

异构计算与 MIC、CUDA 编译使用简介

张文帅 (wszhang@ustc.edu.cn)

中国科学技术大学超级计算中心

2016 年 10 月 14 日

① 异构众核计算

并行计算的现状与结构分类

异构计算的实现方式

② Intel 众核 (MIC) 计算

MIC 简介

MIC 计算运行方式

③ CUDA 计算与程序编译运行

CUDA 计算简介

CUDA 编译环境搭建

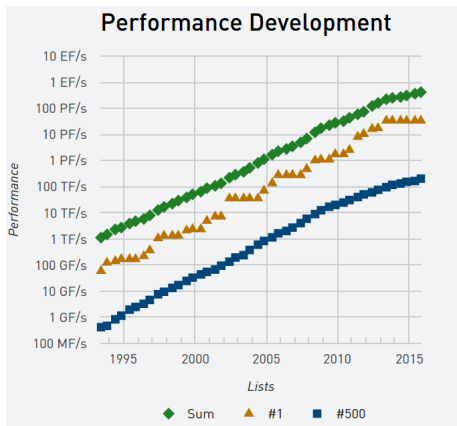
CUDA 程序编译

cuBLAS 函数库的使用

作业提交

并行计算现状

众核异构并行计算发展迅速，近期 Top 500 超级计算机中有 104 台使用众核异构加速卡
[66 Nvidia GPU (CUDA) + 29 Intel Xeon Phi (MIC) + 4 Nvidia & Xeon Phi + 3 ATI Radeon + 2 PEZY-SC(1024 cores)]

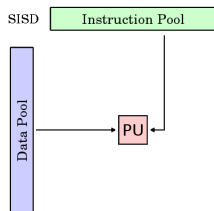


Top500 <http://top500.org>

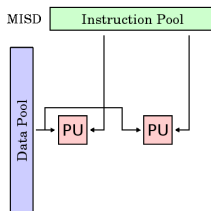
Tianhe-2 包含 16000 节点，每节点 2 Ivy Bridge chips and 3 Xeon Phi chips.

计算体系结构分类

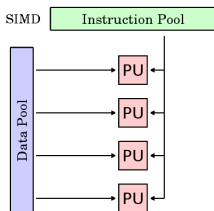
Single Instruction Single Data (SISD)



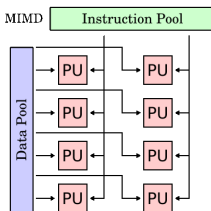
Multiple Instruction Single Data (MISD)



Single Instruction Multiple Data (SIMD)



Multiple Instruction Multiple Data (MIMD)



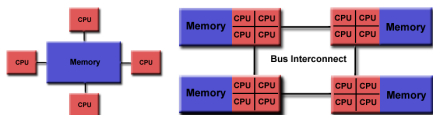
[https://computing.llnl.gov/tutorials/parallel_comp]

异构计算可以较好的实现 MIMD，在 CUDA 编程中，一个 warp(包含 32 线程) 为一个基本调度单元，只能执行同一个指令任务，属于 SIMD；同时多个不同 block 任务可以在不同的 SM 流多处理器中执行，构成 MIMD。

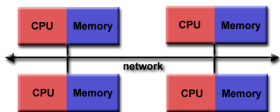
并行计算分类

主流的并行计算框架按照存储方式，可以分成

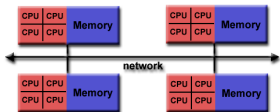
- 共享式存储，单一地址空间（物理或逻辑）



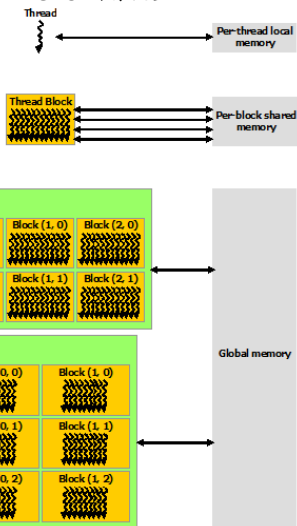
- 分布式存储



- 混合分布式共享存储（主流）



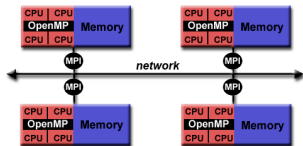
GPU 内存模型



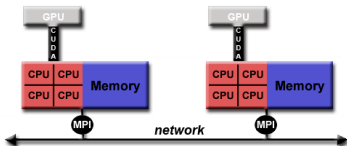
并行计算分类

按照计算机系统的硬件架构组成，区分为

- 通用架构并行
 - 同构多核并行 (CPU)

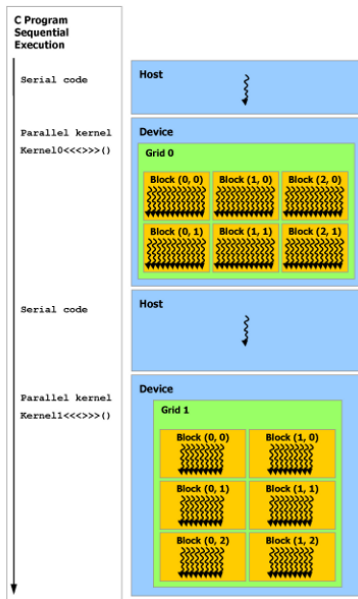


- 异构众核并行 (CPU+GPU / CPU+MIC)



- 专用架构并行
 - CPU+FPGA 异构

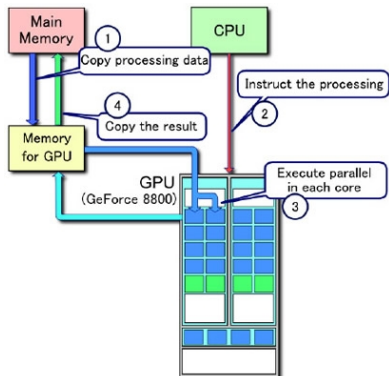
异构计算并行设计 (CUDA 为例)



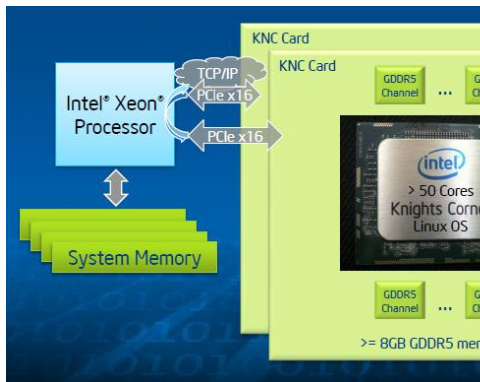
异构计算执行流程

通常异构计算中的加速卡部分执行流程分为四个部分，如 CUDA 程序的执行流程：

CUDA



MIC



对于简单使用编译指导语句的 MIC 计算，在 offload 模式下数据的拷贝操作被隐藏，但依然存在，故而存储带宽经常成为制约速度的主要因素，未来与主机高速共享内存是正在发展的解决方案之一。此外，MIC 程序还具有 native 执行模式，可以登录到卡上 linux 系统执行程序，具有更好的易用性和通用性。

① 异构众核计算

并行计算的现状与结构分类

异构计算的实现方式

② Intel 众核 (MIC) 计算

MIC 简介

MIC 计算运行方式

③ CUDA 计算与程序编译运行

CUDA 计算简介

CUDA 编译环境搭建

CUDA 程序编译

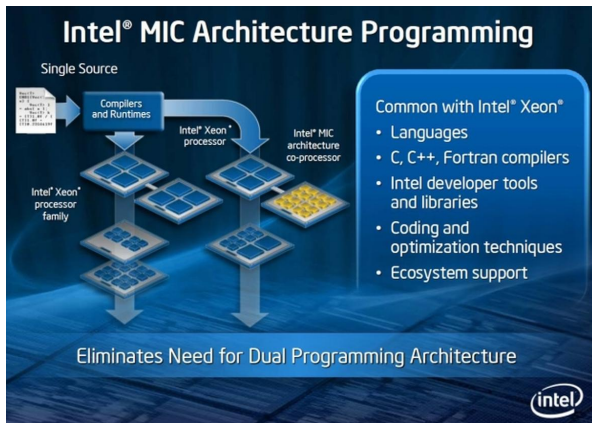
cuBLAS 函数库的使用

作业提交

Mang Intergrated Core, 译名：集成众核。

目前已经发展了两代：

- 2013 年发布 Knights Corner (KNC)
- 2016 年发布 Knight Landing (KNL)



目前超算中心具有 4 个 KNC 加速卡，对比 CPU 指标为：

表 1: CPU和MIC的指标对比

Intel Xeon	CPU	MIC
型号	Xeon E5-2650v2	Xeon Phi 7110P
核心数目	8	61
线程数目	16	244
核心频率	2.6GHz	1.238GHz
L1数据cache	8*32KB	61*32KB
L1指令cache	8*32KB	61*32KB
L2 cache	8*256KB	61*512KB
LLC Cache	20MB	/
内存	外置DDR3	8GB GDDR5
设计最大功耗	95W	300W

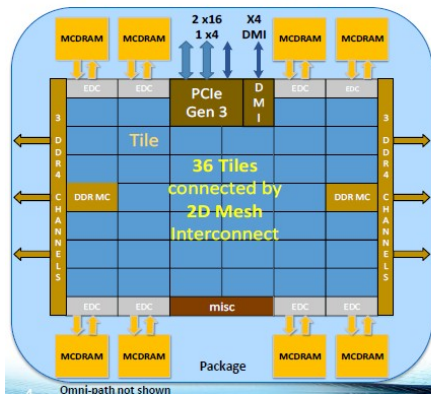
MIC 第二代: KNL

相比第一代 KNC 需要协助其他 CPU 完成加速, 第二代 KNL 可以作为独立的 CPU 使用。

Knights Landing Overview

TILE

2 VPU	CHA	2 VPU
Core	1MB L2	Core



Chip: 36 Tiles interconnected by 2D Mesh

Tile: 2 Cores + 2 VPU/core + 1 MB L2

Memory: MCDRAM: 16 GB on-package; High BW

DDR4: 6 channels @ 2400 up to 384GB

IO: 36 lanes PCIe Gen3. 4 lanes of DMI for chipset

Node: 1-Socket only

Fabric: Omni-Path on-package (not shown)

Vector Peak Perf: 3+TF DP and 6+TF SP Flops

Scalar Perf: ~3x over Knights Corner

Streams Triad (GB/s): MCDRAM : 400+; DDR: 90+

Source Intel: All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. 1 Binary Compatible with Intel Xeon processors using Haswell architecture. 2 Not supported. PEX: Bandwidth numbers are based on STREAM-like memory access pattern with 100MB of on-chip memory. Results have been estimated based on internal Intel analysis and are not intended to represent actual performance. 3 Purpose only. Any difference in system configuration or implementation may affect performance.

native 模式

登录到 MIC 卡中，之后操作与 CPU 没有关系，将 MIC 如同 CPU 一般工作运行：

```
[wszhang@node46 ~]$ssh mic0  
[wszhang@node46-mic0 ~]$cat /proc/cpuinfo | grep  
processor| wc -l  
244
```

但是在 MIC 系统中不能编译，编译需要在 CPU 端进行，只需编译时添加 `-mmic` 参数即可：

```
#icc -mmic -openmp -o pi-omp-native pi-omp.c
```

注意：编译 OpenMP 时需要链接 MIC 版的 `libiomp5.so`

offload 模式

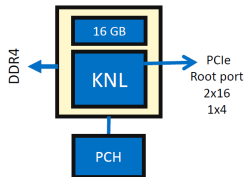
类似于使用 GPU 作为计算设备的程序，即程序在 CPU 端启动并运行，中间遇到需要加速计算的部分则转移到 MIC 上运行，运行结束后再返回 CPU 端，不需要登录 MIC 卡即可使用。使用的编程方法类似 OpenMP，添加一些指定 MIC 作为加速卡的语句如：

```
#pragma offload target(mic)
```

中科大的 MIC 节点位于 ChinaGrid 集群的 node45 与 node46。提交任务并无不同，利用“`bsub -q mic`”选项将任务提交至 mic 队列即可，如下：

```
bsub -q mic -o %j.log -e %j.err executable
```

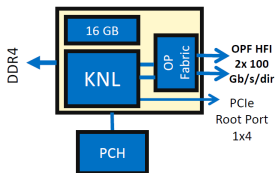
Knights Landing Products



KNL

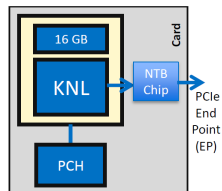
DDR Channels: 6
MCDRAM: up to 16 GB
Gen3 PCIe (Root port): 36 lanes

Self Boot Socket



KNL with Omni-Path

DDR Channels: 6
MCDRAM: up to 16 GB
Gen3 PCIe (Root port): 4 lanes
Omni-Path Fabric: 200 Gb/s/dir



KNL Card

No DDR Channels
MCDRAM: up to 16 GB
Gen3 PCIe (End point): 16 lanes
NTB Chip to create PCIe EP

PCIe Card

Potential future options subject to change without notice. Codenames.

All timeframes, features, products and dates are preliminary forecasts and subject to change without further notice.

① 异构众核计算

并行计算的现状与结构分类

异构计算的实现方式

② Intel 众核 (MIC) 计算

MIC 简介

MIC 计算运行方式

③ CUDA 计算与程序编译运行

CUDA 计算简介

CUDA 编译环境搭建

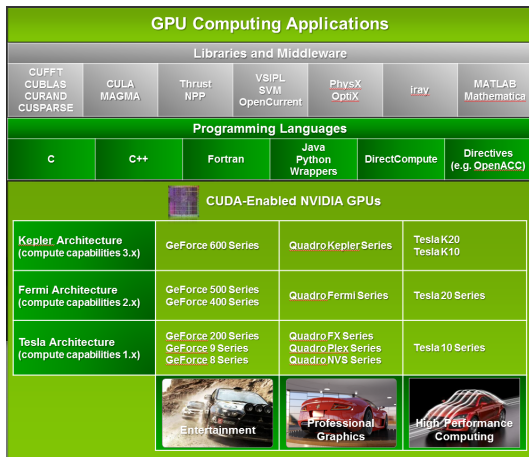
CUDA 程序编译

cuBLAS 函数库的使用

作业提交

GPU 应用

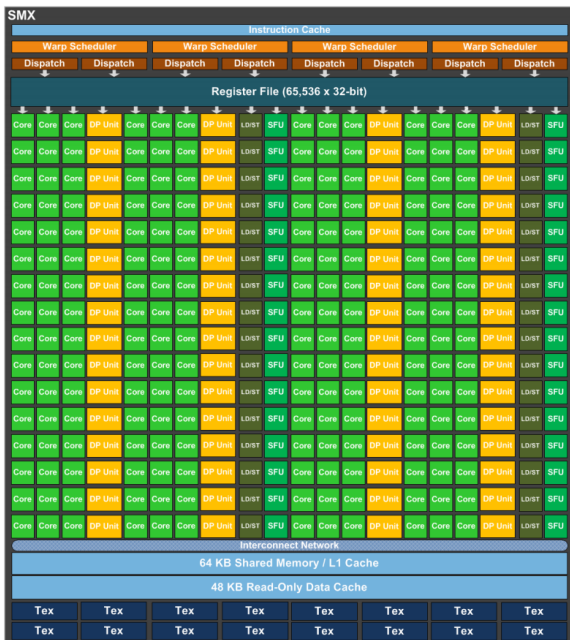
GPU 应用领域广泛，如教育科研，影视创作，智能机器学习等。在本校专注最多的量子化学领域，也陆续支持 GPU 加速，如主流的 Abinit, Quantum Espress, 以及 VASP (Now 2.5 to 4X Faster on Tesla K80)。



GPU 硬件微架构 (Kepler GK110)

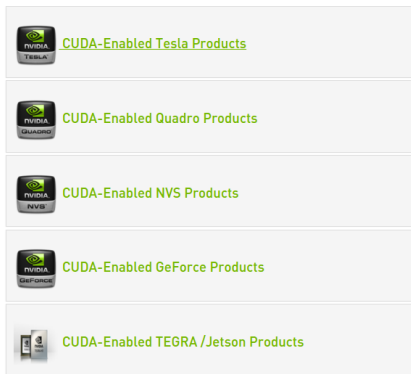


GPU SMX 192 SP and 64 DP units, 32 special function and 32 load/store units



CUDA 安装准备

- 确保具有 root 权限
- 确认已经安装 gcc 编译器
- 确认硬件支持 CUDA: <https://developer.nvidia.com/cuda-gpus>



CUDA 安装准备

- 确认系统支持 CUDA: <https://developer.nvidia.com/cuda-toolkit-archive>, 选择希望安装的 CUDA 版本, 下载自己系统对应的软件包, 或者下载通用的.run 后缀的软件包。

CUDA 7.0 Downloads

Please Note: There is a recommended patch for CUDA 7.0 which resolves an issue in the cuFFT library that can lead to incorrect results for certain inputs sizes less than or equal to 1920 in any dimension when cuFFTSetStream() is passed a **non-blocking stream** (e.g., one created using the cudaStreamNonBlocking flag of the CUDA Runtime API or the **CU_STREAM_NON_BLOCKING** flag of the CUDA Driver API).

	Windows	Linux x86	Linux POWER8	Mac OSX	
					Documentation
					Release Notes
					End User License Agreement
					CUDA Toolkit Overview
					Checksums
Version	Network Installer	Local Package Installer	Runfile Installer		
Fedora 21	RPM (3KB)	RPM (1GB)	RUN (1.1GB)		
OpenSUSE 13.2	RPM (3KB)	RPM (1GB)	RUN (1.1GB)		
OpenSUSE 13.1	RPM (3KB)	RPM (1GB)	RUN (1.1GB)		
RHEL 7 CentOS 7	RPM (10KB)	RPM (1GB)	RUN (1.1GB)		
RHEL 6 CentOS 6	RPM (18KB)	RPM (1GB)	RUN (1.1GB)		
CLFS 12	RPM (3KB)	RPM (1GB)	RUN		

1 禁用图形显示界面

- # service lightdm stop
或
- # init 3

2 运行安装程序

- # sh cuda_7.0.28_linux.run

3 按提示选择安装组件，并设置安装路径

- 可以全部 yes，并默认路径

① 配置 PATH

- # export PATH=\$PATH:/path-to-cuda/bin

② 配置 LD_LIBRARY_PATH

- # export LD_LIBRARY_PATH=\$LD_LIBRARY_PATH:/path-to-cuda/lib64
- # export LD_LIBRARY_PATH=\$LD_LIBRARY_PATH:/path-to-cuda/lib64/stubs (含-lcuda 库)

如果选择通用软件包，并执行

```
./cuda_7.0.28_linux.run -extract=~/.Path
```

可以看到会产生三个 CUDA 安装组件：

- 设备驱动
NVIDIA-Linux-x86_64-346.46.run
- CUDA Toolkit 程序开发包
cuda-linux64-rel-7.0.28-19326674.run
- 程序开发事例包
cuda-samples-linux-7.0.28-19326674.run

- Development Tools
 - NVCC, PTXAS, cuobjdump, Nsight Eclipse
- Libraries
 - cuBLAS, cuFFT, cuRAND, cuSPARES, NPP 等
- Tools
 - CUDA-GDB, CUDA-MEMCHECK, Visual Profiler, NVIDIA-SMI, NVML

nvidia-smi

使用此命令可以查看显卡驱动与运行状态，详细参数: `nvidia-smi -a`

```
[wszhang@node48 ~]$ nvidia-smi
Thu Nov 26 13:57:46 2015
+-----+-----+-----+-----+-----+-----+
| NVIDIA-SMI 352.39          | Driver Version: 352.39 | corresponding Driver Version, MPSS Version and Flash Version |
+-----+-----+-----+-----+-----+-----+
| GPU Name                   | Persistence-M | Bus-Id        | Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap| Memory-Usage  | GPU-Util  Compute M. |
+-----+-----+-----+-----+-----+-----+
| 0   Tesla K40c             | Off           | 0000:06:00.0 | Off    | 0                     |
| 23%   38C    P0             | 63W / 235W   | 22MiB / 11519MiB |         | 0%                     |
+-----+-----+-----+-----+-----+-----+
| 1   Tesla K40c             | Off           | 0000:81:00.0 | Off    | 0                     |
| 23%   36C    P0             | 65W / 235W   | 22MiB / 11519MiB |         | 48%                     |
+-----+-----+-----+-----+-----+-----+

Processes:
+-----+-----+-----+-----+-----+-----+
| GPU  PID  Type  Process name                        | GPU Memory Usage |
+-----+-----+-----+-----+-----+-----+
| No running processes found         |                   |
+-----+-----+-----+-----+-----+-----+
[wszhang@node48 ~]$
```

Occupancy Calculator

CUDA 占用率计算器可以计算在某个 CUDA 内核下 GPU 中多处理器的占用情况，即活动 Warp 数与 GPU 支持的 Warp 最大数比率。

Just follow steps 1, 2, and 3 below (or click here for help)

1.) Select Compute Capability (click): [click](#)

1.b) Select Shared Memory Size Config (bytes): [click](#)

1.c) Select Global Load Caching Mode:

2.) Enter your resource usage:

Threads Per Block: [click](#)

Registers Per Thread: [click](#)

Shared Memory Per Block (bytes): [click](#)

(Don't edit anything below this line)

3.) GPU Occupancy Data is displayed here and in the graphs:

Active Threads per Multiprocessor: [click](#)

Active Warps per Multiprocessor: [click](#)

Active Thread Blocks per Multiprocessor: [click](#)

Occupancy of each Multiprocessor: [click](#)

Physical Limits for GPU Compute Capability: 5.3

Threads per Warp	32
Max Warps per Multiprocessor	64
Max Thread Blocks per Multiprocessor	32
Max Threads per Multiprocessor	2048
Maximum Thread Size	1024
Registers per Multiprocessor	85536
Max Registers per Thread Block	32768
Max Registers per Thread	256
Shared Memory per Multiprocessor (bytes)	32768
Max Shared Memory per Block	32768
Register allocation unit size	256
Register allocation granularity	warp
Shared Memory allocation unit size	256
Warp allocation granularity	4

Allocated Resources

	Par Block	Limit Per SM	* Allocatable Blocks Per SM
Ways (Threads Per Block / Threads Per Warp)	8	64	8
Registers (Warp limit per SM due to per-warp reg count)	8	256	4
Shared Memory (bytes)	4096	32768	8

Note: SM is an abbreviation for Streaming Multiprocessor

Maximum Thread Blocks Per Multiprocessor: $\text{Blocks/SM} * \text{Warps/Block} = \text{Warps/SM}$

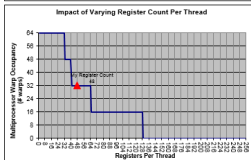
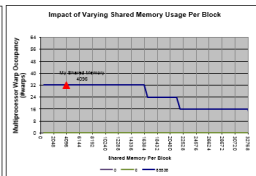
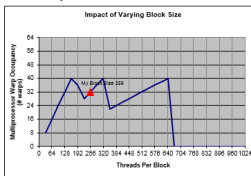
Limited by Max Warps per Multiprocessor	8	64	8
Limited by Registers per Multiprocessor	4	8	32
Limited by Shared Memory per Multiprocessor	8	8	8

Note: Occupancy limits is shown in orange

Physical Max Warps/SM = 64
Occupancy = 32 / 64 = 50%

CUDA Occupancy Calculator
Version: [Copyright and License](#)

Your chosen resource usage is indicated by the red triangle on the graphs. The other data points represent the range of possible block sizes, register counts, and shared memory usage per block.



集成开发界面

The screenshot displays the NVIDIA Nsight IDE interface for a session named "matrixMulCUBLAS - Nsight". The main window shows a timeline view of the application's execution, with various components and their durations visualized. The components include:

- Process "matrixMjIcUBLAS" (17274)
- Thread 2700855168
 - Runtime API
 - Driver API
 - Profiling Overhead
- [0] GeForce GT 740M
 - Context 1 (CUDA)
 - MemCpy (HtoD)
 - MemCpy (DtoH)
 - Compute
 - 100.0% sgemm_sm35_id...
 - Streams
 - DeFault

The Analysis panel at the bottom provides a summary of the results:

Results

- Low Memcpy/Compute Overlap** [0 ns / 1.641 ms = 0%]
The percentage of time when memcpy is being performed in parallel with compute is low. [More...](#)
- Low Kernel Concurency** [0 ns / 59.984 ms = 0%]
The percentage of time when two kernels are being executed in parallel is low. [More...](#)
- Low Memcpy Throughput** [1.622 GB/s avg, for memcpys accounting for 100% of all memcpy time]

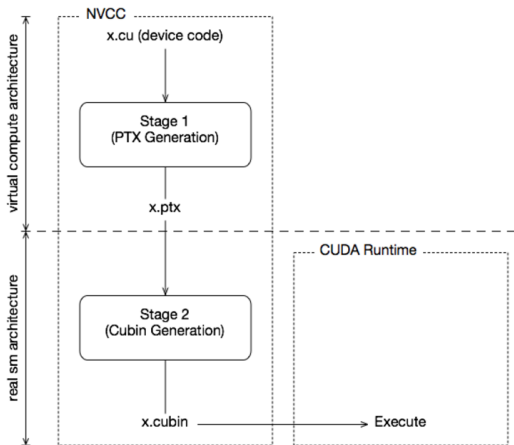
Input File Prefix	Description
.cu	CUDA source file, containing host code and device functions
.c	C source file
.cc, .cxx, .cpp	C++ source file
.gpu	GPU intermediate file
.ptx	PTX intermediate assembly file
.o, .obj	Object file
.a, .lib	Library file
.so	Shared object file

CUDA 编译选项

Phase	short nvcc Option	Default Output File Name
CUDA compilation to C/C++ source file	-cuda	.cpp.ii appended to source file name, as in x.cu.cpp.ii.
C/C++ preprocessing	-E	<result on standard output>
C/C++ compilation to object file	-c	suffix replaced by o on Linux/Mac, obj on Win
Cubin from CUDA source files	-cubin	Source file name with suffix replaced by cubin
Cubin from .gpu intermediate files	-cubin	Source file name with suffix replaced by cubin
Cubin from PTX intermediate files.	-cubin	Source file name with suffix replaced by cubin
PTX from CUDA source files	-ptx	Source file name with suffix replaced by ptx
PTX from .gpu intermediate files	-ptx	Source file name with suffix replaced by ptx
Fatbinary from source, PTX or cubin files	-fatbin	Source file name with suffix replaced by fatbin
GPU C code from CUDA source files	-gpu	Source file name with suffix replaced by gpu
Linking relocatable device code.	-dlink	a_dlink.obj on Win or a_dlink.o on other platforms
Cubin from linked relocatable device code.	-dlink-cubin	a_dlink.cubin
Fatbinary from linked relocatable device code	-dlink-fatbin	a_dlink.fatbin
Linking an executable	<no phase option>	
Constructing an object file archive, or library	-lib	a.lib on Windows or a.a on other platforms
make dependency generation	-M	<result on standard output>
Running an executable	-run	

CUDA 两步编译

为了程序兼容性，CUDA 编译被设计为两步，先在虚拟机架构下编译成类似汇编代码的 PTX 中间文件，然后在第二步中编译成最后的执行文件，PTX 中间件也可以运行时动态的再编译并运行。

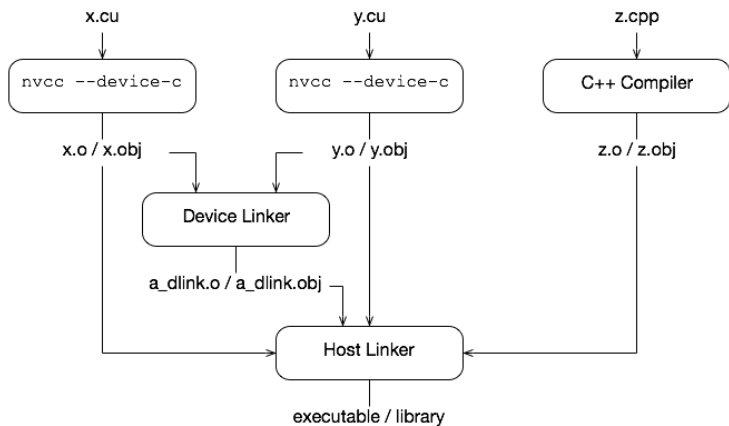


CUDA 中间虚拟架构

	Architecture	Feature
真实架构	sm_20	Basic features + Fermi support
	sm_30 and sm_32	+ Kepler support + Unified memory programming
	sm_35	+ Dynamic parallelism support
	sm_50, sm_52, and sm_53	+ Maxwell support
虚拟架构	compute_20	Basic features + Fermi support
	compute_30 and compute_32	+ Kepler support + Unified memory programming
	compute_35	+ Dynamic parallelism support
	compute_50, compute_52, and compute_53	+ Maxwell support

- `-gpu-architecture` (可简写为 `-arch`) arch
 - Specify the name of the class of NVIDIA virtual GPU architecture for which the CUDA input files must be compiled.
- `-gpu-code` (可简写为 `-code`) code
 - Specify the name of the NVIDIA GPU to assemble and optimize PTX for.

CUDA C/C++ 可以分别独立编译



CUDA 常用编译方法

一般 cuda 程序编译, 虚拟架构版本号需小于等于真实架构版本号

```
nvcc x.cu -gpu-architecture=compute_20 -gpu-code=sm_20
```

多种架构编译

```
nvcc x.cu -gpu-architecture=compute_30 -gpu-code=compute_30,sm_30,sm_35
```

缩写

```
nvcc x.cu -gpu-architecture=sm_35
```

```
nvcc x.cu -gpu-architecture=compute_30
```

分别等价于

```
nvcc x.cu -gpu-architecture=compute_35 -gpu-code=sm_35,compute_35
```

```
nvcc x.cu -gpu-architecture=compute_30 -gpu-code=compute_30
```

CUDA, C/C++ 分别编译后链接

```
nvcc -gpu-architecture=sm_20 -device-c a.cu b.cu
```

```
nvcc -gpu-architecture=sm_20 -device-link a.o b.o -output-file link.o
```

```
g++ a.o b.o link.o -library-path=<path> -library=cudart
```

CUDA 程序编译运行 (ThreadBlock.cu) I

```
nvcc -gencode=arch=compute_30,code=sm_30 ThreadBlock.cu -O2 -o ThreadBlock.o
```

```
1#include <stdio.h>
2//#include <cuda_runtime.h>
3
4__global__ void vadd ( float *a, float *b, float *c ,int nn ) {
5
6    int ii = 3 * blockIdx.x + threadIdx.x;
7    //int ii = blockDim.x * blockIdx.x + threadIdx.x;
8    //int ii = threadIdx.x ;
9    c[ii] = a[ii] + b[ii] ;
10}
11
12int main(void) {
13    cudaError_t err = cudaSuccess;
14
15    int nn = 6 ;
16
17    size_t size = nn*sizeof(float);
18    float *h_a = (float *)malloc(size);
19    float *h_b = (float *)malloc(size);
20    float *h_c = (float *)malloc(size);
21
22    for (int i=0; i< nn; ++i )
23    {
24        h_a[i] = rand()/(float)RAND_MAX;
```

CUDA 程序编译运行 (ThreadBlock.cu) II

```
25     h_b[i] = rand()/(float)RAND_MAX;
26 }
27
28 float *d_a = NULL;
29 err = cudaMalloc( (void **)&d_a, size );
30 float *d_b = NULL;
31 err = cudaMalloc( (void **)&d_b, size );
32 float *d_c = NULL;
33 err = cudaMalloc( (void **)&d_c, size );
34
35 err = cudaMemcpy( d_a, h_a, size, cudaMemcpyHostToDevice );
36 if ( err != cudaSuccess )
37 {
38     fprintf(stderr, "no. 1122311 \n" );
39 }
40
41 err = cudaMemcpy( d_b, h_b, size, cudaMemcpyHostToDevice );
42 if ( err != cudaSuccess )
43 {
44     fprintf(stderr, "no. 1122344 \n" );
45 }
46
47 vadd<<<2,3>>>(d_a,d_b,d_c, nn);
48
49 err = cudaMemcpy(h_c,d_c,size,cudaMemcpyDeviceToHost);
50
51 for ( int i = 0; i < nn; ++i )
52 {
```

CUDA 程序编译运行 (ThreadBlock.cu) III

```
53     if ( fabs(h_a[i] + h_b[i] - h_c[i] ) > 1e-4 )
54     {
55         fprintf( stderr, " failed at element %d! \n", i );
56     } else
57     {
58         fprintf( stderr, " succeed at element %d! \n", i );
59     }
60 }
61
62 err = cudaFree(d_a);
63 err = cudaFree(d_b);
64 err = cudaFree(d_c);
65
66 free(h_a);
67 free(h_b);
68 free(h_c);
69
70 return 0;
71 }
```

源文件

```
fortran.c #This file contains example Fortran bindings for the CUBLAS library  
sgemm_speed.f90
```

```
gcc -O3 -DCUBLAS_USE_THUNKING -I/opt/cuda-7.5/include -c fortran.c  
ifort -o sgemm_speed_cublas -O3 -fpp -DCUBLAS sgemm_speed.f90 fortran.o  
-L/opt/cuda-7.5/lib64 -lcublas
```

sgemm_speed.f90

```
1!  
2! Simple Fortan90 program that multiplies 2 square matrices calling Sgemm  
3! C = alpha A*B + beta C  
4!  
5program matrix_multiply  
6  
7implicit none  
8  
9! Define the floating point kind to be single_precision  
10integer, parameter :: fp_kind = kind(0.0)  
11  
12! Define  
13real (fp_kind), dimension(:,:), allocatable :: A, B, C
```

Fortran 编译链接 cuBLAS 函数库 II

```
14 real ::      time_start,time_end
15 real (fp_kind)::      alpha=1._fp_kind,beta=1._fp_kind, c_right
16 integer::      i,j,m1,m2
17
18
19 do m1=512,10240,512
20
21  allocate(A(m1,m1))
22  allocate(B(m1,m1))
23  allocate(C(m1,m1))
24
25  ! Initialize the matrices A,B and C
26  A=1._fp_kind
27  B=2._fp_kind
28  C=3._fp_kind
29
30  ! With the prescribed inputs, each element of the C matrix should be equal to
      c_right
31  c_right= 2._fp_kind*m1+3._fp_kind
32
33 ! Compute the matrix product computation
34 call cpu_time(time_start)
35
36 #ifdef CUBLAS
37  call cublas_SGEMM ('n', 'n', m1,m1,m1, alpha ,A,m1,B,m1, beta ,C,m1)
38 #else
39  call SGEMM ('n', 'n', m1,m1,m1, alpha ,A,m1,B,m1, beta ,C,m1)
40 #endif
```



```
41
42 call cpu_time(time_end)
43
44! Print timing information
45 print "(i5,1x,a,1x,f8.4,2x,a,f12.4)", m1, " time =",time_end-time_start, "
      MFLOPS=",1.e-6*2._fp_kind*m1*m1*m1/(time_end-time_start)
46
47! check the result
48   do j=1,m1
49     do i=1,m1
50       if ( abs(c(i,j)- c_right ) .gt. 1.d-8 ) then
51         print *, "sgemm failed", i,j, abs(c(i,j)- c_right )
52         exit
53       end if
54     end do
55   end do
56
57 deallocate(A,B,C)
58end do
59
60end program matrix_multiply
```

C 程序调用 cuBLAS 函数库 I

cusblas.h

```
1/* _____ CUDA _____ */
2#ifdef __CUDA
3
4/* Includes, cuda */
5#include <cuda_runtime.h>
6#include <cusblas_v2.h>
7#include <host_defines.h>
8#include "cuda.h" // not necessary when compiling with nvcc
9#include "cuComplex.h"
10
11#ifdef __single
12typedef cuComplex cucomplex;
13#define makecucomplex make_cuComplex;
14#else
15typedef cuDoubleComplex cucomplex;
16#define makecucomplex make_cuDoubleComplex;
17#endif
18
19 // Error handling macro
20#define CUDA_CHECK(call) \
21     if((call) != cudaSuccess) { \
22         cudaError_t err = cudaGetLastError(); \
23         cerr << "CUDA error calling \""#call"\", code is " << err << endl; \
```

```
24         my_abort(err); }
25
26 #define cuCHECK(call) \
27     if((call) != cudaSuccess) { \
28         cudaError_t err = cudaGetLastError(); \
29         cerr << "CUDA error calling \""#call"\", code is " << err << endl; \
30         my_abort(err); }
31
32 #define cublacsCHECK(call) \
33     if((call) != CUBLAS_STATUS_SUCCESS) { \
34         cerr << "CUDA error calling \""#call"\n" " << endl; }
35 #endif
```

C 程序调用 cuBLAS 函数库 I

cublas.cpp

```
1
2#ifdef __CUDA
3
4void cuPAexpm( usrcomplex * h_expMp, usrcomplex * h_Mp, cucomplex * d_expMp,
   cucomplex * d_Mp, cucomplex * d_carray,
5   cucomplex * d_czero, cucomplex * d_cone, cucomplex * d_cngone, const int
   & mval)
6{
7
8   const cuDoubleComplex ddone      = make_cuDoubleComplex(1.0,0.0);
9   const cuDoubleComplex ddzero     = make_cuDoubleComplex(0.0,0.0);
10
11   cublasHandle_t handle1;
12   cublasStatus_t status;
13   cublasCHECK(cublasCreate_v2(&handle1) );
14
15#ifdef __CUDA_NORM1
16   cudaMemcpy(d_Mp , h_Mp , MMsizes, cudaMemcpyHostToDevice);
17
18#endif
19
20   cucomplex * M2p;
21   cucomplex * M4p;
```

C 程序调用 cuBLAS 函数库 II

```
22     cucomplex * M6p;
23     cucomplex * uM;
24     cucomplex * vM;
25
26
27     cudaMalloc( (void **)&M2p, MMSize);
28     cudaMalloc( (void **)&M4p, MMSize);
29     cudaMalloc( (void **)&M6p, MMSize);
30     cudaMalloc( (void **)&uM, MMSize);
31     cudaMalloc( (void **)&vM, MMSize);
32
33     cublasCHECK(cublasZgemm(handle1,CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
34                          MMdim, MMdim,
35                          d_cone, d_Mp, MMdim, d_Mp, MMdim, d_czero, M2p, MMdim) );
36     cublasCHECK(cublasZgemm(handle1,CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
37                          MMdim, MMdim,
38                          d_cone, M2p, MMdim, M2p, MMdim, d_czero, M4p, MMdim) );
39     cublasCHECK(cublasZgemm(handle1,CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
40                          MMdim, MMdim,
41                          d_cone, M4p, MMdim, M2p, MMdim, d_czero, M6p, MMdim) );
42
43     cucomplex * VmUM;
44     cucomplex * invVmUM ;
45     cucomplex * VpUM ;
46     for (int i=0; i<1; i++ )
47     {
48         cudaMalloc( (void **)&(VmUM), MMSize );
49         cudaMalloc( (void **)&(invVmUM), MMSize );
50     }
```

C 程序调用 cuBLAS 函数库 III

```
47         cudaMalloc( (void **)&(VpUM),      MMSize );
48     }
49
50 #ifdef __single
51 #else
52     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+12, M6p, 1, VmUM,
53         1) );
54     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+10, M4p, 1, VmUM,
55         1) );
56     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+8,  M2p, 1, VmUM,
57         1) );
58     cublasCHECK(cublasZgemm(handle1, CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
59         MMdim, MMdim,
60         d_cone, M6p,  MMdim, VmUM, MMdim, d_czero, vM, MMdim) );
61     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+6,  M6p,  1, vM,
62         1) );
63     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+4,  M4p,  1, vM,
64         1) );
65     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+2,  M2p,  1, vM,
66         1) );
67     cublasCHECK(cublasZaxpy(handle1, MMdim , d_carray+0,  d_IMp, MMdim
68         +1, vM, MMdim+1) );
69     cublasCHECK(cublasZcopy(handle1, MMdim2, vM, 1, VmUM, 1) );
70 #endif
71 #ifdef __single
```

C 程序调用 cuBLAS 函数库 IV

```
67 #else
68     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+13, M6p, 1, uM, 1)
69                 );
70     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+11, M4p, 1, uM, 1)
71                 );
72     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+9, M2p, 1, uM, 1)
73                 );
74     cublasCHECK(cublasZgemm(handle1, CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
75                             MMdim, MMdim,
76                             d_cone, M6p, MMdim, uM, MMdim, d_czero, VpUM, MMdim) );
77     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+7, M6p, 1, VpUM
78                 , 1) );
79     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+5, M4p, 1, VpUM
80                 , 1) );
81     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_carray+3, M2p, 1, VpUM
82                 , 1) );
83     cublasCHECK(cublasZaxpy(handle1, MMdim, d_carray+1, d_Imp, MMdim
84                 +1, VpUM, MMdim+1) );
85     cublasCHECK(cublasZcopy(handle1, MMdim2, vM, 1, VmUM, 1) );
86     cublasCHECK(cublasZgemm(handle1, CUBLAS_OP_N, CUBLAS_OP_N, MMdim,
87                             MMdim, MMdim,
88                             d_cone, d_Mp, MMdim, VpUM, MMdim, d_czero, uM,
89                             MMdim) );
90     cublasCHECK(cublasZcopy(handle1, MMdim2, uM, 1, VpUM, 1) );
91 #endif
92
93
```

C 程序调用 cuBLAS 函数库 V

```
84     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_cngone, uM, 1, VmUM, 1)
85     );
86     cublasCHECK(cublasZaxpy(handle1, MMdim2, d_cone, vM, 1, VpUM, 1)
87     );
88
89     cublasCHECK( cublasGetVector(MMdim2, sizeof(cucomplex) , VmUM, 1,
90     h_expMp, 1) );
91
92     #ifdef __single
93     #else
94     const int batchsize = 1;
95     const int NNdim=MMdim;
96
97     int *PArray, *InfoArray;
98     cudacall(cudaMalloc(&PArray, sizeof(int) * batchsize * NNdim ));
99     cudacall(cudaMalloc(&InfoArray, sizeof(int) * batchsize ));
100     int h_InfoArray[batchsize];
101
102     //Random matrix with full pivots
103     cucomplex full_pivot[4*4] = { 0.5, 3, 4, 4,
104     1, 3, 10, 4,
105     4, 9, 16, 4,
106     4, 4, 4, 4};
107
108     //Almost same as above matrix with first pivot zero
109     cucomplex zero_pivot[3*3] = { 0, 3, 4,
110     1, 3, 10,
111     4, 9, 16 };
```


C 程序调用 cuBLAS 函数库 VI

```
009     cucomplex ** AA_hh = (cucomplex **)malloc( sizeof(cucomplex *) *
010             batchsize );
011     cucomplex ** AA_hd = (cucomplex **)malloc( sizeof(cucomplex *) *
012             batchsize );
013     cucomplex *  AA_d ;
014     cucomplex ** AA_dd;
015     AA_d = VmUM;
016     cudaMalloc( &AA_dd, sizeof(cucomplex*) * batchsize
017             );
018     AA_hd[0] = AA_d;
019     for (int i = 1; i < batchsize; i++)
020         AA_hd[i] = AA_hd[i-1] + NNdim*NNdim;
021
022     cudaMemcpy(AA_dd, AA_hd, sizeof(cucomplex *) * batchsize,
023             cudaMemcpyHostToDevice );
024
025     cucomplex ** CC_hh = (cucomplex **)malloc(sizeof(cucomplex *) *
026             batchsize );
027     cucomplex ** CC_hd = (cucomplex **)malloc(sizeof(cucomplex *) *
028             batchsize );
029     cucomplex *  CC_h  = (cucomplex * )malloc(sizeof(cucomplex ) * NNdim
030             *NNdim * batchsize);
031     cucomplex *  CC_d ;
032     cucomplex ** CC_dd;
033     cudaMalloc( &CC_d,  sizeof(cucomplex ) * batchsize *NNdim*NNdim );
034     cudaMalloc( &CC_dd, sizeof(cucomplex*) * batchsize
035             );
036     for (int i = 0; i < batchsize; i++)
037         CC_hh[i] = CC_h + (i*NNdim*NNdim);
```

C 程序调用 cuBLAS 函数库 VII

```
131 CC_hd[0] = CC_d;
132 for (int i = 1; i < batchsize; i++)
133     CC_hd[i] = CC_hd[i-1] + i*NNdim*NNdim;
134 cudaMemcpy(CC_dd, CC_hd, sizeof(cucomplex *) * batchsize,
135             cudaMemcpyHostToDevice );
136
137 cublasZgetrfBatched(handle1, NNdim, AA_dd, NNdim, PArray, InfoArray,
138                     batchsize);
139 cudaMemcpy(h_InfoArray, InfoArray, batchsize*sizeof(int),
140            cudaMemcpyDeviceToHost);
141 for (int i=0; i<batchsize; i++)
142 {
143     if(h_InfoArray[i] != 0 )
144     {
145         cout << "cublas*getrfBatched is failed! i = " <<
146             i << endl;
147         cudaDeviceReset();
148     }
149 }
150
151 cublasZgetriBatched(handle1, NNdim, (const cucomplex **)AA_dd, NNdim,
152                     PArray, CC_dd, NNdim, InfoArray, batchsize);
153 cudaMemcpy(h_InfoArray, InfoArray, batchsize*sizeof(int),
154            cudaMemcpyDeviceToHost);
155 for (int i=0; i<batchsize; i++)
156 {
157     if(h_InfoArray[i] != 0 )
158     {
```

C 程序调用 cuBLAS 函数库 VIII

```
153         cout << "cublas*getriBatched is failed! i = " <<
154             i << endl;
155         cudaDeviceReset();
156     }
157 }
158 cublasZgemm(handle1, CUBLAS_OP_N, CUBLAS_OP_N, MMdim, MMdim, MMdim,
159             d_cone, CC_d, MMdim, VpUM, MMdim, d_czero,
160             d_expMp, MMdim);
161 #endif
162
163     cudaMemcpy( h_expMp, d_expMp, batchsize*MMSize,
164               cudaMemcpyDeviceToHost);
165
166     cudaFree(PArray); cudaFree(InfoArray);
167     cudaFree(AA_d); cudaFree(AA_dd); free(AA_hd);
168     cudaFree(CC_d); cudaFree(CC_dd); free(CC_hd); free(CC_hh);
169
170     cudaFree(M2p); cudaFree(M4p); cudaFree(M6p); cudaFree(vM); cudaFree(
171         uM);
172     cudaFree(VmUM); cudaFree(invVmUM); cudaFree(VpUM);
173     cublasDestroy_v2(handle1);
174     cudaDeviceReset();
175 }
176 #endif
```

提交作业 (多 GPU 程序) I

```
bsub -m nodename -q k40 -e %j.err -o %j.log ./simpleMultiGPU
```

```
1/*
2 * Copyright 1993–2015 NVIDIA Corporation. All rights reserved.
3 *
4 * Please refer to the NVIDIA end user license agreement (EULA) associated
5 * with this source code for terms and conditions that govern your use of
6 * this software. Any use, reproduction, disclosure, or distribution of
7 * this software and related documentation outside the terms of the EULA
8 * is strictly prohibited.
9 *
10 */
11
12/*
13 * This application demonstrates how to use the CUDA API to use multiple GPUs,
14 * with an emphasis on simple illustration of the techniques (not on performance)
15 *
16 * Note that in order to detect multiple GPUs in your system you have to disable
17 * SLI in the nvidia control panel. Otherwise only one GPU is visible to the
18 * application. On the other side, you can still extend your desktop to screens
19 * attached to both GPUs.
20 */
21
22// System includes
23#include <stdio.h>
```

提交作业 (多 GPU 程序) II

```
24#include <assert.h>
25
26// CUDA runtime
27#include <cuda_runtime.h>
28
29// helper functions and utilities to work with CUDA
30#include <helper_functions.h>
31#include <helper_cuda.h>
32#include <timer.h>
33
34#ifndef MAX
35#define MAX(a,b) (a > b ? a : b)
36#endif
37
38#include "simpleMultiGPU.h"
39
40////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
41// Data configuration
42////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
43const int MAX_GPU_COUNT = 32;
44const int DATA_N       = 1048576 * 32;
45
46////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
47// Simple reduction kernel.
48// Refer to the 'reduction' CUDA Sample describing
49// reduction optimization strategies
50////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////////
51__global__ static void reduceKernel(float *d_Result, float *d_Input, int N)
```

提交作业 (多 GPU 程序) III

```
52{
53    const int      tid = blockIdx.x * blockDim.x + threadIdx.x;
54    const int threadN = gridDim.x * blockDim.x;
55    float sum = 0;
56
57    for (int pos = tid; pos < N; pos += threadN)
58        sum += d_Input[pos];
59
60    d_Result[tid] = sum;
61}
62
63 ////////////////////////////////////////////////////////////////////////////////////////////
64 // Program main
65 ////////////////////////////////////////////////////////////////////////////////////////////
66 int main(int argc, char **argv)
67 {
68     // Solver config
69     TGPUplan      plan[MAX_GPU_COUNT];
70
71     // GPU reduction results
72     float      h_SumGPU[MAX_GPU_COUNT];
73
74     float sumGPU;
75     double sumCPU, diff;
76
77     int i, j, gpuBase, GPU_N;
78
79     const int  BLOCK_N = 32;
```

提交作业 (多 GPU 程序) IV

```
80  const int THREAD_N = 256;
81  const int  ACCUM_N = BLOCK_N * THREAD_N;
82
83  printf("Starting simpleMultiGPU\n");
84  checkCudaErrors(cudaGetDeviceCount(&GPU_N));
85
86  if (GPU_N > MAX_GPU_COUNT)
87  {
88      GPU_N = MAX_GPU_COUNT;
89  }
90
91  printf("CUDA-capable device count: %i\n", GPU_N);
92
93  printf("Generating input data...\n\n");
94
95  //Subdividing input data across GPUs
96  //Get data sizes for each GPU
97  for (i = 0; i < GPU_N; i++)
98  {
99      plan[i].dataN = DATA_N / GPU_N;
100 }
101
102 //Take into account "odd" data sizes
103 for (i = 0; i < DATA_N % GPU_N; i++)
104 {
105     plan[i].dataN++;
106 }
107
```

提交作业 (多 GPU 程序) V

```
08 //Assign data ranges to GPUs
09 gpuBase = 0;
10
11 for (i = 0; i < GPU_N; i++)
12 {
13     plan[i].h_Sum = h_SumGPU + i;
14     gpuBase += plan[i].dataN;
15 }
16
17 //Create streams for issuing GPU command asynchronously and allocate memory (
18 //GPU and System page-locked)
19 for (i = 0; i < GPU_N; i++)
20 {
21     checkCudaErrors(cudaSetDevice(i));
22     checkCudaErrors(cudaStreamCreate(&plan[i].stream));
23     //Allocate memory
24     checkCudaErrors(cudaMalloc((void **)&plan[i].d_Data, plan[i].dataN *
25         sizeof(float)));
26     checkCudaErrors(cudaMalloc((void **)&plan[i].d_Sum, ACCUM_N * sizeof(
27         float)));
28     checkCudaErrors(cudaMallocHost((void **)&plan[i].h_Sum_from_device,
29         ACCUM_N * sizeof(float)));
30     checkCudaErrors(cudaMallocHost((void **)&plan[i].h_Data, plan[i].dataN *
31         sizeof(float)));
32
33     for (j = 0; j < plan[i].dataN; j++)
34     {
35         plan[i].h_Data[j] = (float)rand() / (float)RAND_MAX;
36     }
37 }
```


提交作业 (多 GPU 程序) VI

```
131     }
132 }
133
134 //Start timing and compute on GPU(s)
135 printf("Computing with %d GPUs...\n", GPU_N);
136 StartTimer();
137
138 //Copy data to GPU, launch the kernel and copy data back. All asynchronously
139 for (i = 0; i < GPU_N; i++)
140 {
141     //Set device
142     checkCudaErrors(cudaSetDevice(i));
143
144     //Copy input data from CPU
145     checkCudaErrors(cudaMemcpyAsync(plan[i].d_Data, plan[i].h_Data, plan[i].
146         dataN * sizeof(float), cudaMemcpyHostToDevice, plan[i].stream));
147
148     //Perform GPU computations
149     reduceKernel<<<BLOCK_N, THREAD_N, 0, plan[i].stream>>>(plan[i].d_Sum,
150         plan[i].d_Data, plan[i].dataN);
151     getLastCudaError("reduceKernel() execution failed.\n");
152
153     //Read back GPU results
154     checkCudaErrors(cudaMemcpyAsync(plan[i].h_Sum_from_device, plan[i].d_Sum,
155         ACCUM_N * sizeof(float), cudaMemcpyDeviceToHost, plan[i].stream));
156 }
157
158 //Process GPU results
```

提交作业 (多 GPU 程序) VII

```
056 for (i = 0; i < GPU_N; i++)
057 {
058     float sum;
059
060     //Set device
061     checkCudaErrors(cudaSetDevice(i));
062
063     //Wait for all operations to finish
064     cudaStreamSynchronize(plan[i].stream);
065
066     //Finalize GPU reduction for current subvector
067     sum = 0;
068
069     for (j = 0; j < ACCUM_N; j++)
070     {
071         sum += plan[i].h_Sum_from_device[j];
072     }
073
074     *(plan[i].h_Sum) = (float)sum;
075
076     //Shut down this GPU
077     checkCudaErrors(cudaFreeHost(plan[i].h_Sum_from_device));
078     checkCudaErrors(cudaFree(plan[i].d_Sum));
079     checkCudaErrors(cudaFree(plan[i].d_Data));
080     checkCudaErrors(cudaStreamDestroy(plan[i].stream));
081 }
082
083 sumGPU = 0;
```

提交作业 (多 GPU 程序) VIII

```
184
185 for (i = 0; i < GPU_N; i++)
186 {
187     sumGPU += h_SumGPU[i];
188 }
189
190 printf(" GPU Processing time: %f (ms)\n\n", GetTimer());
191
192 // Compute on Host CPU
193 printf("Computing with Host CPU...\n\n");
194
195 sumCPU = 0;
196
197 for (i = 0; i < GPU_N; i++)
198 {
199     for (j = 0; j < plan[i].dataN; j++)
200     {
201         sumCPU += plan[i].h_Data[j];
202     }
203 }
204
205 // Compare GPU and CPU results
206 printf("Comparing GPU and Host CPU results...\n");
207 diff = fabs(sumCPU - sumGPU) / fabs(sumCPU);
208 printf(" GPU sum: %f\n CPU sum: %f\n", sumGPU, sumCPU);
209 printf(" Relative difference: %E \n\n", diff);
210
211 // Cleanup and shutdown
```

提交作业 (多 GPU 程序) IX

```
12 for (i = 0; i < GPU_N; i++)
13 {
14     checkCudaErrors(cudaSetDevice(i));
15     checkCudaErrors(cudaFreeHost(plan[i].h_Data));
16
17     // cudaDeviceReset causes the driver to clean up all state. While
18     // not mandatory in normal operation, it is good practice. It is also
19     // needed to ensure correct operation when the application is being
20     // profiled. Calling cudaDeviceReset causes all profile data to be
21     // flushed before the application exits
22     cudaDeviceReset();
23 }
24
25 exit((diff < 1e-5) ? EXIT_SUCCESS : EXIT_FAILURE);
26 }
```

作业状态输出 |

```
1Sender: LSF System <lsfadmin@node48>
2Subject: Job 27734: <./simpleMultiGPU> Done
3
4Job <./simpleMultiGPU> was submitted from host <node48> by user <wszhang> in
  cluster <chinagrid>.
5Job was executed on host(s) <node48>, in queue <k40>, as user <wszhang> in
  cluster <chinagrid>.
6</home/nic/wszhang> was used as the home directory.
7</home/nic/wszhang/OPT/expm/cuda_samples/0_Simple/simpleMultiGPU> was used as the
  working directory.
8Started at Wed Nov 25 21:03:01 2015
9Results reported at Wed Nov 25 21:03:15 2015
10
11Your job looked like:
12
13-----
14# LSBATCH: User input
15./simpleMultiGPU
16-----
17
18Successfully completed.
19
20Resource usage summary:
21
22CPU time      :          4.56 sec.
23Max Memory    :           1 MB
24Max Swap      :           32 MB
```

作业状态输出 II

```
25
26 Max Processes : 1
27 Max Threads : 1
28
29 The output (if any) follows:
30
31 Starting simpleMultiGPU
32 CUDA-capable device count: 2
33 Generating input data...
34
35 Computing with 2 GPUs...
36 GPU Processing time: 10.902000 (ms)
37
38 Computing with Host CPU...
39
40 Comparing GPU and Host CPU results...
41 GPU sum: 16777280.000000
42 CPU sum: 16777294.395033
43 Relative difference: 8.580068E-07
44
45 PS:
46 Read file <27734.err> for stderr output of this job.
```

- **中国科学技术大学超算中心:**

办公室科大东区新图书馆一楼东侧 126 室

电话:0551-63602248

信箱:sccadmin@ustc.edu.cn

主页:<http://scc.ustc.edu.cn>