val, offset);
15
16
        }
        return val;
17
18
```

Special Function Units (SFUs)

Special Function Units (SFUs). 除了SPs,每个SM还包括Special Function Units (SFUs),专门处理复杂的数学运算。这些单元加速了诸如正弦、余弦、平方根和插值等函数的计算,这些函数在图形和科学应用中经常使用。

- SFUs有专用的硬件来计算超越函数和插值函数。这种专用硬件使得这些操作比使用通用的ALUs或 FPUs实现要快得多。
- 通过将复杂的数学运算分配给SFUs, SPs可以专注于更通用的计算, 从而提高GPU的整体效率和性能。

Shared Memory: 16KB Read/Write Memory Managed by Software

每个SM提供一个小型的、低延迟的共享内存空间,通常为16KB,所有线程都可以在一个块内访问。与硬件管理的缓存不同,共享内存是软件管理的数据存储,允许程序员显式地控制数据共享和同步。

- **Shared Memory的结构**: 共享内存被组织成多个bank,每个bank允许多个线程同时访问,只要没有bank冲突(即多个线程访问同一bank内的不同位置)。
- **延迟和带宽**: 共享内存比全局内存快得多,其延迟可与寄存器访问相媲美。这使得它非常适合存储需要被块内多个线程频繁访问的数据。

高效使用共享内存可以显著提升并行应用的性能,因为减少了对较慢的全局内存的访问需求。共享内存常用于缓存数据、执行归约操作和促进线程间通信。在矩阵乘法中,共享内存可以用来加载子矩阵(tiles),使得多个线程可以同时工作,减少对全局内存的访问次数。

Multithreading Issuing Unit: Dispatches Instructions

在一个Streaming Multiprocessor (SM)内的多线程发出单元是一个关键组件,负责管理多个warps (32 个线程的组)同时执行。其主要目标是确保GPU核心保持忙碌,并有效地隐藏延迟。该单元在GPU的计算模型中起着至关重要的作用,通过优化执行流程和最大化资源利用率来提高性能。

Dispatching Instructions to SPs. 多线程发出单元负责将指令分配给Streaming Processors (SPs)。它决定哪些指令要执行,并管理执行顺序。这样确保了SPs始终有指令可执行,减少了空闲时间,增加了吞叶量。

- Instruction Fetch Unit: 从指令缓存中取指令。
- Decode Unit: 解码取出的指令,以确定所需的操作和资源。
- Scheduler: 确定执行顺序,并将解码后的指令分配给相应的SPs。
- SPs: 执行指令并进行计算。

Warp Scheduling

Scheduling的分类:

Dynamic Scheduling. 发出单元根据warps的准备情况动态调度。这涉及检查一个warp是否准备好执行下一条指令或是否在等待内存数据。动态调度通过切换到另一个准备好执行的warp,帮助隐藏由于内存访问或其他停滞造成的延迟。

- Warp Scheduler: 监控每个warp的状态,并决定下一个调度的warp。
- Execution Control Unit: 管理每个warp的状态,并根据需要在warps之间切换。
- Memory Access Unit: 处理内存操作,并将状态反馈给warp调度器。

Instruction Dispatch. 指令以轮询或优先级方式分配给SPs。这确保了所有warps都有公平的执行机会,并且高优先级的warps可以在需要时更快地执行。这种机制有助于平衡所有SPs的负载,防止任何单一warp垄断计算资源。

- Round-Robin Dispatcher: 通过轮询确保公平调度。
- Priority Dispatcher: 根据预定的标准优先调度warps,确保高优先级任务及时执行。

CUDA设备中的warp调度器负责决定在任何给定时间执行哪个warp。其主要目标是通过动态调度warps,基于它们的准备情况来最大化Streaming Processors (SPs)的利用率。

1. Monitoring Warp Status

- 。 调度器持续监控每个warp的状态。warps可以处于不同的状态,如准备执行、等待内存访问或 因其他依赖关系(如同步点)而停滞。
- o Ready to Execute: warp拥有执行下一条指令所需的所有数据和资源。
- o Waiting for Memory: warp正在等待从全局内存或其他内存层次获取数据。
- o Stalled: warp在等待其他依赖关系,例如块内线程的同步。

2. Dynamic Scheduling

- 。 调度器根据warp的状态动态选择下一个执行的warp。这类似于操作系统管理进程调度,但侧重于最大化并行执行和隐藏延迟。
- Round-Robin Scheduling: 调度器循环选择可用的warps,确保每个warp都有公平的执行机会。这确保没有任何单一warp垄断资源。
- o Priority Scheduling: 更关键或接近完成任务的warps可能会被赋予更高的优先级。这在某些任务需要比其他任务更快完成的情况下很有帮助。
- Latency-Hiding: 当一个warp在等待内存访问时,调度器切换到另一个准备好执行的warp。
 这通过保持SPs忙碌来帮助隐藏内存延迟。

3. Execution Control

- 一旦选择了warp,调度器就会将其下一条指令分配给SPs。执行控制单元管理warp的状态,包括程序计数器和寄存器值。
- o **Context Switching**: 调度器维护每个warp的上下文,类似于操作系统在中断期间处理上下文切换。这包括在warp调度进出执行时保存和恢复warp的状态。

Instruction and Constant Caches

Instruction Cache

Streaming Multiprocessors (SMs)中的指令缓存是一种小型高速内存,用于存储正在运行的kernels的指令。虽然这种缓存对程序员不可见,但它在kernel执行效率中起着至关重要的作用。

结构和大小:

- 指令缓存通常较小,设计用于存储当前执行的kernels中最常访问的指令。其大小经过优化,以平衡快速访问的需求与芯片空间和功耗的限制。
- **例子**: 缓存大小可能从几千字节到几十千字节不等,足以存储多个kernel指令和反复执行的循环。

性能作用:

- 通过将常用指令存储在靠近SPs的位置,指令缓存将从全局内存获取指令的延迟最小化。这很重要,因为访问全局内存比访问缓存要慢得多。
- **性能影响**: 减少指令获取延迟直接提高了kernel的总体执行速度。例如,在一个紧凑的循环中,反复执行相同的指令,缓存中的这些指令避免了重复从全局内存获取的开销。

硬件实现:

- **缓存层次结构**: 指令缓存是GPU多级缓存层次结构的一部分。它位于全局内存和执行单元之间,为指令获取提供高速缓冲。
- **关联性和替换策略**: 缓存可能使用组相联映射,并采用LRU(最近最少使用)替换策略来管理哪些指令保留在缓存中。这确保了最常访问的指令留在缓存中,而不常使用的指令被移出。
- **与获取和解码阶段的集成**: 指令缓存与指令流水线的获取和解码阶段紧密集成。指令从缓存中取到解码单元,减少了开始执行每条指令所需的时间。

Constant Cache

SMs中的常量缓存存储在kernel执行期间线程频繁访问的只读数据。该缓存优化了数据广播到多个线程的效率。

• 结构和访问:

- 常量缓存设计用于处理在kernel执行期间保持不变的只读数据,如数学公式中的系数或配置参数。它的大小通常为这些小而频繁访问的数据集进行优化。
- 广播机制:如果warp中的多个线程读取常量内存中的相同位置,数据只从常量缓存中获取一次,并广播给所有请求的线程。这减少了内存带宽的使用和延迟,因为多个线程可以共享一次获取操作。
- **例子**: 如果warp中的32个线程需要相同的常量值,该值只获取一次并同时分发给所有线程,而不是每个线程单独获取该值。

常量内存非常适合存储频繁读取但不修改的数据,如查找表、计算中使用的常数和配置参数。常量缓存确保这些值对SPs快速可用。考虑一个使用查找表系数的kernel:

```
1
   __constant__ float coeffs[10]; // 存储在常量内存中
2
3
    __global__ void computeKernel(float *data, int N) {
4
       int idx = threadIdx.x + blockIdx.x * blockDim.x;
5
       if (idx < N) {
            float value = data[idx];
6
7
            float result = value * coeffs[0] + value * coeffs[1]; // 访问常量内存
8
           data[idx] = result;
9
       }
10
   }
```

硬件实现:

- **缓存设计**: 常量缓存设计得很小但对于其特定用例非常高效。它通常是全相联的,意味着任何内存地址都可以映射到任何缓存行,这有助于最小化缓存未命中。
- **访问优化**: 常量缓存针对warp中的多个线程读取相同数据的场景进行了优化。这允许一次获取操作为所有线程服务,大大减少了对内存带宽的需求。
- **广播引擎**: 硬件包括一个广播引擎,确保从常量缓存获取的数据高效地分发给warp中所有请求的线程。

CUDA编程结构

集成的host-device应用程序

一个CUDA程序在单个应用程序中集成了主机(CPU)代码和设备(GPU)代码。主机代码管理整体控制流、内存分配和数据传输,而设备代码(kernels)在GPU上执行计算密集型任务。

host 代码:

- 使用标准的C/C++编写,并带有CUDA扩展。
- 在CPU上运行,管理GPU上的内存分配和释放。
- 在主机和设备内存之间传输数据。
- 在GPU上启动kernel函数。
- 同步执行并处理任何错误。

```
1 float *d_A;
   cudaMalloc((void**)&d_A, size);
4
   // 从主机复制数据到设备
   cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
 5
7
   // 启动kernel
   vectorAdd<<<<numBlocks, numThreadsPerBlock>>>(d_A, d_B, d_C, N);
8
9
10
   // 将结果复制回主机
11
   cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
12
13
   // 释放设备内存
14 cudaFree(d_A);
```

device 代码:

- 使用CUDA C/C++编写。
- 在GPU上执行。
- 包含设计为在多个线程上并行运行的kernel函数。

```
1  __global__ void vectorAdd(float *A, float *B, float *C, int N) {
2   int idx = threadIdx.x + blockDim.x * blockIdx.x;
3   if (idx < N) {
4        C[idx] = A[idx] + B[idx];
5   }
6 }</pre>
```

Kernel调用语法和示例结构

从主机代码中启动kernel函数使用一种特殊的语法,指定网格和块的维度。这种语法称为执行配置。

语法: kernelFunction<<<numBlocks, numThreadsPerBlock>>>(arg1, arg2, ...)

- numBlocks: 网格中的块数量。
- numThreadsPerBlock:每个块中的线程数量。

```
__global__ void vectorAdd(float *A, float *B, float *C, int N) {
1
        int idx = threadIdx.x + blockDim.x * blockIdx.x;
2
3
        if (idx < N) {
            C[idx] = A[idx] + B[idx];
4
5
        }
6
   }
   // 启动kernel的主机代码
8
   int N = 1024;
9
10 | int numThreadsPerBlock = 256;
11
   int numBlocks = (N + numThreadsPerBlock - 1) / numThreadsPerBlock;
12 vectorAdd<<<numBlocks, numThreadsPerBlock>>>(d_A, d_B, d_C, N);
```

并行线程数组

CUDA kernels使用并行线程数组执行,遵循单程序多数据(SPMD)模型。每个线程运行相同的代码,但操作不同的数据元素。一个块中的每个线程可以使用 threadIdx 唯一识别,一个网格中的每个块可以使用 blockIdx 唯一识别,块和网格的维度可以使用 blockDim和 gridDim访问。线程使用其唯一的索引来计算内存地址并进行控制决策,确保每个线程处理数据的不同部分。

```
1  __global__ void vectorAdd(float *A, float *B, float *C, int N) {
2    int idx = threadIdx.x + blockDim.x * blockIdx.x;
3    if (idx < N) {
4        C[idx] = A[idx] + B[idx];
5    }
6 }</pre>
```

块和网格组织

块和网格是CUDA中管理并行执行的基本组织单位。

- 线程ID: threadIdx.x, threadIdx.y, threadIdx.z
- 块ID: blockIdx.x, blockIdx.y, blockIdx.z
- 块维度: blockDim.x, blockDim.y, blockDim.z
- 网络维度: gridDim.x, gridDim.y, gridDim.z

块内的线程和网格内的块可以组织成1D、2D或3D配置。这种灵活性简化了多维数据结构(如矩阵和体积)的寻址。

```
__global__ void matrixAdd(float *A, float *B, float *C, int width, int
2
        int x = blockIdx.x * blockDim.x + threadIdx.x;
        int y = blockIdx.y * blockDim.y + threadIdx.y;
3
4
        int idx = y * width + x;
5
        if (x < width && y < height) {
            C[idx] = A[idx] + B[idx];
6
7
        }
8
   }
9
    // 启动kernel的主机代码
10
   dim3 numThreadsPerBlock(16, 16);
11
12
    dim3 numBlocks((width + numThreadsPerBlock.x - 1) / numThreadsPerBlock.x,
                   (height + numThreadsPerBlock.y - 1) / numThreadsPerBlock.y);
13
14
    matrixAdd<<<numBlocks, numThreadsPerBlock>>>(d_A, d_B, d_C, width, height);
```

线程层次结构和同步

层次化组织

CUDA的层次化组织设计用于高效管理并行计算,通过将任务划分为更小、更易管理的单元。这一层次结构包括线程、线程块、协同线程数组(CTAs)和网格。

Threads: CUDA中最小的执行单元。每个线程独立执行kernel函数,但可以与同一块内的其他线程协作。线程通过 threadIdx 标识,可以是1D、2D或3D。

```
int tx = threadIdx.x;
int ty = threadIdx.y;
int tz = threadIdx.z;
```

Thread Blocks:可以通过共享内存和同步执行来协作的一组线程。每个块通过 blockIdx 标识,并由 blockDim 指定其维度。块内的线程可以通过 __syncthreads() 进行通信和同步。

```
int bx = blockIdx.x;
int by = blockIdx.y;
int bz = blockIdx.z;
int blockSize = blockDim.x * blockDim.y * blockDim.z;
```

Cooperative Thread Arrays (CTAs): 线程块的另一种说法,强调其在整体任务中的协作能力。每个CTA独立执行,可以相对于其他CTAs以任何顺序调度。CTAs对于将大问题分解为可以并行执行的小块至关重要。

Grids:执行给定kernel的一组线程块。整个网格的块同时启动,每个块在GPU上运行。网格可以是1D、2D或3D,允许灵活组织并行任务。

```
1 dim3 gridDim(16, 16, 1); // 16x16的块网格
2 dim3 blockDim(16, 16, 1); // 每个块有16x16的线程
3 kernelFunction<<<gridDim, blockDim>>>(...);
```

块内的同步和共享内存

块内的同步和共享内存对于线程间高效协作至关重要。

Synchronization:块内的线程可以使用___syncthreads() 同步其执行。这个内建函数确保块内的所有 线程在继续执行之前到达同步点。这对于协调对共享资源的访问和避免竞争条件非常有用。

```
__global__ void kernel(float *data) {
1
 2
        __shared__ float sharedData[256];
 3
        int idx = threadIdx.x;
        sharedData[idx] = data[idx];
 4
 5
        __syncthreads();
 6
        if (idx < 128) {
 7
            sharedData[idx] += sharedData[idx + 128];
8
        }
 9
        __syncthreads();
10
        data[idx] = sharedData[idx];
    }
11
```

Shared Memory:共享内存是块内所有线程可访问的低延迟内存。用于存储需要快速访问和线程间共享的数据。使用__shared__关键字在kernel内声明共享内存。适当使用共享内存可以显著提高CUDA应用的性能,因为它减少了对较慢的全局内存访问的需求。

```
__global__ void matrixMul(float *A, float *B, float *C, int N) {
1
 2
        __shared__ float sharedA[16][16];
 3
         __shared__ float sharedB[16][16];
 4
        int tx = threadIdx.x;
 5
        int ty = threadIdx.y;
 6
        int row = blockIdx.y * blockDim.y + ty;
 7
        int col = blockIdx.x * blockDim.x + tx;
        float value = 0;
 8
        for (int k = 0; k < (N / 16); ++k) {
 9
10
            sharedA[ty][tx] = A[row * N + (k * 16 + tx)];
            sharedB[ty][tx] = B[(k * 16 + ty) * N + col];
11
12
            __syncthreads();
13
            for (int n = 0; n < 16; ++n) {
14
                value += sharedA[ty][n] * sharedB[n][tx];
15
            }
16
            __syncthreads();
17
18
        C[row * N + col] = value;
19
    }
```

运行时网格的动态调度

网格的动态调度使CUDA程序的执行更灵活和高效,特别适用于负载不均匀或执行过程中动态变化的工作负载。

Kernel 启动. GPU可以根据运行时条件动态启动新的kernel。这允许更复杂的执行流和自适应负载平衡。

```
__global__ void parentKernel() {
2
       // 一些计算
3
       if (condition) {
4
           // 动态启动子kernel
5
            childKernel<<<gridDim, blockDim>>>(...);
6
       }
7
    }
8
9
   __global__ void childKernel() {
10
       // 子kernel计算
11 }
```

网格和块维度. 网格和块的维度可以在运行时动态设置,以适应应用的具体需求。这允许更有效地利用 GPU资源。

```
int N = ...; // 输入的大小
int numThreadsPerBlock = 256;
int numBlocks = (N + numThreadsPerBlock - 1) / numThreadsPerBlock;
kernel<<<numBlocks, numThreadsPerBlock>>>(...);
```

Stream 和 Event 同步. CUDA流和事件可以用于动态同步kernel执行,允许更细粒度地控制执行顺序和kernels之间的依赖关系。

```
1 cudaStream_t stream1, stream2;
 2
    cudaStreamCreate(&stream1);
 3
    cudaStreamCreate(&stream2);
 4
 5
    kernelA<<<gridDim, blockDim, 0, stream1>>>(...);
    kernelB<<<qridDim, blockDim, 0, stream2>>>(...);
 6
 7
8
    cudaStreamSynchronize(stream1);
9
    cudaStreamSynchronize(stream2);
10
11 cudaStreamDestroy(stream1);
12
    cudaStreamDestroy(stream2);
```

CUDA 内存模型和变量管理

理解CUDA内存模型和变量管理对于优化CUDA应用程序的性能至关重要。CUDA的内存模型旨在提供灵活且高效的方式来处理各种类型的内存,每种内存都有其特定的作用范围、生命周期和用途。

变量类型限定符

CUDA提供了几种类型限定符来指定变量的内存范围和生命周期。这些限定符包括device、shared、constant和local。

device:

- Scope: 全局内存 (所有块中的所有线程都可以访问)。
- Lifetime: 应用程序的生命周期。
- Usage: 用于需要从多个kernel或主机代码访问的变量。

```
1 __device__ float deviceVar; // 所有线程都可以访问的全局变量
```

shared:

- Scope: 共享内存 (同一块内的所有线程都可以访问)。
- Lifetime: 块的生命周期。
- Usage: 用于需要在同一块内线程间共享的变量,提供低延迟访问。

```
__global___ void kernel() {
__shared__ float sharedvar[256]; // 块内共享变量
int idx = threadIdx.x;
sharedvar[idx] = idx; // 每个线程访问其对应的元素
__syncthreads(); // 同步块内所有线程
}
```

constant:

- Scope: 常量内存 (只读内存, 所有线程都可以访问)。
- Lifetime: 应用程序的生命周期。
- Usage: 用于只读且在kernel执行期间保持不变的变量。优化了warp内线程的广播读取。

```
1 __constant__ float constVar[10]; // 所有线程都可以访问的常量变量
```

local:

- Scope: 局部内存(每个线程私有,用于自动变量)。
- Lifetime: 线程的生命周期。
- Usage: 用于每个线程私有的变量,如函数中的局部变量。

```
1 __global__ void kernel() {
2    int localVar = threadIdx.x; // 每个线程私有的局部变量
3 }
```

内置向量类型和dim3类型

CUDA提供内置的向量类型和 dim3 类型,以简化内存管理和线程/块的组织。

Vector Types:

• CUDA支持内置的向量类型,如 int2 、 int4 、 float2 、 float4 等。这些类型将多个标量值组合成一个变量,允许更高效的内存访问和操作。

• 构造函数: 提供了构造函数如 make_int4 和 make_float4 来创建这些向量类型。

```
1 int4 vec = make_int4(1, 2, 3, 4);
2 float4 fvec = make_float4(1.0f, 2.0f, 3.0f, 4.0f);
```

dim3 Type:

- dim3 类型用于指定CUDA kernel启动时线程块和网格的维度。它可以表示最多三个维度(x, y, z),适合组织复杂的数据结构,如2D和3D数组。
- 用法: dim3 通常用于定义每块的线程数和每个网格的块数。

```
1 dim3 numBlocks(16, 16, 1); // 16x16的块网格
2 dim3 blockDim(16, 16, 1); // 每个块有16x16的线程
3 kernelFunction<<<<numBlocks, blockDim>>>(...);
```

CUDA函数声明

CUDA中函数根据其执行位置和可调用上下文分为三类: device函数、global函数和host函数。

device函数:在GPU设备上执行,只能由另一个device函数或global函数调用,用于封装可在设备上复用的代码片段,提高代码的模块化和可维护性。

```
1 __device__ float deviceFunction(float a, float b) {
2    return a + b;
3 }
```

global函数:在GPU设备上执行,由主机代码调用,启动时指定线程块和网格的维度,用于定义需要在GPU上并行执行的核函数。

```
1  __global__ void kernelFunction(float *a, float *b, float *c, int N) {
2    int idx = threadIdx.x + blockDim.x * blockIdx.x;
3    if (idx < N) {
4        c[idx] = a[idx] + b[idx];
5    }
6  }</pre>
```

host函数:在CPU主机上执行,可以由其他主机函数调用,不能在设备上执行。用于定义仅在主机上执行的函数,与普通C/C++函数类似。

对设备执行函数的限制

在CUDA编程中,设备执行的函数(即device和global函数)有一些特定的限制:

- 不能有递归调用。
- 不能使用可变长度数组。
- 不能有静态变量(每个线程都有自己的变量副本)。
- 不能使用标准库中的某些函数 (如malloc和free) 。

Kernel执行配置:在调用**global**函数时,必须指定网格和块的维度。使用 dim3 类型可以方便地定义 1D、2D或3D的块和网格。

```
1 dim3 gridDim(16, 16); // 16x16的网格
2 dim3 blockDim(16, 16); // 每块有16x16的线程
3 kernelFunction<<<gridDim, blockDim>>>(d_a, d_b, d_c, N);
```

内存管理和数据传输

使用 cudaMalloc() 分配设备内存,使用 cudaFree()释放设备内存。

```
1 float *d_ptr;
2 size_t size = 1024 * sizeof(float);
3 cudaMalloc(&d_ptr, size); // 分配内存
4 cudaFree(d_ptr); // 释放内存
```

使用 cudaMemcpy() 在主机和设备之间复制数据,使用 cudaMemcpyAsync() 进行异步数据传输,允许重叠计算和数据传输以提高效率。

```
1 float *h_data = (float*)malloc(size);
   float *d_data;
2
   cudaMalloc(&d_data, size);
3
4
   // 从主机复制到设备
6
   cudaMemcpy(d_data, h_data, size, cudaMemcpyHostToDevice);
 7
8
   // 从设备复制回主机
9
   cudaMemcpy(h_data, d_data, size, cudaMemcpyDeviceToHost);
10
11
   // 异步数据传输
12
   cudaMemcpyAsync(d_data, h_data, size, cudaMemcpyHostToDevice);
13
14
   free(h_data);
15 cudaFree(d_data);
```

线程同步和死锁

__syncthreads()用于块内线程的同步。确保块内所有线程都到达同步点后才继续执行。这种同步机制对于协调对共享资源的访问和避免竞争条件非常重要。

```
__global__ void exampleKernel(float *data) {
1
2
        __shared__ float sharedData[256];
 3
        int idx = threadIdx.x;
 4
       sharedData[idx] = data[idx];
 5
        __syncthreads();
 6
        if (idx < 128) {
 7
            sharedData[idx] += sharedData[idx + 128];
8
9
        __syncthreads();
10
        data[idx] = sharedData[idx];
11 }
```

如果一个线程在 __syncthreads() 之前退出,可能会导致其他线程永远等待,从而引发死锁。

```
1 __global___ void faultyKernel(float *data) {
2    int idx = threadIdx.x;
3    if (idx < 128) {
4        return; // 如果线程提前退出,可能会导致死锁
5    }
6    __syncthreads();
7    data[idx] = data[idx] * 2;
8  }
```

示例应用程序: 矩阵乘法

在CUDA中实现矩阵乘法是一个典型的并行计算示例,通过这个示例可以清晰地理解内存管理和线程管理的具体应用。

1. **定义矩阵维度和分配内存**:设定矩阵的大小(M×N和N×P),并分配主机和设备内存。

```
#define M 1024
#define N 1024

#define P 1024

float *h_A = (float*)malloc(M * N * sizeof(float));
float *h_B = (float*)malloc(N * P * sizeof(float));
float *h_C = (float*)malloc(M * P * sizeof(float));

float *d_A, *d_B, *d_C;
cudaMalloc(&d_A, M * N * sizeof(float));
cudaMalloc(&d_B, N * P * sizeof(float));
cudaMalloc(&d_C, M * P * sizeof(float));

cudaMalloc(&d_C, M * P * sizeof(float));
```

2. 初始化矩阵并复制到设备:在主机端初始化矩阵A和B,并将其复制到设备内存。

```
// 初始化矩阵A和B
for (int i = 0; i < M * N; ++i) h_A[i] = static_cast<float>(rand()) /
RAND_MAX;
for (int i = 0; i < N * P; ++i) h_B[i] = static_cast<float>(rand()) /
RAND_MAX;

cudaMemcpy(d_A, h_A, M * N * sizeof(float), cudaMemcpyHostToDevice);
cudaMemcpy(d_B, d_B, N * P * sizeof(float), cudaMemcpyHostToDevice);
```

3. **定义Kernel函数**:编写用于矩阵乘法的Kernel函数。

```
1
    __global__ void matrixMulKernel(float *A, float *B, float *C, int M, int
    N, int P) {
 2
        int row = blockIdx.y * blockDim.y + threadIdx.y;
 3
        int col = blockIdx.x * blockDim.x + threadIdx.x;
        float value = 0.0f;
 4
 5
        if (row < M && col < P) {
 6
 7
            for (int k = 0; k < N; ++k) {
 8
                value += A[row * N + k] * B[k * P + col];
9
10
            C[row * P + col] = value;
11
12 | }
```

4. **配置Kernel执行参数并启动Kernel**:设置线程块和网格的维度,启动Kernel进行矩阵乘法计算。

```
dim3 blockDim(16, 16); // 每个块包含16x16个线程
dim3 gridDim((P + blockDim.x - 1) / blockDim.x, (M + blockDim.y - 1) / blockDim.y);

matrixMulKernel<<<<gri>gridDim, blockDim>>>(d_A, d_B, d_C, M, N, P);
cudaDeviceSynchronize();
```

5. 复制结果到主机并释放内存:将计算结果从设备内存复制回主机内存,并释放分配的内存。

```
cudaMemcpy(h_C, d_C, M * P * sizeof(float), cudaMemcpyDeviceToHost);

// 释放内存
free(h_A);
free(h_B);
free(h_C);
cudaFree(d_A);
cudaFree(d_B);
cudaFree(d_C);
```

以下是用于矩阵乘法的Kernel函数示例。该函数计算矩阵 $C = A \times B$,其中A的尺寸为 $M \times N$,B的尺寸为 $N \times P$,C的尺寸为 $M \times P$ 。

```
__global__ void matrixMulKernel(float *A, float *B, float *C, int M, int N,
    int P) {
 2
        int row = blockIdx.y * blockDim.y + threadIdx.y;
 3
        int col = blockIdx.x * blockDim.x + threadIdx.x;
        float value = 0.0f;
 4
 5
 6
        if (row < M && col < P) {
 7
            for (int k = 0; k < N; ++k) {
8
                value += A[row * N + k] * B[k * P + col];
9
            C[row * P + col] = value;
10
        }
11
12 }
```

使用分块计算处理任意大小的矩阵

为了处理任意大小的矩阵,可以使用分块(tiled)计算方法。这种方法通过将矩阵分割成较小的子矩阵(tiles),使得每个线程块处理一个子矩阵,从而提高数据局部性和内存访问效率。

- 1. **定义分块大小**:选择合适的块大小(如16×16),定义共享内存来存储子矩阵。#define TILE_SIZE 16
- 2. **修改Kernel函数**: 修改Kernel函数以使用共享内存进行分块计算。

```
cppCopy code__global__ void matrixMulTiledKernel(float *A, float *B,
    float *C, int M, int N, int P) {
 2
        __shared__ float sharedA[TILE_SIZE][TILE_SIZE];
 3
        __shared__ float sharedB[TILE_SIZE][TILE_SIZE];
 5
        int tx = threadIdx.x;
 6
        int ty = threadIdx.y;
 7
        int row = blockIdx.y * TILE_SIZE + ty;
        int col = blockIdx.x * TILE_SIZE + tx;
 8
        float value = 0.0f;
 9
10
11
        for (int k = 0; k < (N + TILE_SIZE - 1) / TILE_SIZE; ++k) {
             if (row < M && k * TILE_SIZE + tx < N)
12
13
                 sharedA[ty][tx] = A[row * N + k * TILE_SIZE + tx];
14
            else
15
                 sharedA[ty][tx] = 0.0;
16
17
            if (col < P \&\& k * TILE_SIZE + ty < N)
                 sharedB[ty][tx] = B[(k * TILE_SIZE + ty) * P + col];
18
            else
19
20
                 sharedB[ty][tx] = 0.0;
21
22
            __syncthreads();
23
24
            for (int n = 0; n < TILE_SIZE; ++n)</pre>
25
                value += sharedA[ty][n] * sharedB[n][tx];
26
            __syncthreads();
27
28
        }
29
30
        if (row < M \&\& col < P)
31
            C[row * P + col] = value;
32 }
```

3. 配置并启动Kernel:设置适当的块和网格维度,启动Kernel进行分块计算。

```
cppCopy codedim3 blockDim(TILE_SIZE, TILE_SIZE);
dim3 gridDim((P + TILE_SIZE - 1) / TILE_SIZE, (M + TILE_SIZE - 1) /
TILE_SIZE);

matrixMulTiledKernel<<<gridDim, blockDim>>>(d_A, d_B, d_C, M, N, P);
cudaDeviceSynchronize();
```

通过分块计算方法,不仅可以处理任意大小的矩阵,还可以显著提高计算效率,减少对全局内存的访问,提高共享内存的利用率。

编译和调试工具

NVCC 编译器

NVCC(NVIDIA CUDA Compiler)是编译CUDA程序的主要工具。它将包含主机(CPU)和设备(GPU)代码的CUDA代码转换为可执行的二进制文件。以下是其编译过程和输出的详细说明:

编译过程:

- 1. **源代码解析**。NVCC首先解析CUDA源代码,该代码通常包含主机代码(用标准C/C++编写)和设备 代码(用CUDA C/C++编写)。设备代码包括将在GPU上执行的kernel函数。
- 2. **主机代码和设备代码的分离**。NVCC将主机代码和设备代码分开。主机代码使用标准的C/C++编译器(如GCC)编译,而设备代码由NVCC编译为PTX(Parallel Thread Execution)中间表示或直接编译为目标GPU架构的二进制代码。
- 3. **主机编译**。主机代码使用主机编译器编译为目标文件(.o或.obj)。这些目标文件包含管理和启动 CUDA kernels的CPU指令。

4. 设备编译:

- 。 设备代码编译为PTX代码,这是可以进一步优化的中间表示。NVCC也可以将PTX代码编译为GPU可直接执行的二进制代码(cubin文件)。
- o 设备代码也可以编译为fat binaries,包含用于未来GPU的PTX代码和当前GPU的二进制代码。
- 5. **链接**。将主机和设备目标文件链接在一起形成最终的可执行文件。主机代码包含管理内存分配、数据传输和在GPU上启动kernel的必要指令。

输出:

- 可执行二进制文件:可以在具有CUDA功能的GPU系统上运行的最终可执行文件。
- 目标文件: 主机和设备代码的中间目标文件。
- PTX文件:设备代码的中间表示,可用于进一步优化或编译为特定GPU架构的二进制代码。
- Cubin文件: GPU可直接执行的二进制代码。

使用NVCC编译CUDA程序的示例: nvcc -o vectorAdd vectorAdd.cu

该命令将 vectorAdd.cu CUDA源文件编译为名为 vectorAdd 的可执行文件。

调试工具

NVIDIA提供了一些强大的调试工具,帮助开发人员识别和修复CUDA程序中的问题。这些工具包括 NVIDIA cuda-gdb、cuda-memcheck、CUDA profiler和NVIDIA Nsight工具。

NVIDIA cuda-gdb

cuda-gdb是GNU Debugger (GDB)的扩展,支持调试CUDA应用程序中的主机(CPU)和设备(GPU) 代码。

- 在主机和设备代码中设置断点。
- 检查和修改主机和设备内存中的变量。
- 单步执行CPU和GPU代码。
- 查看主机和GPU代码的调用栈和回溯。

```
1cuda-gdb ./vectorAdd2(cuda-gdb) break vectorAdd.cu:20 # 在vectorAdd.cu第20行设置断点3(cuda-gdb) run # 启动程序4(cuda-gdb) bt # 显示调用栈5(cuda-gdb) step # 单步执行代码
```

cuda-memcheck

cuda-memcheck是一组用于检测CUDA应用程序中的内存错误的工具,类似于CPU程序的Valgrind。

• Memcheck: 检测越界内存访问、未对齐内存访问和内存泄漏。

• Racecheck: 检测共享内存中的数据竞争。

• Initcheck: 识别未初始化内存的使用。

• Synccheck: 报告同步原语使用中的错误。

```
1 cuda-memcheck ./vectorAdd
```

CUDA profiler

CUDA profiler提供CUDA应用程序的详细性能分析,帮助识别性能瓶颈并优化GPU资源的使用。

- 基于事件的分析:测量特定事件,如内存访问、指令计数和执行时间。
- **Visual Profiler (nvvp)**:一个GUI工具,提供性能指标、kernel执行时间线和内存传输时间线的可视化表示。

```
1 | nvprof ./vectorAdd
```

该命令分析 vectorAdd 可执行文件, 生成时间线和详细的性能指标。

Nsight Systems

Nsight Systems提供系统范围的性能分析工具,能够可视化应用的算法并识别优化机会。

- 可视化CPU和GPU活动。
- 将GPU活动与CPU函数相关联。
- 识别瓶颈并优化整体系统性能。

```
1 | nsys profile --trace=cuda ./vectorAdd
```

Nsight Compute

Nsight Compute是一个交互式kernel分析器,提供CUDA应用的详细性能指标和API调试。

- 收集单个kernel的详细性能指标。
- 提供用户友好的界面来分析性能问题。
- 支持自定义指标和引导分析。

```
1 | ncu ./vectorAdd
```

通过这些编译和调试工具,开发人员可以有效地开发、优化和调试CUDA应用程序,从而充分利用GPU的强大计算能力。