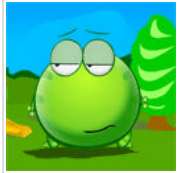登录 | 注册

# 网络资源是无限的

目录视图    摘要视图    RSS 订阅

**个人资料**

**fengbingchun**

访问：2252522次
积分：25003
等级：BLOG 7
排名：第202名

原创：341篇 转载：144篇
译文：0篇 评论：1434条

**文章分类**

Android (9)
ActiveX (18)
Bar Code (16)
Caffe (20)
C# (5)
CImg (4)
Contour Detection (9)
CxImage (6)
Code::Blocks (3)
Cloud Computing (1)
C/C++ (82)
CUDA (10)
CMake (3)
Design Patterns (25)
Database/Dataset (4)
Deep Learning (9)
Eclipse (3)
Emgu CV (1)
Eigen (1)
FFmpeg (1)
Feature Extraction (1)
FreeType (1)
Face (8)
GPU (3)
Git (3)
GCC (1)
GDAL (5)

CSDN学院招募微信小程序讲师啦    程序员简历优化指南！    【观点】移动原生App开发 PK HTML 5开发    云端应用征文大赛，秀绝招，赢无人机！

## 《GPU高性能编程CUDA实战》中代码整理

2015-05-24 19:35    5081人阅读    评论(2)    收藏    举报

分类：    CUDA（10）

CUDA架构专门为GPU计算设计了一种全新的模块，目的是减轻早期GPU计算中存在的一些限制，而正是这些限制使得之前的GPU在通用计算中没有得到广泛的应用。

使用**CUDA C**来编写代码的前提条件包括：(1)、支持CUDA的图形处理器，即由NVIDIA推出的GPU显卡，要求显存超过256MB；(2)、NVIDIA设备驱动程序，用于实现应用程序与支持CUDA的硬件之间的通信，确保安装最新的驱动程序，注意选择与开发环境相符的图形卡和操作系统；(3)、CUDA开发工具箱即CUDA Toolkit，此工具箱中包括一个编译GPU代码的编译器；(4)、标准C编译器，即CPU编译器。CUDA C应用程序将在两个不同的处理器上执行计算，因此需要两个编译器。其中一个编译器为GPU编译代码，而另一个为CPU编译代码。

一般，将**CPU**以及系统的内存称为主机**(Host)**，而将**GPU**及其内存称为设备**(Device)**。在**GPU**设备上执行的函数通常称为核函数**(Kernel)**。

cudaMalloc函数使用限制总结：(1)、可以将cudaMalloc()分配的指针传递给在设备上执行的函数；(2)、可以在设备代码中使用cudaMalloc()分配的指针进行内存读/写操作；(3)、可以将cudaMalloc()分配的指针传递给在主机上执行的函数；(4)、不能在主机代码中使用cudaMalloc()分配的指针进行内存读/写操作。

不能使用标准C的free()函数来释放cudaMalloc()分配的内存；要释放cudaMalloc()分配的内存，需要调用cudaFree()。

设备指针的使用方式与标准**C**中指针的使用方式完全一样。主机指针只能访问主机代码中的内存，而设备指针也只能访问设备代码中的内存。

在主机代码中可以通过调用cudaMemcpy()来访问设备上的内存。

有可能在单块卡上包含了两个或多个GPU。

在集成的**GPU**上运行代码，可以与**CPU**共享内存。

计算功能集的版本为**1.3**或者更高的显卡才能支持双精度浮点数的计算。

尖括号的第一个参数表示设备在执行核函数时使用的并行线程块的数量，

并行线程块集合也称为一个线程格(Grid)。线程格既可以是一维的线程块集合，也可以是二维的线程块集合。

**GPU**有着完善的内存管理机制，它将强行结束所有违反内存访问规则的进程。

在启动线程块数组时，数组每一维的最大数量都不能超过65535.这是一种硬件限制,如果启动的线程块数量超过了这个限值,那么程序将运行失败。

CUDA运行时将线程块(Block)分解为多个线程。当需要启动多个并行线程块时，只需将尖括号中的第一个参数由1改为想要启动的线程块数量。在尖括号中，第二个参数表示CUDA运行时在每个线程块中创建的线程数量。

**Free Codes**

硬件将线程块的数量限制为不超过**65535**.同样，对于启动核函数时每个线程块中的线程数量，硬件也进行了限制。具体来说，最大的线程数量不能超过设备属性结构中maxThreadsPerBlock域的值。这个值并不固定，有的是512，有的是1024.

内置变量**blockDim**，对于所有线程块来说，这个变量是一个常数，保存的是线程块中每一维的线程数量。

内置变量**gridDim**，对于所有线程块来说，这个变量是一个常数，用来保存线程格每一维的大小，即每个线程格中线程块的数量。

内置变量**blockIdx**，变量中包含的值就是当前执行设备代码的线程块的索引。

内置变量**threadIdx**，变量中包含的值就是当前执行设备代码的线程索引。

CUDA运行时允许启动一个二维线程格，并且线程格中的每个线程块都是一个三维的线程数组。

**CUDA C**支持共享内存：可以将CUDA C的关键字__share__添加到变量声明中，这将使这个变量驻留在共享内存中。CUDA C编译器对共享内存中的变量与普通变量将分别采取不同的处理方式。

CUDA架构将确保，除非线程块中的每个线程都执行了__syncthreads()，否则没有任何线程能执行__syncthreads()之后的指令。

由于在GPU上包含有数百个数学计算单元，因此性能瓶颈通常并不在于芯片的数学计算吞吐量，而是在于芯片的内存带宽。

常量内存用于保存在核函数执行期间不会发生变化的数据。NVIDIA硬件提供了64KB的常量内存，并且对常量内存采取了不同于标准全局内存的处理方式。在某些情况下，用常量内存来替换全局内存能有效地减少内存带宽。要使用常量内存，需在变量前面加上__constant__关键字。

在CUDA架构中，线程束是指一个包含**32**个线程的集合，这个线程集合被"编织在一起"并且以"步调一致(Lockstep)"的形式执行。在程序中的每一行，线程束中的每个线程都将在不同的数据上执行相同的指令。

纹理内存是在**CUDA C**程序中可以使用的另一种只读内存。与常量内存类似的是，纹理内存同样缓存在芯片上，因此在某些情况中，它能够减少对内存的请求并提供更高效的内存带宽。纹理缓存是专门为那些在内存访问模式中存在大量空间局部性(Spatial Locality)的图形应用程序而设计的。

NVIDIA将GPU支持的各种功能统称为计算功能集(Compute Capability)。高版本计算功能集是低版本计算功能集的超集。

只有**1.1**或者更高版本的**GPU**计算功能集才能支持全局内存上的原子操作。此外，只有1.2或者更高版本的GPU计算功能集才能支持共享内存上的原子操作。CUDA C支持多种原子操作。

C库函数malloc函数将分配标准的，可分页的(Pageble)主机内存。而cudaHostAlloc函数将分配页锁定的主机内存。页锁定内存也称为固定内存(Pinned Memory)或者不可分页内存，它有一个重要的属性：操作系统将不会对这块内存分页并交换到磁盘上，从而确保了该内存始终驻留在物理内存中。因此，操作系统能够安全地使某个应用程序访问该内存的物理地址，因为这块内存将不会被破坏或者重新定位。

固定内存是一把双刃剑。当使用固定内存时，你将失去虚拟内存的所有功能。特别是，在应用程序中使用每个页锁定内存时都需要分配物理内存，因为这些内存不能交换到磁盘上。这意味着，与使用标准的malloc函数调用相比，系统将更快地耗尽内存。因此，应用程序在物理内存较少的机器上会运行失败，而且意味着应用程序将影响在系统上运行的其它应用程序的性能。

建议，仅对cudaMemcpy()调用中的源内存或者目标内存，才使用页锁定内存，并且在不再需要使用它们时立即释放，而不是等到应用程序关闭时才释放。

关闭

**CUDA**流表示一个**GPU**操作队列，并且该队列中的操作将以指定的顺序执行。

通过使用零拷贝内存，可以避免**CPU**和**GPU**之间的显式复制操作。

对于零拷贝内存，独立GPU和集成GPU，带来的性能提升是不同的。对于集成GPU，使用零拷贝内存通常都会带来性能提升，因为内存在物理上与主机是共享的。将缓冲区声明为零拷贝内存的唯一作用就是避免不必要的数据复制。但是，所有类型的固定内存都存在一定的局限性，零拷贝内存同样不例外。每个固定内存都会占用系统的可用物理内存，这最终将降低系统的性能。对于独立GPU，当输入内存和输出内存都只能使用一次时，那么在独立GPU上使用零拷贝内存将带来性能提升。但由于GPU不会缓存零拷贝内存的内容，如果多次读取内存，那么最终将

得不偿失，还不如一开始就将数据复制到GPU。

　　CUDA工具箱(CUDAToolkit)包含了两个重要的工具库：(1)、**CUFFT**(Fast FourierTransform，快速傅里叶变换)库；(2)、**CUBLAS**(Basic Linear Algebra Subprograms,BLAS)是一个线性代数函数库。

　　**NPP**(NVIDIA Performance Primitives)称为NVIDIA性能原语，它是一个函数库，用来执行基于CUDA加速的数据处理操作，它的基本功能集合主要侧重于图像处理和视频处理。

　　新建一个基于CUDA的测试工程testCUDA，此工程中除了包括common文件外，还添加了另外三个文件，分别为testCUDA.cu、funset.cu、funset.cuh，这三个文件包括了书中绝大部分的测试代码：

testCUDA.cu:

```cpp
01.   #include "funset.cuh"
02.   #include <iostream>
03.   #include "book.h"
04.   #include "cpu_bitmap.h"
05.   #include "gpu_anim.h"
06.
07.   using namespace std;
08.
09.   int test1();//简单的两数相加
10.   int test2();//获取GPU设备相关属性
11.   int test3();//通过线程块索引来计算两个矢量和
12.   int test4();//Julia的CUDA实现
13.   int test5();//通过线程索引来计算两个矢量和
14.   int test6();//通过线程块索引和线程索引来计算两个矢量和
15.   int test7();//ripple的CUDA实现
16.   int test8();//点积运算的CUDA实现
17.   int test9();//Julia的CUDA实现，加入了线程同步函数__syncthreads()
18.   int test10();//光线跟踪(Ray Tracing)实现，没有常量内存+使用事件来计算GPU运行时间
19.   int test11();//光线跟踪(Ray Tracing)实现，使用常量内存+使用事件来计算GPU运行时间
20.   int test12();//模拟热传导，使用纹理内存，有些问题
21.   int test13();//模拟热传导，使用二维纹理内存，有些问题
22.   int test14();//ripple的CUDA+OpenGL实现
23.   int test15();//模拟热传导,CUDA+OpenGL实现，有些问题
24.   int test16();//直方图计算，利用原子操作函数atomicAdd实现
25.   int test17();//固定内存的使用
26.   int test18();//单个stream的使用
27.   int test19();//多个stream的使用
28.   int test20();//通过零拷贝内存的方式实现点积运算
29.   int test21();//使用多个GPU实现点积运算
30.
31.   int main(int argc, char* argv[])
32.   {
33.       test21();
34.
35.       cout<<"ok!"<<endl;
36.       return 0;
37.   }
38.
39.   int test1()
40.   {
41.       int a = 2, b = 3, c = 0;
42.       int* dev_c = NULL;
43.       HANDLE_ERROR(cudaMalloc((void**)&dev_c, sizeof(int)));
44.
45.       //尖括号表示要将一些参数传递给CUDA编译器和运行时系统
46.       //尖括号中这些参数并不是传递给设备代码的参数，而是告诉运行时如何启动设备代码，
47.       //传递给设备代码本身的参数是放在圆括号中传递的，就像标准的函数调用一样
48.       add<<<1, 1>>>(a, b, dev_c);
49.       HANDLE_ERROR(cudaMemcpy(&c, dev_c, sizeof(int), cudaMemcpyDeviceToHost));
50.
51.       printf("%d + %d = %d\n", a, b, c);
52.       cudaFree(dev_c);
53.
54.       return 0;
55.   }
56.
57.   int test2()
58.   {
59.       int count = -1;
60.       HANDLE_ERROR(cudaGetDeviceCount(&count));
61.       printf("device count: %d\n", count);
```

关闭

```cpp
62.
63.        cudaDeviceProp prop;
64.        for (int i = 0; i < count; i++) {
65.            HANDLE_ERROR(cudaGetDeviceProperties(&prop, i));
66.
67.            printf("   --- General Information for device %d ---\n", i);
68.            printf("Name:  %s\n", prop.name);
69.            printf("Compute capability:  %d.%d\n", prop.major, prop.minor);
70.            printf("Clock rate:  %d\n", prop.clockRate);
71.            printf("Device copy overlap:  ");
72.            if (prop.deviceOverlap) printf("Enabled\n");
73.            else printf("Disabled\n");
74.            printf("Kernel execution timeout :  ");
75.            if (prop.kernelExecTimeoutEnabled) printf("Enabled\n");
76.            else printf("Disabled\n");
77.
78.            printf("   --- Memory Information for device %d ---\n", i);
79.            printf("Total global mem:  %ld\n", prop.totalGlobalMem);
80.            printf("Total constant Mem:  %ld\n", prop.totalConstMem);
81.            printf("Max mem pitch:  %ld\n", prop.memPitch);
82.            printf("Texture Alignment:  %ld\n", prop.textureAlignment);
83.
84.            printf("   --- MP Information for device %d ---\n", i);
85.            printf("Multiprocessor count:  %d\n", prop.multiProcessorCount);
86.            printf("Shared mem per mp:  %ld\n", prop.sharedMemPerBlock);
87.            printf("Registers per mp:  %d\n", prop.regsPerBlock);
88.            printf("Threads in warp:  %d\n", prop.warpSize);
89.            printf("Max threads per block:  %d\n", prop.maxThreadsPerBlock);
90.            printf("Max thread dimensions:  (%d, %d, %d)\n", prop.maxThreadsDim[0],
91.                prop.maxThreadsDim[1], prop.maxThreadsDim[2]);
92.            printf("Max grid dimensions:  (%d, %d, %d)\n", prop.maxGridSize[0],
93.                prop.maxGridSize[1], prop.maxGridSize[2]);
94.            printf("\n");
95.        }
96.
97.        int dev;
98.
99.        HANDLE_ERROR(cudaGetDevice(&dev));
100.        printf("ID of current CUDA device:  %d\n", dev);
101.
102.        memset(&prop, 0, sizeof(cudaDeviceProp));
103.        prop.major = 1;
104.        prop.minor = 3;
105.        HANDLE_ERROR(cudaChooseDevice(&dev, &prop));
106.        printf("ID of CUDA device closest to revision %d.%d:  %d\n", prop.major, prop.minor, 
107.
108.        HANDLE_ERROR(cudaSetDevice(dev));
109.
110.        return 0;
111.    }
112.
113.    int test3()
114.    {
115.        int a[NUM] = {0}, b[NUM] = {0}, c[NUM] = {0};
116.        int *dev_a = NULL, *dev_b = NULL, *dev_c = NULL;
117.
118.        //allocate the memory on the GPU
119.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, NUM * sizeof(int)));
120.        HANDLE_ERROR(cudaMalloc((void**)&dev_b, NUM * sizeof(int)));
121.        HANDLE_ERROR(cudaMalloc((void**)&dev_c, NUM * sizeof(int)));
122.
123.        //fill the arrays 'a' and 'b' on the CPU
124.        for (int i=0; i<NUM; i++) {
125.            a[i] = -i;
126.            b[i] = i * i;
127.        }
128.
129.        //copy the arrays 'a' and 'b' to the GPU
130.        HANDLE_ERROR(cudaMemcpy(dev_a, a, NUM * sizeof(int), cudaMemcpyHostToDevice));
131.        HANDLE_ERROR(cudaMemcpy(dev_b, b, NUM * sizeof(int), cudaMemcpyHostToDevice));
132.
133.        //尖括号中的第一个参数表示设备在执行核函数时使用的并行线程块的数量
134.        add_blockIdx<<<NUM,1>>>( dev_a, dev_b, dev_c );
135.
136.        //copy the array 'c' back from the GPU to the CPU
137.        HANDLE_ERROR(cudaMemcpy(c, dev_c, NUM * sizeof(int), cudaMemcpyDeviceToHost));
138.
139.        //display the results
140.        for (int i=0; i<NUM; i++) {
```

关闭

```
141.          printf( "%d + %d = %d\n", a[i], b[i], c[i] );
142.      }
143.
144.      //free the memory allocated on the GPU
145.      HANDLE_ERROR(cudaFree(dev_a));
146.      HANDLE_ERROR(cudaFree(dev_b));
147.      HANDLE_ERROR(cudaFree(dev_c));
148.
149.      return 0;
150.  }
151.
152.  int test4()
153.  {
154.      //globals needed by the update routine
155.      struct DataBlock {
156.          unsigned char* dev_bitmap;
157.      };
158.
159.      DataBlock    data;
160.      CPUBitmap bitmap(DIM, DIM, &data);
161.      unsigned char* dev_bitmap;
162.
163.      HANDLE_ERROR(cudaMalloc((void**)&dev_bitmap, bitmap.image_size()));
164.      data.dev_bitmap = dev_bitmap;
165.
166.      //声明一个二维的线程格
167.      //类型dim3表示一个三维数组，可以用于指定启动线程块的数量
168.      //当用两个值来初始化dim3类型的变量时，CUDA运行时将自动把第3维的大小指定为1
169.      dim3 grid(DIM, DIM);
170.      kernel_julia<<<grid,1>>>(dev_bitmap);
171.
172.      HANDLE_ERROR(cudaMemcpy(bitmap.get_ptr(), dev_bitmap,
173.          bitmap.image_size(), cudaMemcpyDeviceToHost));
174.
175.      HANDLE_ERROR(cudaFree(dev_bitmap));
176.
177.      bitmap.display_and_exit();
178.
179.      return 0;
180.  }
181.
182.  int test5()
183.  {
184.      int a[NUM], b[NUM], c[NUM];
185.      int *dev_a = NULL, *dev_b = NULL, *dev_c = NULL;
186.
187.      //在GPU上分配内存
188.      HANDLE_ERROR(cudaMalloc((void**)&dev_a, NUM * sizeof(int)));
189.      HANDLE_ERROR(cudaMalloc((void**)&dev_b, NUM * sizeof(int)));
190.      HANDLE_ERROR(cudaMalloc((void**)&dev_c, NUM * sizeof(int)));
191.
192.      //在CPU上为数组'a'和'b'赋值
193.      for (int i = 0; i < NUM; i++) {
194.          a[i] = i;
195.          b[i] = i * i;
196.      }
197.
198.      //将数组'a'和'b'复制到GPU
199.      HANDLE_ERROR(cudaMemcpy(dev_a, a, NUM * sizeof(int), cudaMemcpyHostToDevice));
200.      HANDLE_ERROR(cudaMemcpy(dev_b, b, NUM * sizeof(int), cudaMemcpyHostToDevice));
201.
202.      add_threadIdx<<<1, NUM>>>(dev_a, dev_b, dev_c);
203.
204.      //将数组'c'从GPU复制到CPU
205.      HANDLE_ERROR(cudaMemcpy(c, dev_c, NUM * sizeof(int), cudaMemcpyDeviceToHost));
206.
207.      //显示结果
208.      for (int i = 0; i < NUM; i++) {
209.          printf("%d + %d = %d\n", a[i], b[i], c[i]);
210.      }
211.
212.      //释放在GPU分配的内存
213.      cudaFree(dev_a);
214.      cudaFree(dev_b);
215.      cudaFree(dev_c);
216.
217.      return 0;
218.  }
219.
```

关闭

```
220.    int test6()
221.    {
222.        int a[NUM], b[NUM], c[NUM];
223.        int *dev_a = NULL, *dev_b = NULL, *dev_c = NULL;
224.
225.        //在GPU上分配内存
226.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, NUM * sizeof(int)));
227.        HANDLE_ERROR(cudaMalloc((void**)&dev_b, NUM * sizeof(int)));
228.        HANDLE_ERROR(cudaMalloc((void**)&dev_c, NUM * sizeof(int)));
229.
230.        //在CPU上为数组'a'和'b'赋值
231.        for (int i = 0; i < NUM; i++) {
232.            a[i] = i;
233.            b[i] = i * i / 10;
234.        }
235.
236.        //将数组'a'和'b'复制到GPU
237.        HANDLE_ERROR(cudaMemcpy(dev_a, a, NUM * sizeof(int), cudaMemcpyHostToDevice));
238.        HANDLE_ERROR(cudaMemcpy(dev_b, b, NUM * sizeof(int), cudaMemcpyHostToDevice));
239.
240.        add_blockIdx_threadIdx<<<128, 128>>>(dev_a, dev_b, dev_c);
241.
242.        //将数组'c'从GPU复制到CPU
243.        HANDLE_ERROR(cudaMemcpy(c, dev_c, NUM * sizeof(int), cudaMemcpyDeviceToHo
244.
245.        //验证GPU确实完成了我们要求的工作
246.        bool success = true;
247.        for (int i = 0; i < NUM; i++) {
248.            if ((a[i] + b[i]) != c[i]) {
249.                printf("error: %d + %d != %d\n", a[i], b[i], c[i]);
250.                success = false;
251.            }
252.        }
253.
254.        if (success)
255.            printf("we did it!\n");
256.
257.        //释放在GPU分配的内存
258.        cudaFree(dev_a);
259.        cudaFree(dev_b);
260.        cudaFree(dev_c);
261.
262.        return 0;
263.    }
264.
265.    int test7()
266.    {
267.        DataBlock data;
268.        CPUAnimBitmap bitmap(DIM, DIM, &data);
269.        data.bitmap = &bitmap;
270.
271.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_bitmap, bitmap.image_size()));
272.
273.        bitmap.anim_and_exit((void(*)(void*,int))generate_frame, (void(*)(void*))cleanup);
274.
275.        return 0;
276.    }
277.
278.    void generate_frame(DataBlock *d, int ticks)
279.    {
280.        dim3 blocks(DIM/16, DIM/16);
281.        dim3 threads(16, 16);
282.        ripple_kernel<<<blocks,threads>>>(d->dev_bitmap, ticks);
283.
284.        HANDLE_ERROR(cudaMemcpy(d->bitmap->get_ptr(), d->dev_bitmap, d->bitmap->image_size(),
285.    }
286.
287.    //clean up memory allocated on the GPU
288.    void cleanup(DataBlock *d)
289.    {
290.        HANDLE_ERROR(cudaFree(d->dev_bitmap));
291.    }
292.
293.    int test8()
294.    {
295.        float *a, *b, c, *partial_c;
296.        float *dev_a, *dev_b, *dev_partial_c;
297.
298.        //allocate memory on the cpu side
```

关闭

```
299.        a = (float*)malloc(NUM * sizeof(float));
300.        b = (float*)malloc(NUM * sizeof(float));
301.        partial_c = (float*)malloc(blocksPerGrid * sizeof(float));
302.
303.        //allocate the memory on the GPU
304.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, NUM * sizeof(float)));
305.        HANDLE_ERROR(cudaMalloc((void**)&dev_b, NUM * sizeof(float)));
306.        HANDLE_ERROR(cudaMalloc((void**)&dev_partial_c, blocksPerGrid*sizeof(float)));
307.
308.        //fill in the host memory with data
309.        for (int i = 0; i < NUM; i++) {
310.            a[i] = i;
311.            b[i] = i*2;
312.        }
313.
314.        //copy the arrays 'a' and 'b' to the GPU
315.        HANDLE_ERROR(cudaMemcpy(dev_a, a, NUM * sizeof(float), cudaMemcpyHostToDevice));
316.        HANDLE_ERROR(cudaMemcpy(dev_b, b, NUM * sizeof(float), cudaMemcpyHostToDevice));
317.
318.        dot_kernel<<<blocksPerGrid,threadsPerBlock>>>(dev_a, dev_b, dev_partial_c);
319.
320.        //copy the array 'c' back from the GPU to the CPU
321.        HANDLE_ERROR(cudaMemcpy(partial_c, dev_partial_c, blocksPerGrid * sizeof(float), cuda
322.
323.        //finish up on the CPU side
324.        c = 0;
325.        for (int i = 0; i < blocksPerGrid; i++) {
326.            c += partial_c[i];
327.        }
328.
329.        //点积计算结果应该是从0到NUM-1中每个数值的平方再乘以2
330.        //闭合形式解
331.  #define sum_squares(x)  (x * (x + 1) * (2 * x + 1) / 6)
332.        printf("Does GPU value %.6g = %.6g?\n", c, 2 * sum_squares((float)(NUM - 1)));
333.
334.        //free memory on the gpu side
335.        HANDLE_ERROR(cudaFree(dev_a));
336.        HANDLE_ERROR(cudaFree(dev_b));
337.        HANDLE_ERROR(cudaFree(dev_partial_c));
338.
339.        //free memory on the cpu side
340.        free(a);
341.        free(b);
342.        free(partial_c);
343.
344.        return 0;
345.  }
346.
347.  int test9()
348.  {
349.        DataBlock data;
350.        CPUBitmap bitmap(DIM, DIM, &data);
351.        unsigned char *dev_bitmap;
352.
353.        HANDLE_ERROR(cudaMalloc((void**)&dev_bitmap, bitmap.image_size()));
354.        data.dev_bitmap = dev_bitmap;
355.
356.        dim3 grids(DIM / 16, DIM / 16);
357.        dim3 threads(16,16);
358.        julia_kernel<<<grids, threads>>>(dev_bitmap);
359.
360.        HANDLE_ERROR(cudaMemcpy(bitmap.get_ptr(), dev_bitmap, bitmap.image_size(), cudaMemcpyD
361.
362.        HANDLE_ERROR(cudaFree(dev_bitmap));
363.
364.        bitmap.display_and_exit();                                                        关闭
365.
366.        return 0;
367.  }
368.
369.  int test10()
370.  {
371.        DataBlock data;
372.        //capture the start time
373.        cudaEvent_t start, stop;
374.        HANDLE_ERROR(cudaEventCreate(&start));
375.        HANDLE_ERROR(cudaEventCreate(&stop));
376.        HANDLE_ERROR(cudaEventRecord(start, 0));
377.
```

```
378.        CPUBitmap bitmap(DIM, DIM, &data);
379.        unsigned char *dev_bitmap;
380.        Sphere *s;
381.
382.        //allocate memory on the GPU for the output bitmap
383.        HANDLE_ERROR(cudaMalloc((void**)&dev_bitmap, bitmap.image_size()));
384.        //allocate memory for the Sphere dataset
385.        HANDLE_ERROR(cudaMalloc((void**)&s, sizeof(Sphere) * SPHERES));
386.
387.        //allocate temp memory, initialize it, copy to memory on the GPU, then free our temp r
388.        Sphere *temp_s = (Sphere*)malloc(sizeof(Sphere) * SPHERES);
389.        for (int i = 0; i < SPHERES; i++) {
390.            temp_s[i].r = rnd(1.0f);
391.            temp_s[i].g = rnd(1.0f);
392.            temp_s[i].b = rnd(1.0f);
393.            temp_s[i].x = rnd(1000.0f) - 500;
394.            temp_s[i].y = rnd(1000.0f) - 500;
395.            temp_s[i].z = rnd(1000.0f) - 500;
396.            temp_s[i].radius = rnd(100.0f) + 20;
397.        }
398.
399.        HANDLE_ERROR(cudaMemcpy( s, temp_s, sizeof(Sphere) * SPHERES, cudaMemcpyHostToDevice)
400.        free(temp_s);
401.
402.        //generate a bitmap from our sphere data
403.        dim3 grids(DIM / 16, DIM / 16);
404.        dim3 threads(16, 16);
405.        RayTracing_kernel<<<grids, threads>>>(s, dev_bitmap);
406.
407.        //copy our bitmap back from the GPU for display
408.        HANDLE_ERROR(cudaMemcpy(bitmap.get_ptr(), dev_bitmap, bitmap.image_size(), cudaMemcpyI
409.
410.        //get stop time, and display the timing results
411.        HANDLE_ERROR(cudaEventRecord(stop, 0));
412.        HANDLE_ERROR(cudaEventSynchronize(stop));
413.        float elapsedTime;
414.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
415.        printf("Time to generate:  %3.1f ms\n", elapsedTime);
416.
417.        HANDLE_ERROR(cudaEventDestroy(start));
418.        HANDLE_ERROR(cudaEventDestroy(stop));
419.
420.        HANDLE_ERROR(cudaFree(dev_bitmap));
421.        HANDLE_ERROR(cudaFree(s));
422.
423.        // display
424.        bitmap.display_and_exit();
425.
426.        return 0;
427.    }
428.
429.    int test11()
430.    {
431.        DataBlock data;
432.        //capture the start time
433.        cudaEvent_t start, stop;
434.        HANDLE_ERROR(cudaEventCreate(&start));
435.        HANDLE_ERROR(cudaEventCreate(&stop));
436.        HANDLE_ERROR(cudaEventRecord(start, 0));
437.
438.        CPUBitmap bitmap(DIM, DIM, &data);
439.        unsigned char *dev_bitmap;
440.
441.        //allocate memory on the GPU for the output bitmap
442.        HANDLE_ERROR(cudaMalloc((void**)&dev_bitmap, bitmap.image_size()));
443.
444.        //allocate temp memory, initialize it, copy to constant memory on the GPU, then free t
445.        Sphere *temp_s = (Sphere*)malloc(sizeof(Sphere) * SPHERES);
446.        for (int i = 0; i < SPHERES; i++) {
447.            temp_s[i].r = rnd(1.0f);
448.            temp_s[i].g = rnd(1.0f);
449.            temp_s[i].b = rnd(1.0f);
450.            temp_s[i].x = rnd(1000.0f) - 500;
451.            temp_s[i].y = rnd(1000.0f) - 500;
452.            temp_s[i].z = rnd(1000.0f) - 500;
453.            temp_s[i].radius = rnd(100.0f) + 20;
454.        }
455.
456.        HANDLE_ERROR(cudaMemcpyToSymbol(s, temp_s, sizeof(Sphere) * SPHERES));
```

关闭

```
457.        free(temp_s);
458.
459.        //generate a bitmap from our sphere data
460.        dim3 grids(DIM / 16, DIM / 16);
461.        dim3 threads(16, 16);
462.        RayTracing_kernel<<<grids, threads>>>(dev_bitmap);
463.
464.        //copy our bitmap back from the GPU for display
465.        HANDLE_ERROR(cudaMemcpy(bitmap.get_ptr(), dev_bitmap, bitmap.image_size(), cudaMemcpyD
466.
467.        //get stop time, and display the timing results
468.        HANDLE_ERROR(cudaEventRecord(stop, 0));
469.        HANDLE_ERROR(cudaEventSynchronize(stop));
470.        float elapsedTime;
471.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
472.        printf("Time to generate:  %3.1f ms\n", elapsedTime);
473.
474.        HANDLE_ERROR(cudaEventDestroy(start));
475.        HANDLE_ERROR(cudaEventDestroy(stop));
476.
477.        HANDLE_ERROR(cudaFree(dev_bitmap));
478.
479.        //display
480.        bitmap.display_and_exit();
481.
482.        return 0;
483.    }
484.
485.    int test12()
486.    {
487.        Heat_DataBlock data;
488.        CPUAnimBitmap bitmap(DIM, DIM, &data);
489.        data.bitmap = &bitmap;
490.        data.totalTime = 0;
491.        data.frames = 0;
492.
493.        HANDLE_ERROR(cudaEventCreate(&data.start));
494.        HANDLE_ERROR(cudaEventCreate(&data.stop));
495.
496.        int imageSize = bitmap.image_size();
497.
498.        HANDLE_ERROR(cudaMalloc((void**)&data.output_bitmap, imageSize));
499.
500.        //assume float == 4 chars in size (ie rgba)
501.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_inSrc, imageSize));
502.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_outSrc, imageSize));
503.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_constSrc, imageSize));
504.
505.        HANDLE_ERROR(cudaBindTexture(NULL, texConstSrc, data.dev_constSrc, imageSize));
506.        HANDLE_ERROR(cudaBindTexture(NULL, texIn, data.dev_inSrc, imageSize));
507.        HANDLE_ERROR(cudaBindTexture(NULL, texOut, data.dev_outSrc, imageSize));
508.
509.        //intialize the constant data
510.        float *temp = (float*)malloc(imageSize);
511.
512.        for (int i = 0; i < DIM*DIM; i++) {
513.            temp[i] = 0;
514.            int x = i % DIM;
515.            int y = i / DIM;
516.            if ((x>300) && (x<600) && (y>310) && (y<601))
517.                temp[i] = MAX_TEMP;
518.        }
519.
520.        temp[DIM * 100 + 100] = (MAX_TEMP + MIN_TEMP) / 2;
521.        temp[DIM * 700 + 100] = MIN_TEMP;
522.        temp[DIM * 300 + 300] = MIN_TEMP;
523.        temp[DIM * 200 + 700] = MIN_TEMP;
524.
525.        for (int y = 800; y < 900; y++) {
526.            for (int x = 400; x < 500; x++) {
527.                temp[x + y * DIM] = MIN_TEMP;
528.            }
529.        }
530.
531.        HANDLE_ERROR(cudaMemcpy(data.dev_constSrc, temp, imageSize, cudaMemcpyHostToDevice));
532.
533.        //initialize the input data
534.        for (int y = 800; y < DIM; y++) {
535.            for (int x = 0; x < 200; x++) {
```

关闭

```
536.              temp[x+y*DIM] = MAX_TEMP;
537.          }
538.      }
539.
540.      HANDLE_ERROR(cudaMemcpy(data.dev_inSrc, temp,imageSize, cudaMemcpyHostToDevice));
541.      free(temp);
542.
543.      bitmap.anim_and_exit((void (*)(void*,int))Heat_anim_gpu, (void (*)(void*))Heat_anim_ex
544.
545.      return 0;
546.  }
547.
548.  int test13()
549.  {
550.      Heat_DataBlock data;
551.      CPUAnimBitmap bitmap(DIM, DIM, &data);
552.      data.bitmap = &bitmap;
553.      data.totalTime = 0;
554.      data.frames = 0;
555.      HANDLE_ERROR(cudaEventCreate(&data.start));
556.      HANDLE_ERROR(cudaEventCreate(&data.stop));
557.
558.      int imageSize = bitmap.image_size();
559.
560.      HANDLE_ERROR(cudaMalloc((void**)&data.output_bitmap, imageSize));
561.
562.      //assume float == 4 chars in size (ie rgba)
563.      HANDLE_ERROR(cudaMalloc((void**)&data.dev_inSrc, imageSize));
564.      HANDLE_ERROR(cudaMalloc((void**)&data.dev_outSrc, imageSize));
565.      HANDLE_ERROR(cudaMalloc((void**)&data.dev_constSrc, imageSize));
566.
567.      cudaChannelFormatDesc desc = cudaCreateChannelDesc<float>();
568.      HANDLE_ERROR(cudaBindTexture2D(NULL, texConstSrc2, data.dev_constSrc, desc, DIM, DIM,
569.      HANDLE_ERROR(cudaBindTexture2D(NULL, texIn2, data.dev_inSrc, desc, DIM, DIM, sizeof(f:
570.      HANDLE_ERROR(cudaBindTexture2D(NULL, texOut2, data.dev_outSrc, desc, DIM, DIM, sizeof
571.
572.      //initialize the constant data
573.      float *temp = (float*)malloc(imageSize);
574.      for (int i = 0; i < DIM*DIM; i++) {
575.          temp[i] = 0;
576.          int x = i % DIM;
577.          int y = i / DIM;
578.          if ((x > 300) && ( x < 600) && (y > 310) && (y < 601))
579.              temp[i] = MAX_TEMP;
580.      }
581.
582.      temp[DIM * 100 + 100] = (MAX_TEMP + MIN_TEMP) / 2;
583.      temp[DIM * 700 + 100] = MIN_TEMP;
584.      temp[DIM * 300 + 300] = MIN_TEMP;
585.      temp[DIM * 200 + 700] = MIN_TEMP;
586.
587.      for (int y = 800; y < 900; y++) {
588.          for (int x = 400; x < 500; x++) {
589.              temp[x + y * DIM] = MIN_TEMP;
590.          }
591.      }
592.
593.      HANDLE_ERROR(cudaMemcpy(data.dev_constSrc, temp, imageSize, cudaMemcpyHostToDevice));
594.
595.      //initialize the input data
596.      for (int y = 800; y < DIM; y++) {
597.          for (int x = 0; x < 200; x++) {
598.              temp[x + y * DIM] = MAX_TEMP;
599.          }
600.      }
601.
602.      HANDLE_ERROR(cudaMemcpy(data.dev_inSrc, temp,imageSize, cudaMemcpyHostToDevice));
603.      free(temp);
604.
605.      bitmap.anim_and_exit((void (*)(void*,int))anim_gpu, (void (*)(void*))anim_exit);
606.
607.      return 0;
608.  }
609.
610.  void Heat_anim_gpu(Heat_DataBlock *d, int ticks)
611.  {
612.      HANDLE_ERROR(cudaEventRecord(d->start, 0));
613.
614.      dim3 blocks(DIM / 16, DIM / 16);
```

关闭

```
615.        dim3 threads(16, 16);
616.        CPUAnimBitmap *bitmap = d->bitmap;
617.
618.        //since tex is global and bound, we have to use a flag to
619.        //select which is in/out per iteration
620.        volatile bool dstOut = true;
621.
622.        for (int i = 0; i < 90; i++) {
623.            float *in, *out;
624.            if (dstOut) {
625.                in  = d->dev_inSrc;
626.                out = d->dev_outSrc;
627.            } else {
628.                out = d->dev_inSrc;
629.                in  = d->dev_outSrc;
630.            }
631.
632.            Heat_copy_const_kernel<<<blocks, threads>>>(in);
633.            Heat_blend_kernel<<<blocks, threads>>>(out, dstOut);
634.            dstOut = !dstOut;
635.        }
636.
637.        float_to_color<<<blocks, threads>>>(d->output_bitmap, d->dev_inSrc);
638.
639.        HANDLE_ERROR(cudaMemcpy(bitmap->get_ptr(), d->output_bitmap, bitmap->image_size(), cud
640.
641.        HANDLE_ERROR(cudaEventRecord(d->stop, 0));
642.        HANDLE_ERROR(cudaEventSynchronize(d->stop));
643.        float elapsedTime;
644.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, d->start, d->stop));
645.        d->totalTime += elapsedTime;
646.        ++d->frames;
647.
648.        printf( "Average Time per frame:  %3.1f ms\n", d->totalTime/d->frames );
649.    }
650.
651.    void anim_gpu(Heat_DataBlock *d, int ticks)
652.    {
653.        HANDLE_ERROR(cudaEventRecord(d->start, 0));
654.        dim3 blocks(DIM / 16, DIM / 16);
655.        dim3 threads(16, 16);
656.        CPUAnimBitmap  *bitmap = d->bitmap;
657.
658.        //since tex is global and bound, we have to use a flag to
659.        //select which is in/out per iteration
660.        volatile bool dstOut = true;
661.        for (int i = 0; i < 90; i++) {
662.            float *in, *out;
663.            if (dstOut) {
664.                in  = d->dev_inSrc;
665.                out = d->dev_outSrc;
666.            } else {
667.                out = d->dev_inSrc;
668.                in  = d->dev_outSrc;
669.            }
670.            copy_const_kernel<<<blocks, threads>>>(in);
671.            blend_kernel<<<blocks, threads>>>(out, dstOut);
672.            dstOut = !dstOut;
673.        }
674.
675.        float_to_color<<<blocks, threads>>>(d->output_bitmap, d->dev_inSrc);
676.
677.        HANDLE_ERROR(cudaMemcpy(bitmap->get_ptr(), d->output_bitmap, bitmap->image_size(), cud
678.
679.        HANDLE_ERROR(cudaEventRecord(d->stop, 0));
680.        HANDLE_ERROR(cudaEventSynchronize(d->stop));
681.        float elapsedTime;
682.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, d->start, d->stop));
683.        d->totalTime += elapsedTime;
684.        ++d->frames;
685.        printf("Average Time per frame:  %3.1f ms\n", d->totalTime/d->frames);
686.    }
687.
688.    void Heat_anim_exit(Heat_DataBlock *d)
689.    {
690.        cudaUnbindTexture(texIn);
691.        cudaUnbindTexture(texOut);
692.        cudaUnbindTexture(texConstSrc);
693.
```

关闭

```
694.        HANDLE_ERROR(cudaFree(d->dev_inSrc));
695.        HANDLE_ERROR(cudaFree(d->dev_outSrc));
696.        HANDLE_ERROR(cudaFree(d->dev_constSrc));
697.
698.        HANDLE_ERROR(cudaEventDestroy(d->start));
699.        HANDLE_ERROR(cudaEventDestroy(d->stop));
700.    }
701.
702.    //clean up memory allocated on the GPU
703.    void anim_exit(Heat_DataBlock *d)
704.    {
705.        cudaUnbindTexture(texIn2);
706.        cudaUnbindTexture(texOut2);
707.        cudaUnbindTexture(texConstSrc2);
708.        HANDLE_ERROR(cudaFree(d->dev_inSrc));
709.        HANDLE_ERROR(cudaFree(d->dev_outSrc));
710.        HANDLE_ERROR(cudaFree(d->dev_constSrc));
711.
712.        HANDLE_ERROR(cudaEventDestroy(d->start));
713.        HANDLE_ERROR(cudaEventDestroy(d->stop));
714.    }
715.
716.    int test14()
717.    {
718.        GPUAnimBitmap  bitmap(DIM, DIM, NULL);
719.
720.        bitmap.anim_and_exit((void (*)(uchar4*, void*, int))generate_frame_opengl, NULL);
721.
722.        return 0;
723.    }
724.
725.    int test15()
726.    {
727.        DataBlock_opengl data;
728.        GPUAnimBitmap bitmap(DIM, DIM, &data);
729.        data.totalTime = 0;
730.        data.frames = 0;
731.        HANDLE_ERROR(cudaEventCreate(&data.start));
732.        HANDLE_ERROR(cudaEventCreate(&data.stop));
733.
734.        int imageSize = bitmap.image_size();
735.
736.        //assume float == 4 chars in size (ie rgba)
737.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_inSrc, imageSize));
738.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_outSrc, imageSize));
739.        HANDLE_ERROR(cudaMalloc((void**)&data.dev_constSrc, imageSize));
740.
741.        HANDLE_ERROR(cudaBindTexture(NULL, texConstSrc ,data.dev_constSrc, imageSize));
742.        HANDLE_ERROR(cudaBindTexture(NULL, texIn, data.dev_inSrc, imageSize));
743.        HANDLE_ERROR(cudaBindTexture(NULL, texOut, data.dev_outSrc, imageSize));
744.
745.        //intialize the constant data
746.        float *temp = (float*)malloc(imageSize);
747.        for (int i = 0; i < DIM*DIM; i++) {
748.            temp[i] = 0;
749.            int x = i % DIM;
750.            int y = i / DIM;
751.            if ((x>300) && (x<600) && (y>310) && (y<601))
752.                temp[i] = MAX_TEMP;
753.        }
754.
755.        temp[DIM*100+100] = (MAX_TEMP + MIN_TEMP)/2;
756.        temp[DIM*700+100] = MIN_TEMP;
757.        temp[DIM*300+300] = MIN_TEMP;
758.        temp[DIM*200+700] = MIN_TEMP;
759.
760.        for (int y = 800; y < 900; y++) {
761.            for (int x = 400; x < 500; x++) {
762.                temp[x+y*DIM] = MIN_TEMP;
763.            }
764.        }
765.
766.        HANDLE_ERROR(cudaMemcpy(data.dev_constSrc, temp, imageSize, cudaMemcpyHostToDevice));
767.
768.        //initialize the input data
769.        for (int y = 800; y < DIM; y++) {
770.            for (int x = 0; x < 200; x++) {
771.                temp[x+y*DIM] = MAX_TEMP;
772.            }
```

关闭

```
773.            }
774.
775.            HANDLE_ERROR(cudaMemcpy(data.dev_inSrc, temp, imageSize, cudaMemcpyHostToDevice));
776.            free(temp);
777.
778.            bitmap.anim_and_exit((void (*)(uchar4*, void*, int))anim_gpu_opengl, (void (*)(void*)
779.
780.            return 0;
781.    }
782.
783.    void anim_gpu_opengl(uchar4* outputBitmap, DataBlock_opengl *d, int ticks)
784.    {
785.            HANDLE_ERROR(cudaEventRecord(d->start, 0));
786.            dim3 blocks(DIM / 16, DIM / 16);
787.            dim3 threads(16, 16);
788.
789.            //since tex is global and bound, we have to use a flag to select which is in/out per :
790.            volatile bool dstOut = true;
791.            for (int i = 0; i < 90; i++) {
792.                float *in, *out;
793.                if (dstOut) {
794.                    in  = d->dev_inSrc;
795.                    out = d->dev_outSrc;
796.                } else {
797.                    out = d->dev_inSrc;
798.                    in  = d->dev_outSrc;
799.                }
800.                Heat_copy_const_kernel_opengl<<<blocks, threads>>>(in);
801.                Heat_blend_kernel_opengl<<<blocks, threads>>>(out, dstOut);
802.                dstOut = !dstOut;
803.            }
804.
805.            float_to_color<<<blocks, threads>>>(outputBitmap, d->dev_inSrc);
806.
807.            HANDLE_ERROR(cudaEventRecord(d->stop, 0));
808.            HANDLE_ERROR(cudaEventSynchronize(d->stop));
809.            float elapsedTime;
810.            HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, d->start, d->stop));
811.            d->totalTime += elapsedTime;
812.            ++d->frames;
813.            printf("Average Time per frame:  %3.1f ms\n", d->totalTime/d->frames);
814.    }
815.
816.    void anim_exit_opengl(DataBlock_opengl *d)
817.    {
818.            HANDLE_ERROR(cudaUnbindTexture(texIn));
819.            HANDLE_ERROR(cudaUnbindTexture(texOut));
820.            HANDLE_ERROR(cudaUnbindTexture(texConstSrc));
821.            HANDLE_ERROR(cudaFree(d->dev_inSrc));
822.            HANDLE_ERROR(cudaFree(d->dev_outSrc));
823.            HANDLE_ERROR(cudaFree(d->dev_constSrc));
824.
825.            HANDLE_ERROR(cudaEventDestroy(d->start));
826.            HANDLE_ERROR(cudaEventDestroy(d->stop));
827.    }
828.
829.    int test16()
830.    {
831.            unsigned char *buffer = (unsigned char*)big_random_block(SIZE);
832.
833.            //capture the start time starting the timer here so that we include the cost of
834.            //all of the operations on the GPU.  if the data were already on the GPU and we just
835.            //timed the kernel the timing would drop from 74 ms to 15 ms.  Very fast.
836.            cudaEvent_t start, stop;
837.            HANDLE_ERROR( cudaEventCreate( &start ) );
838.            HANDLE_ERROR( cudaEventCreate( &stop ) );
839.            HANDLE_ERROR( cudaEventRecord( start, 0 ) );
840.
841.            // allocate memory on the GPU for the file's data
842.            unsigned char *dev_buffer;
843.            unsigned int *dev_histo;
844.            HANDLE_ERROR(cudaMalloc((void**)&dev_buffer, SIZE));
845.            HANDLE_ERROR(cudaMemcpy(dev_buffer, buffer, SIZE, cudaMemcpyHostToDevice));
846.
847.            HANDLE_ERROR(cudaMalloc((void**)&dev_histo, 256 * sizeof(int)));
848.            HANDLE_ERROR(cudaMemset(dev_histo, 0, 256 * sizeof(int)));
849.
850.            //kernel launch - 2x the number of mps gave best timing
851.            cudaDeviceProp prop;
```

关闭

```
852.        HANDLE_ERROR(cudaGetDeviceProperties(&prop, 0));
853.        int blocks = prop.multiProcessorCount;
854.        histo_kernel<<<blocks*2, 256>>>(dev_buffer, SIZE, dev_histo);
855.
856.        unsigned int histo[256];
857.        HANDLE_ERROR(cudaMemcpy(histo, dev_histo, 256 * sizeof(int), cudaMemcpyDeviceToHost));
858.
859.        //get stop time, and display the timing results
860.        HANDLE_ERROR(cudaEventRecord(stop, 0));
861.        HANDLE_ERROR(cudaEventSynchronize(stop));
862.        float elapsedTime;
863.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
864.        printf("Time to generate:  %3.1f ms\n", elapsedTime);
865.
866.        long histoCount = 0;
867.        for (int i=0; i<256; i++) {
868.            histoCount += histo[i];
869.        }
870.        printf("Histogram Sum:  %ld\n", histoCount);
871.
872.        //verify that we have the same counts via CPU
873.        for (int i = 0; i < SIZE; i++)
874.            histo[buffer[i]]--;
875.        for (int i = 0; i < 256; i++) {
876.            if (histo[i] != 0)
877.                printf("Failure at %d!\n", i);
878.        }
879.
880.        HANDLE_ERROR(cudaEventDestroy(start));
881.        HANDLE_ERROR(cudaEventDestroy(stop));
882.        cudaFree(dev_histo);
883.        cudaFree(dev_buffer);
884.        free(buffer);
885.
886.        return 0;
887.    }
888.
889.    float cuda_malloc_test(int size, bool up)
890.    {
891.        cudaEvent_t start, stop;
892.        int *a, *dev_a;
893.        float elapsedTime;
894.
895.        HANDLE_ERROR(cudaEventCreate(&start));
896.        HANDLE_ERROR(cudaEventCreate(&stop));
897.
898.        a = (int*)malloc(size * sizeof(*a));
899.        HANDLE_NULL(a);
900.        HANDLE_ERROR(cudaMalloc((void**)&dev_a,size * sizeof(*dev_a)));
901.
902.        HANDLE_ERROR(cudaEventRecord(start, 0));
903.
904.        for (int i=0; i<100; i++) {
905.            if (up)
906.                HANDLE_ERROR(cudaMemcpy(dev_a, a, size * sizeof( *dev_a ), cudaMemcpyHostToDev
907.            else
908.                HANDLE_ERROR(cudaMemcpy(a, dev_a, size * sizeof(*dev_a), cudaMemcpyDeviceToHos
909.        }
910.        HANDLE_ERROR(cudaEventRecord(stop, 0));
911.        HANDLE_ERROR(cudaEventSynchronize(stop));
912.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
913.
914.        free(a);
915.        HANDLE_ERROR(cudaFree(dev_a));
916.        HANDLE_ERROR(cudaEventDestroy(start));
917.        HANDLE_ERROR(cudaEventDestroy(stop));                                    关闭
918.
919.        return elapsedTime;
920.    }
921.
922.    float cuda_host_alloc_test(int size, bool up)
923.    {
924.        cudaEvent_t start, stop;
925.        int *a, *dev_a;
926.        float elapsedTime;
927.
928.        HANDLE_ERROR(cudaEventCreate(&start));
929.        HANDLE_ERROR(cudaEventCreate(&stop));
930.
```

```
931.        HANDLE_ERROR(cudaHostAlloc((void**)&a,size * sizeof(*a), cudaHostAllocDefault));
932.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, size * sizeof(*dev_a)));
933.
934.        HANDLE_ERROR(cudaEventRecord(start, 0));
935.
936.        for (int i=0; i<100; i++) {
937.            if (up)
938.                HANDLE_ERROR(cudaMemcpy(dev_a, a,size * sizeof(*a), cudaMemcpyHostToDevice));
939.            else
940.                HANDLE_ERROR(cudaMemcpy(a, dev_a,size * sizeof(*a), cudaMemcpyDeviceToHost));
941.        }
942.        HANDLE_ERROR(cudaEventRecord(stop, 0));
943.        HANDLE_ERROR(cudaEventSynchronize(stop));
944.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
945.
946.        HANDLE_ERROR(cudaFreeHost(a));
947.        HANDLE_ERROR(cudaFree(dev_a));
948.        HANDLE_ERROR(cudaEventDestroy(start));
949.        HANDLE_ERROR(cudaEventDestroy(stop));
950.
951.        return elapsedTime;
952.    }
953.
954.    int test17()
955.    {
956.        float elapsedTime;
957.        float MB = (float)100 * SIZE * sizeof(int) / 1024 / 1024;
958.
959.        //try it with cudaMalloc
960.        elapsedTime = cuda_malloc_test(SIZE, true);
961.        printf("Time using cudaMalloc:  %3.1f ms\n", elapsedTime);
962.        printf("\tMB/s during copy up:  %3.1f\n", MB/(elapsedTime/1000));
963.
964.        elapsedTime = cuda_malloc_test(SIZE, false);
965.        printf("Time using cudaMalloc:  %3.1f ms\n", elapsedTime);
966.        printf("\tMB/s during copy down:  %3.1f\n", MB/(elapsedTime/1000));
967.
968.        //now try it with cudaHostAlloc
969.        elapsedTime = cuda_host_alloc_test(SIZE, true);
970.        printf("Time using cudaHostAlloc:  %3.1f ms\n", elapsedTime);
971.        printf("\tMB/s during copy up:  %3.1f\n", MB/(elapsedTime/1000));
972.
973.        elapsedTime = cuda_host_alloc_test(SIZE, false);
974.        printf("Time using cudaHostAlloc:  %3.1f ms\n", elapsedTime);
975.        printf("\tMB/s during copy down:  %3.1f\n", MB/(elapsedTime/1000));
976.
977.        return 0;
978.    }
979.
980.    int test18()
981.    {
982.        cudaDeviceProp prop;
983.        int whichDevice;
984.        HANDLE_ERROR(cudaGetDevice(&whichDevice));
985.        HANDLE_ERROR(cudaGetDeviceProperties(&prop, whichDevice));
986.
987.        if (!prop.deviceOverlap) {
988.            printf("Device will not handle overlaps, so no speed up from streams\n");
989.            return 0;
990.        }
991.
992.        cudaEvent_t start, stop;
993.        float elapsedTime;
994.        cudaStream_t stream;
995.        int *host_a, *host_b, *host_c;
996.        int *dev_a, *dev_b, *dev_c;                                                          关闭
997.
998.        //start the timers
999.        HANDLE_ERROR(cudaEventCreate(&start));
1000.       HANDLE_ERROR(cudaEventCreate(&stop));
1001.
1002.       //initialize the stream
1003.       HANDLE_ERROR(cudaStreamCreate(&stream));
1004.
1005.       //allocate the memory on the GPU
1006.       HANDLE_ERROR(cudaMalloc((void**)&dev_a, NUM * sizeof(int)));
1007.       HANDLE_ERROR(cudaMalloc((void**)&dev_b, NUM * sizeof(int)));
1008.       HANDLE_ERROR(cudaMalloc((void**)&dev_c, NUM * sizeof(int)));
1009.
```

```
1010.        //allocate host locked memory, used to stream
1011.        HANDLE_ERROR(cudaHostAlloc((void**)&host_a, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
1012.        HANDLE_ERROR(cudaHostAlloc((void**)&host_b, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
1013.        HANDLE_ERROR(cudaHostAlloc((void**)&host_c, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
1014.
1015.        for (int i=0; i<FULL_DATA_SIZE; i++) {
1016.            host_a[i] = rand();
1017.            host_b[i] = rand();
1018.        }
1019.
1020.        HANDLE_ERROR(cudaEventRecord(start, 0));
1021.        //now loop over full data, in bite-sized chunks
1022.        for (int i=0; i<FULL_DATA_SIZE; i+= NUM) {
1023.            //copy the locked memory to the device, async
1024.            HANDLE_ERROR(cudaMemcpyAsync(dev_a, host_a+i, NUM * sizeof(int), cudaMemcpyHostToI
1025.            HANDLE_ERROR(cudaMemcpyAsync(dev_b, host_b+i, NUM * sizeof(int), cudaMemcpyHostToI
1026.
1027.            singlestream_kernel<<<NUM/256, 256, 0, stream>>>(dev_a, dev_b, dev_c);
1028.
1029.            //copy the data from device to locked memory
1030.            HANDLE_ERROR(cudaMemcpyAsync(host_c+i, dev_c, NUM * sizeof(int), cudaMemcpyDevice
1031.
1032.        }
1033.
1034.        // copy result chunk from locked to full buffer
1035.        HANDLE_ERROR(cudaStreamSynchronize(stream));
1036.
1037.        HANDLE_ERROR(cudaEventRecord(stop, 0));
1038.        HANDLE_ERROR(cudaEventSynchronize(stop));
1039.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime, start, stop));
1040.        printf("Time taken:  %3.1f ms\n", elapsedTime);
1041.
1042.        //cleanup the streams and memory
1043.        HANDLE_ERROR(cudaFreeHost(host_a));
1044.        HANDLE_ERROR(cudaFreeHost(host_b));
1045.        HANDLE_ERROR(cudaFreeHost(host_c));
1046.        HANDLE_ERROR(cudaFree(dev_a));
1047.        HANDLE_ERROR(cudaFree(dev_b));
1048.        HANDLE_ERROR(cudaFree(dev_c));
1049.        HANDLE_ERROR(cudaStreamDestroy(stream));
1050.
1051.        return 0;
1052.    }
1053.
1054.    int test19()
1055.    {
1056.        cudaDeviceProp prop;
1057.        int whichDevice;
1058.        HANDLE_ERROR(cudaGetDevice(&whichDevice));
1059.        HANDLE_ERROR(cudaGetDeviceProperties(&prop, whichDevice));
1060.        if (!prop.deviceOverlap) {
1061.            printf( "Device will not handle overlaps, so no speed up from streams\n" );
1062.            return 0;
1063.        }
1064.
1065.        //start the timers
1066.        cudaEvent_t start, stop;
1067.        HANDLE_ERROR(cudaEventCreate(&start));
1068.        HANDLE_ERROR(cudaEventCreate(&stop));
1069.
1070.        //initialize the streams
1071.        cudaStream_t stream0, stream1;
1072.        HANDLE_ERROR(cudaStreamCreate(&stream0));
1073.        HANDLE_ERROR(cudaStreamCreate(&stream1));
1074.
1075.        int *host_a, *host_b, *host_c;                                            关闭
1076.        int *dev_a0, *dev_b0, *dev_c0;//为第0个流分配的GPU内存
1077.        int *dev_a1, *dev_b1, *dev_c1;//为第1个流分配的GPU内存
1078.
1079.        //allocate the memory on the GPU
1080.        HANDLE_ERROR(cudaMalloc((void**)&dev_a0, NUM * sizeof(int)));
1081.        HANDLE_ERROR(cudaMalloc((void**)&dev_b0, NUM * sizeof(int)));
1082.        HANDLE_ERROR(cudaMalloc((void**)&dev_c0, NUM * sizeof(int)));
1083.        HANDLE_ERROR(cudaMalloc((void**)&dev_a1, NUM * sizeof(int)));
1084.        HANDLE_ERROR(cudaMalloc((void**)&dev_b1, NUM * sizeof(int)));
1085.        HANDLE_ERROR(cudaMalloc((void**)&dev_c1, NUM * sizeof(int)));
1086.
1087.        //allocate host locked memory, used to stream
1088.        HANDLE_ERROR(cudaHostAlloc((void**)&host_a, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
```

```
1089.        HANDLE_ERROR(cudaHostAlloc((void**)&host_b, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
1090.        HANDLE_ERROR(cudaHostAlloc((void**)&host_c, FULL_DATA_SIZE * sizeof(int), cudaHostAll(
1091.
1092.        for (int i=0; i<FULL_DATA_SIZE; i++) {
1093.            host_a[i] = rand();
1094.            host_b[i] = rand();
1095.        }
1096.
1097.        HANDLE_ERROR(cudaEventRecord(start, 0));
1098.
1099.        //now loop over full data, in bite-sized chunks
1100.        for (int i=0; i<FULL_DATA_SIZE; i+= NUM*2) {
1101.            //enqueue copies of a in stream0 and stream1
1102.            //将锁定内存以异步方式复制到设备上
1103.            HANDLE_ERROR(cudaMemcpyAsync(dev_a0, host_a+i, NUM * sizeof(int), cudaMemcpyHostT(
1104.            HANDLE_ERROR(cudaMemcpyAsync(dev_a1, host_a+i+NUM, NUM * sizeof(int), cudaMemcpyH(
1105.            //enqueue copies of b in stream0 and stream1
1106.            HANDLE_ERROR(cudaMemcpyAsync(dev_b0, host_b+i, NUM * sizeof(int), cudaMemcpyHostT(
1107.            HANDLE_ERROR(cudaMemcpyAsync(dev_b1, host_b+i+NUM, NUM * sizeof(int), cudaMemcpyH(
1108.
1109.            //enqueue kernels in stream0 and stream1
1110.            singlestream_kernel<<<NUM/256, 256, 0, stream0>>>(dev_a0, dev_b0, dev_c0);
1111.            singlestream_kernel<<<NUM/256, 256, 0, stream1>>>(dev_a1, dev_b1, dev_c1);
1112.
1113.            //enqueue copies of c from device to locked memory
1114.            HANDLE_ERROR(cudaMemcpyAsync(host_c+i, dev_c0, NUM * sizeof(int), cudaMemcpyDevice
1115.            HANDLE_ERROR(cudaMemcpyAsync(host_c+i+NUM, dev_c1, NUM * sizeof(int), cudaMemcpyD(
1116.        }
1117.
1118.        float elapsedTime;
1119.
1120.        HANDLE_ERROR(cudaStreamSynchronize(stream0));
1121.        HANDLE_ERROR(cudaStreamSynchronize(stream1));
1122.
1123.        HANDLE_ERROR(cudaEventRecord(stop, 0));
1124.
1125.        HANDLE_ERROR(cudaEventSynchronize(stop));
1126.        HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime,start, stop));
1127.        printf( "Time taken:  %3.1f ms\n", elapsedTime );
1128.
1129.        //cleanup the streams and memory
1130.        HANDLE_ERROR(cudaFreeHost(host_a));
1131.        HANDLE_ERROR(cudaFreeHost(host_b));
1132.        HANDLE_ERROR(cudaFreeHost(host_c));
1133.        HANDLE_ERROR(cudaFree(dev_a0));
1134.        HANDLE_ERROR(cudaFree(dev_b0));
1135.        HANDLE_ERROR(cudaFree(dev_c0));
1136.        HANDLE_ERROR(cudaFree(dev_a1));
1137.        HANDLE_ERROR(cudaFree(dev_b1));
1138.        HANDLE_ERROR(cudaFree(dev_c1));
1139.        HANDLE_ERROR(cudaStreamDestroy(stream0));
1140.        HANDLE_ERROR(cudaStreamDestroy(stream1));
1141.
1142.        return 0;
1143.    }
1144.
1145.    float malloc_test(int size)
1146.    {
1147.        cudaEvent_t start, stop;
1148.        float *a, *b, c, *partial_c;
1149.        float *dev_a, *dev_b, *dev_partial_c;
1150.        float elapsedTime;
1151.
1152.        HANDLE_ERROR(cudaEventCreate(&start));
1153.        HANDLE_ERROR(cudaEventCreate(&stop));
1154.
1155.        //allocate memory on the CPU side
1156.        a = (float*)malloc(size * sizeof(float));
1157.        b = (float*)malloc(size * sizeof(float));
1158.        partial_c = (float*)malloc(blocksPerGrid * sizeof(float));
1159.
1160.        //allocate the memory on the GPU
1161.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, size * sizeof(float)));
1162.        HANDLE_ERROR(cudaMalloc((void**)&dev_b, size * sizeof(float)));
1163.        HANDLE_ERROR(cudaMalloc((void**)&dev_partial_c, blocksPerGrid * sizeof(float)));
1164.
1165.        //fill in the host memory with data
1166.        for (int i=0; i<size; i++) {
1167.            a[i] = i;
```

关闭

```
1168.          b[i] = i * 2;
1169.      }
1170.
1171.      HANDLE_ERROR(cudaEventRecord(start, 0));
1172.      //copy the arrays 'a' and 'b' to the GPU
1173.      HANDLE_ERROR(cudaMemcpy(dev_a, a, size * sizeof(float), cudaMemcpyHostToDevice));
1174.      HANDLE_ERROR(cudaMemcpy(dev_b, b, size * sizeof(float), cudaMemcpyHostToDevice));
1175.
1176.      dot_kernel<<<blocksPerGrid, threadsPerBlock>>>(size, dev_a, dev_b, dev_partial_c);
1177.      //copy the array 'c' back from the GPU to the CPU
1178.      HANDLE_ERROR(cudaMemcpy(partial_c, dev_partial_c,blocksPerGrid*sizeof(float), cudaMemc
1179.
1180.      HANDLE_ERROR(cudaEventRecord(stop, 0));
1181.      HANDLE_ERROR(cudaEventSynchronize(stop));
1182.      HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime,start, stop));
1183.
1184.      //finish up on the CPU side
1185.      c = 0;
1186.      for (int i=0; i<blocksPerGrid; i++) {
1187.          c += partial_c[i];
1188.      }
1189.
1190.      HANDLE_ERROR(cudaFree(dev_a));
1191.      HANDLE_ERROR(cudaFree(dev_b));
1192.      HANDLE_ERROR(cudaFree(dev_partial_c));
1193.
1194.      //free memory on the CPU side
1195.      free(a);
1196.      free(b);
1197.      free(partial_c);
1198.
1199.      //free events
1200.      HANDLE_ERROR(cudaEventDestroy(start));
1201.      HANDLE_ERROR(cudaEventDestroy(stop));
1202.
1203.      printf("Value calculated:  %f\n", c);
1204.
1205.      return elapsedTime;
1206.  }
1207.
1208.  float cuda_host_alloc_test(int size)
1209.  {
1210.      cudaEvent_t start, stop;
1211.      float *a, *b, c, *partial_c;
1212.      float *dev_a, *dev_b, *dev_partial_c;
1213.      float elapsedTime;
1214.
1215.      HANDLE_ERROR(cudaEventCreate(&start));
1216.      HANDLE_ERROR(cudaEventCreate(&stop));
1217.
1218.      //allocate the memory on the CPU
1219.      HANDLE_ERROR(cudaHostAlloc((void**)&a, size*sizeof(float), cudaHostAllocWriteCombined
1220.      HANDLE_ERROR(cudaHostAlloc((void**)&b, size*sizeof(float), cudaHostAllocWriteCombined
1221.      HANDLE_ERROR(cudaHostAlloc((void**)&partial_c, blocksPerGrid*sizeof(float), cudaHostA
1222.
1223.      //find out the GPU pointers
1224.      HANDLE_ERROR(cudaHostGetDevicePointer(&dev_a, a, 0));
1225.      HANDLE_ERROR(cudaHostGetDevicePointer(&dev_b, b, 0));
1226.      HANDLE_ERROR( cudaHostGetDevicePointer(&dev_partial_c, partial_c, 0));
1227.
1228.      //fill in the host memory with data
1229.      for (int i=0; i<size; i++) {
1230.          a[i] = i;
1231.          b[i] = i*2;
1232.      }
1233.
1234.      HANDLE_ERROR(cudaEventRecord(start, 0));
1235.
1236.      dot_kernel<<<blocksPerGrid, threadsPerBlock>>>(size, dev_a, dev_b, dev_partial_c);
1237.
1238.      HANDLE_ERROR(cudaThreadSynchronize());
1239.      HANDLE_ERROR(cudaEventRecord(stop, 0));
1240.      HANDLE_ERROR(cudaEventSynchronize(stop));
1241.      HANDLE_ERROR(cudaEventElapsedTime(&elapsedTime,start, stop));
1242.
1243.      //finish up on the CPU side
1244.      c = 0;
1245.      for (int i=0; i<blocksPerGrid; i++) {
1246.          c += partial_c[i];
```

关闭

```
1247.            }
1248.
1249.        HANDLE_ERROR(cudaFreeHost(a));
1250.        HANDLE_ERROR(cudaFreeHost(b));
1251.        HANDLE_ERROR(cudaFreeHost(partial_c));
1252.
1253.        // free events
1254.        HANDLE_ERROR(cudaEventDestroy(start));
1255.        HANDLE_ERROR(cudaEventDestroy(stop));
1256.
1257.        printf("Value calculated:  %f\n", c);
1258.
1259.        return elapsedTime;
1260.    }
1261.
1262.    int test20()
1263.    {
1264.        cudaDeviceProp prop;
1265.        int whichDevice;
1266.        HANDLE_ERROR(cudaGetDevice(&whichDevice));
1267.        HANDLE_ERROR(cudaGetDeviceProperties(&prop, whichDevice));
1268.        if (prop.canMapHostMemory != 1) {
1269.            printf( "Device can not map memory.\n" );
1270.            return 0;
1271.        }
1272.
1273.        HANDLE_ERROR(cudaSetDeviceFlags(cudaDeviceMapHost));
1274.
1275.        //try it with malloc
1276.        float elapsedTime = malloc_test(NUM);
1277.        printf("Time using cudaMalloc:  %3.1f ms\n", elapsedTime);
1278.
1279.        //now try it with cudaHostAlloc
1280.        elapsedTime = cuda_host_alloc_test(NUM);
1281.        printf("Time using cudaHostAlloc:  %3.1f ms\n", elapsedTime);
1282.
1283.        return 0;
1284.    }
1285.
1286.    void* routine(void *pvoidData)
1287.    {
1288.        DataStruct *data = (DataStruct*)pvoidData;
1289.        HANDLE_ERROR(cudaSetDevice(data->deviceID));
1290.
1291.        int size = data->size;
1292.        float *a, *b, c, *partial_c;
1293.        float *dev_a, *dev_b, *dev_partial_c;
1294.
1295.        //allocate memory on the CPU side
1296.        a = data->a;
1297.        b = data->b;
1298.        partial_c = (float*)malloc(blocksPerGrid * sizeof(float));
1299.
1300.        //allocate the memory on the GPU
1301.        HANDLE_ERROR(cudaMalloc((void**)&dev_a, size * sizeof(float)));
1302.        HANDLE_ERROR(cudaMalloc((void**)&dev_b, size * sizeof(float)));
1303.        HANDLE_ERROR(cudaMalloc((void**)&dev_partial_c, blocksPerGrid*sizeof(float)));
1304.
1305.        //copy the arrays 'a' and 'b' to the GPU
1306.        HANDLE_ERROR(cudaMemcpy(dev_a, a, size*sizeof(float), cudaMemcpyHostToDevice));
1307.        HANDLE_ERROR(cudaMemcpy(dev_b, b, size*sizeof(float), cudaMemcpyHostToDevice));
1308.
1309.        dot_kernel<<<blocksPerGrid, threadsPerBlock>>>(size, dev_a, dev_b, dev_partial_c);
1310.        //copy the array 'c' back from the GPU to the CPU
1311.        HANDLE_ERROR(cudaMemcpy( partial_c, dev_partial_c, blocksPerGrid * sizeof(float), cuda
1312.
1313.        //finish up on the CPU side
1314.        c = 0;
1315.        for (int i=0; i<blocksPerGrid; i++) {
1316.            c += partial_c[i];
1317.        }
1318.
1319.        HANDLE_ERROR(cudaFree(dev_a));
1320.        HANDLE_ERROR(cudaFree(dev_b));
1321.        HANDLE_ERROR(cudaFree(dev_partial_c));
1322.
1323.        //free memory on the CPU side
1324.        free(partial_c);
1325.
```

关闭

```
1326.        data->returnValue = c;
1327.        return 0;
1328.    }
1329.
1330.    int test21()
1331.    {
1332.        int deviceCount;
1333.        HANDLE_ERROR(cudaGetDeviceCount(&deviceCount));
1334.        if (deviceCount < 2) {
1335.            printf("We need at least two compute 1.0 or greater devices, but only found %d\n",
1336.            return 0;
1337.        }
1338.
1339.        float *a = (float*)malloc(sizeof(float) * NUM);
1340.        HANDLE_NULL(a);
1341.        float *b = (float*)malloc(sizeof(float) * NUM);
1342.        HANDLE_NULL(b);
1343.
1344.        //fill in the host memory with data
1345.        for (int i=0; i<NUM; i++) {
1346.            a[i] = i;
1347.            b[i] = i*2;
1348.        }
1349.
1350.        //prepare for multithread
1351.        DataStruct  data[2];
1352.        data[0].deviceID = 0;
1353.        data[0].size = NUM/2;
1354.        data[0].a = a;
1355.        data[0].b = b;
1356.
1357.        data[1].deviceID = 1;
1358.        data[1].size = NUM/2;
1359.        data[1].a = a + NUM/2;
1360.        data[1].b = b + NUM/2;
1361.
1362.        CUTThread thread = start_thread(routine, &(data[0]));
1363.        routine(&(data[1]));
1364.        end_thread(thread);
1365.
1366.        //free memory on the CPU side
1367.        free(a);
1368.        free(b);
1369.
1370.        printf("Value calculated:  %f\n", data[0].returnValue + data[1].returnValue);
1371.
1372.        return 0;
1373.    }
```

funset.cuh:

```cpp
01.    #ifndef _FUNSET_CUH_
02.    #define _FUNSET_CUH_
03.
04.    #include <stdio.h>
05.    #include "cpu_anim.h"
06.
07.    #define NUM  33 * 1024 * 1024//1024*1024//33 * 1024//10
08.    #define DIM 1024//1000
09.    #define PI 3.1415926535897932f
10.    #define imin(a, b) (a < b ? a : b)
11.    const int threadsPerBlock = 256;
12.    const int blocksPerGrid = imin(32, (NUM/2+threadsPerBlock-1) / threadsPerBlock);
       //imin(32, (NUM + threadsPerBlock - 1) / threadsPerBlock);
13.    #define rnd(x) (x * rand() / RAND_MAX)
14.    #define INF 2e10f
15.    #define SPHERES 20
16.    #define MAX_TEMP 1.0f
17.    #define MIN_TEMP 0.0001f
18.    #define SPEED 0.25f
19.    #define SIZE (100*1024*1024)
20.    #define FULL_DATA_SIZE (NUM*20)
21.
22.    //__global__关键字将告诉编译器，函数应该编译为在设备而不是主机上运行
23.    __global__ void add(int a, int b, int* c);
```

关闭

```
 24.   __global__ void add_blockIdx(int* a, int* b, int* c);
 25.   __global__ void add_threadIdx(int* a, int* b, int* c);
 26.   __global__ void add_blockIdx_threadIdx(int* a, int* b, int* c);
 27.
 28.   struct cuComplex {
 29.       float r, i;
 30.
 31.       __device__ cuComplex(float a, float b) : r(a), i(b)  {}
 32.
 33.       __device__ float magnitude2(void)
 34.       {
 35.           return r * r + i * i;
 36.       }
 37.
 38.       __device__ cuComplex operator*(const cuComplex& a)
 39.       {
 40.           return cuComplex(r*a.r - i*a.i, i*a.r + r*a.i);
 41.       }
 42.
 43.       __device__ cuComplex operator+(const cuComplex& a)
 44.       {
 45.           return cuComplex(r+a.r, i+a.i);
 46.       }
 47.   };
 48.
 49.   __device__ int julia(int x, int y);
 50.   __global__ void kernel_julia(unsigned char *ptr);
 51.   __global__ void ripple_kernel(unsigned char *ptr, int ticks);
 52.   struct Sphere;
 53.
 54.   struct DataBlock {
 55.       unsigned char *dev_bitmap;
 56.       CPUAnimBitmap *bitmap;
 57.       Sphere *s;
 58.   };
 59.
 60.   void generate_frame(DataBlock *d, int ticks);
 61.   void cleanup(DataBlock *d);
 62.
 63.   __global__ void dot_kernel(float *a, float *b, float *c);
 64.   __global__ void julia_kernel(unsigned char *ptr);
 65.
 66.   //通过一个数据结构对球面建模
 67.   struct Sphere {
 68.       float r,b,g;
 69.       float radius;
 70.       float x,y,z;
 71.       __device__ float hit(float ox, float oy, float *n)
 72.       {
 73.           float dx = ox - x;
 74.           float dy = oy - y;
 75.           if (dx*dx + dy*dy < radius*radius) {
 76.               float dz = sqrtf(radius*radius - dx*dx - dy*dy);
 77.               *n = dz / sqrtf(radius * radius);
 78.               return dz + z;
 79.           }
 80.           return -INF;
 81.       }
 82.   };
 83.
 84.   //声明为常量内存,__constant__将把变量的访问限制为只读
 85.   __constant__ Sphere s[SPHERES];
 86.
 87.   __global__ void RayTracing_kernel(Sphere *s, unsigned char *ptr);
 88.   __global__ void RayTracing_kernel(unsigned char *ptr);
 89.
 90.   //these exist on the GPU side
 91.   texture<float> texConstSrc, texIn, texOut;
 92.   texture<float, 2> texConstSrc2, texIn2, texOut2;
 93.
 94.   //this kernel takes in a 2-d array of floats it updates the value-of-interest by a
 95.   //scaled value based on itself and its nearest neighbors
 96.   __global__ void Heat_blend_kernel(float *dst, bool dstOut);
 97.   __global__ void blend_kernel(float *dst, bool dstOut);
 98.   __global__ void Heat_copy_const_kernel(float *iptr);
 99.   __global__ void copy_const_kernel(float *iptr);
100.
101.   struct Heat_DataBlock {
102.       unsigned char   *output_bitmap;
```

关闭

```cpp
103.        float          *dev_inSrc;
104.        float          *dev_outSrc;
105.        float          *dev_constSrc;
106.        CPUAnimBitmap  *bitmap;
107.
108.        cudaEvent_t    start, stop;
109.        float          totalTime;
110.        float          frames;
111. };
112.
113. //globals needed by the update routine
114. struct DataBlock_opengl {
115.        float          *dev_inSrc;
116.        float          *dev_outSrc;
117.        float          *dev_constSrc;
118.
119.        cudaEvent_t    start, stop;
120.        float          totalTime;
121.        float          frames;
122. };
123.
124. void Heat_anim_gpu(Heat_DataBlock *d, int ticks);
125. void anim_gpu(Heat_DataBlock *d, int ticks);
126. //clean up memory allocated on the GPU
127. void Heat_anim_exit(Heat_DataBlock *d);
128. void anim_exit(Heat_DataBlock *d);
129. void generate_frame_opengl(uchar4 *pixels, void*, int ticks);
130. __global__ void ripple_kernel_opengl(uchar4 *ptr, int ticks);
131. __global__ void Heat_blend_kernel_opengl(float *dst, bool dstOut);
132. __global__ void Heat_copy_const_kernel_opengl(float *iptr);
133. void anim_gpu_opengl(uchar4* outputBitmap, DataBlock_opengl *d, int ticks);
134. void anim_exit_opengl(DataBlock_opengl *d);
135. __global__ void histo_kernel(unsigned char *buffer, long size, unsigned int *histo);
136. __global__ void singlestream_kernel(int *a, int *b, int *c);
137. __global__ void dot_kernel(int size, float *a, float *b, float *c);
138.
139. struct DataStruct{
140.        int      deviceID;
141.        int      size;
142.        float    *a;
143.        float    *b;
144.        float    returnValue;
145. };
146.
147. #endif //_FUNSET_CUH_
```

funset.cu:

```cpp
[cpp]

01. #include "funset.cuh"
02. #include <stdio.h>
03.
04. __global__ void add(int a, int b, int* c)
05. {
06.     *c = a + b;
07. }
08.
09. //__global__：从主机上调用并在设备上运行
10. __global__ void add_blockIdx(int* a, int* b, int* c)
11. {
12.     //计算该索引处的数据
13.     //变量blockIdx，是一个内置变量，在CUDA运行时已经预先定义了这个变量
14.     //此变量中包含的值就是当前执行设备代码的线程块的索引
15.     int tid = blockIdx.x;//this thread handles the data at its thread id
16.     if (tid < NUM)
17.         c[tid] = a[tid] + b[tid];
18. }
19.
20. //__device__：表示代码将在GPU而不是主机上运行，
21. //由于此函数已声明为__device__函数，因此只能从其它__device__函数或者
22. //从__global__函数中调用它们
23. __device__ int julia(int x, int y)
24. {
25.     const float scale = 1.5;
26.     float jx = scale * (float)(DIM/2 - x)/(DIM/2);
27.     float jy = scale * (float)(DIM/2 - y)/(DIM/2);
```

关闭

```
28.
29.        cuComplex c(-0.8, 0.156);
30.        cuComplex a(jx, jy);
31.
32.        int i = 0;
33.        for (i=0; i<200; i++) {
34.            a = a * a + c;
35.
36.            if (a.magnitude2() > 1000)
37.                return 0;
38.        }
39.
40.        return 1;
41.    }
42.
43.    __global__ void kernel_julia(unsigned char *ptr)
44.    {
45.        //map from blockIdx to pixel position
46.        int x = blockIdx.x;
47.        int y = blockIdx.y;
48.        //gridDim为内置变量，对所有的线程块来说，gridDim是一个常数，用来保存线程格每一维的大小
49.        //此处gridDim的值是(DIM, DIM)
50.        int offset = x + y * gridDim.x;
51.
52.        //now calculate the value at that position
53.        int juliaValue = julia(x, y);
54.
55.        ptr[offset*4 + 0] = 255 * juliaValue;
56.        ptr[offset*4 + 1] = 0;
57.        ptr[offset*4 + 2] = 0;
58.        ptr[offset*4 + 3] = 255;
59.    }
60.
61.    __global__ void add_threadIdx(int* a, int* b, int* c)
62.    {
63.        //使用线程索引来对数据进行索引而非通过线程块索引(blockIdx.x)
64.        int tid = threadIdx.x;
65.
66.        if (tid < NUM)
67.            c[tid] = a[tid] + b[tid];
68.    }
69.
70.    __global__ void add_blockIdx_threadIdx(int* a, int* b, int* c)
71.    {
72.        int tid = threadIdx.x + blockIdx.x * blockDim.x;
73.
74.        if (tid == 0) {
75.            printf("blockDim.x = %d, gridDim.x = %d\n", blockDim.x, gridDim.x);
76.        }
77.
78.        while (tid < NUM) {
79.            c[tid] = a[tid] + b[tid];
80.
81.            tid += blockDim.x * gridDim.x;
82.        }
83.    }
84.
85.    __global__ void ripple_kernel(unsigned char *ptr, int ticks)
86.    {
87.        // map from threadIdx/BlockIdx to pixel position
88.        //将线程和线程块的索引映射到图像坐标
89.        //对x和y的值进行线性化从而得到输出缓冲区中的一个偏移
90.        int x = threadIdx.x + blockIdx.x * blockDim.x;
91.        int y = threadIdx.y + blockIdx.y * blockDim.y;
92.        int offset = x + y * blockDim.x * gridDim.x;
93.
94.        // now calculate the value at that position
95.        //生成一个随时间变化的正弦曲线"波纹"
96.        float fx = x - DIM/2;
97.        float fy = y - DIM/2;
98.        float d = sqrtf(fx * fx + fy * fy);
99.        unsigned char grey = (unsigned char)(128.0f + 127.0f * cos(d/10.0f - ticks/7.0f) / (d
100.
101.        ptr[offset*4 + 0] = grey;
102.        ptr[offset*4 + 1] = grey;
103.        ptr[offset*4 + 2] = grey;
104.        ptr[offset*4 + 3] = 255;
105.    }
106.
```

关闭

```cuda
107.    __global__ void dot_kernel(float *a, float *b, float *c)
108.    {
109.        //声明了一个共享内存缓冲区，它将保存每个线程计算的加和值
110.        __shared__ float cache[threadsPerBlock];
111.        int tid = threadIdx.x + blockIdx.x * blockDim.x;
112.        int cacheIndex = threadIdx.x;
113.
114.        float temp = 0;
115.        while (tid < NUM) {
116.            temp += a[tid] * b[tid];
117.            tid += blockDim.x * gridDim.x;
118.        }
119.
120.        //set the cache values
121.        cache[cacheIndex] = temp;
122.
123.        //synchronize threads in this block
124.        //对线程块中的线程进行同步
125.        //这个函数将确保线程块中的每个线程都执行完__syncthreads()前面的语句后，才会执行下一条语句
126.        __syncthreads();
127.
128.        //for reductions(归
    约), threadsPerBlock must be a power of 2 because of the following code
129.        int i = blockDim.x/2;
130.        while (i != 0) {
131.            if (cacheIndex < i)
132.                cache[cacheIndex] += cache[cacheIndex + i];
133.            //在循环迭代中更新了共享内存变量cache，并且在循环的下一次迭代开始之前，
134.            //需要确保当前迭代中所有线程的更新操作都已经完成
135.            __syncthreads();
136.            i /= 2;
137.        }
138.
139.        if (cacheIndex == 0)
140.            c[blockIdx.x] = cache[0];
141.    }
142.
143.    __global__ void julia_kernel(unsigned char *ptr)
144.    {
145.        //map from threadIdx/BlockIdx to pixel position
146.        int x = threadIdx.x + blockIdx.x * blockDim.x;
147.        int y = threadIdx.y + blockIdx.y * blockDim.y;
148.        int offset = x + y * blockDim.x * gridDim.x;
149.
150.        __shared__ float shared[16][16];
151.
152.        //now calculate the value at that position
153.        const float period = 128.0f;
154.
155.        shared[threadIdx.x]
    [threadIdx.y] = 255 * (sinf(x*2.0f*PI/ period) + 1.0f) *(sinf(y*2.0f*PI/ period) + 1.0f) ,
156.
157.        //removing this syncthreads shows graphically what happens
158.        //when it doesn't exist.this is an example of why we need it.
159.        __syncthreads();
160.
161.        ptr[offset*4 + 0] = 0;
162.        ptr[offset*4 + 1] = shared[15 - threadIdx.x][15 - threadIdx.y];
163.        ptr[offset*4 + 2] = 0;
164.        ptr[offset*4 + 3] = 255;
165.    }
166.
167.    __global__ void RayTracing_kernel(Sphere *s, unsigned char *ptr)
168.    {
169.        //map from threadIdx/BlockIdx to pixel position
170.        int x = threadIdx.x + blockIdx.x * blockDim.x;        关闭
171.        int y = threadIdx.y + blockIdx.y * blockDim.y;
172.        int offset = x + y * blockDim.x * gridDim.x;
173.        float ox = (x - DIM/2);
174.        float oy = (y - DIM/2);
175.
176.        float r=0, g=0, b=0;
177.        float maxz = -INF;
178.
179.        for (int i = 0; i < SPHERES; i++) {
180.            float n;
181.            float t = s[i].hit(ox, oy, &n);
182.            if (t > maxz) {
183.                float fscale = n;
```

```
184.                r = s[i].r * fscale;
185.                g = s[i].g * fscale;
186.                b = s[i].b * fscale;
187.                maxz = t;
188.            }
189.        }
190.
191.     ptr[offset*4 + 0] = (int)(r * 255);
192.     ptr[offset*4 + 1] = (int)(g * 255);
193.     ptr[offset*4 + 2] = (int)(b * 255);
194.     ptr[offset*4 + 3] = 255;
195. }
196.
197. __global__ void RayTracing_kernel(unsigned char *ptr)
198. {
199.     //map from threadIdx/BlockIdx to pixel position
200.     int x = threadIdx.x + blockIdx.x * blockDim.x;
201.     int y = threadIdx.y + blockIdx.y * blockDim.y;
202.     int offset = x + y * blockDim.x * gridDim.x;
203.     float ox = (x - DIM/2);
204.     float oy = (y - DIM/2);
205.
206.     float r=0, g=0, b=0;
207.     float maxz = -INF;
208.
209.     for(int i = 0; i < SPHERES; i++) {
210.         float n;
211.         float t = s[i].hit(ox, oy, &n);
212.         if (t > maxz) {
213.             float fscale = n;
214.             r = s[i].r * fscale;
215.             g = s[i].g * fscale;
216.             b = s[i].b * fscale;
217.             maxz = t;
218.         }
219.     }
220.
221.     ptr[offset*4 + 0] = (int)(r * 255);
222.     ptr[offset*4 + 1] = (int)(g * 255);
223.     ptr[offset*4 + 2] = (int)(b * 255);
224.     ptr[offset*4 + 3] = 255;
225. }
226.
227. __global__ void Heat_blend_kernel(float *dst, bool dstOut)
228. {
229.     //map from threadIdx/BlockIdx to pixel position
230.     int x = threadIdx.x + blockIdx.x * blockDim.x;
231.     int y = threadIdx.y + blockIdx.y * blockDim.y;
232.     int offset = x + y * blockDim.x * gridDim.x;
233.
234.     int left = offset - 1;
235.     int right = offset + 1;
236.     if (x == 0) left++;
237.     if (x == DIM-1) right--;
238.
239.     int top = offset - DIM;
240.     int bottom = offset + DIM;
241.     if (y == 0) top += DIM;
242.     if (y == DIM-1) bottom -= DIM;
243.
244.     float t, l, c, r, b;
245.
246.     if (dstOut) {
247.         //tex1Dfetch是编译器内置函数，从设备内存取纹理
248.         t = tex1Dfetch(texIn, top);
249.         l = tex1Dfetch(texIn, left);
250.         c = tex1Dfetch(texIn, offset);
251.         r = tex1Dfetch(texIn, right);
252.         b = tex1Dfetch(texIn, bottom);
253.
254.     } else {
255.         t = tex1Dfetch(texOut, top);
256.         l = tex1Dfetch(texOut, left);
257.         c = tex1Dfetch(texOut, offset);
258.         r = tex1Dfetch(texOut, right);
259.         b = tex1Dfetch(texOut, bottom);
260.     }
261.
262.     dst[offset] = c + SPEED * (t + b + r + l - 4 * c);
```

关闭

```
263.    }
264.
265.    __global__ void blend_kernel(float *dst, bool dstOut)
266.    {
267.        //map from threadIdx/BlockIdx to pixel position
268.        int x = threadIdx.x + blockIdx.x * blockDim.x;
269.        int y = threadIdx.y + blockIdx.y * blockDim.y;
270.        int offset = x + y * blockDim.x * gridDim.x;
271.
272.        float t, l, c, r, b;
273.        if (dstOut) {
274.            t = tex2D(texIn2, x, y-1);
275.            l = tex2D(texIn2, x-1, y);
276.            c = tex2D(texIn2, x, y);
277.            r = tex2D(texIn2, x+1, y);
278.            b = tex2D(texIn2, x, y+1);
279.        } else {
280.            t = tex2D(texOut2, x, y-1);
281.            l = tex2D(texOut2, x-1, y);
282.            c = tex2D(texOut2, x, y);
283.            r = tex2D(texOut2, x+1, y);
284.            b = tex2D(texOut2, x, y+1);
285.        }
286.        dst[offset] = c + SPEED * (t + b + r + l - 4 * c);
287.    }
288.
289.    __global__ void Heat_copy_const_kernel(float *iptr)
290.    {
291.        //map from threadIdx/BlockIdx to pixel position
292.        int x = threadIdx.x + blockIdx.x * blockDim.x;
293.        int y = threadIdx.y + blockIdx.y * blockDim.y;
294.        int offset = x + y * blockDim.x * gridDim.x;
295.
296.        float c = tex1Dfetch(texConstSrc, offset);
297.        if (c != 0)
298.            iptr[offset] = c;
299.    }
300.
301.    __global__ void copy_const_kernel(float *iptr)
302.    {
303.        //map from threadIdx/BlockIdx to pixel position
304.        int x = threadIdx.x + blockIdx.x * blockDim.x;
305.        int y = threadIdx.y + blockIdx.y * blockDim.y;
306.        int offset = x + y * blockDim.x * gridDim.x;
307.
308.        float c = tex2D(texConstSrc2, x, y);
309.        if (c != 0)
310.            iptr[offset] = c;
311.    }
312.
313.    void generate_frame_opengl(uchar4 *pixels, void*, int ticks)
314.    {
315.        dim3 grids(DIM / 16, DIM / 16);
316.        dim3 threads(16, 16);
317.        ripple_kernel_opengl<<<grids, threads>>>(pixels, ticks);
318.    }
319.
320.    __global__ void ripple_kernel_opengl(uchar4 *ptr, int ticks)
321.    {
322.        //map from threadIdx/BlockIdx to pixel position
323.        int x = threadIdx.x + blockIdx.x * blockDim.x;
324.        int y = threadIdx.y + blockIdx.y * blockDim.y;
325.        int offset = x + y * blockDim.x * gridDim.x;
326.
327.        // now calculate the value at that position
328.        float fx = x - DIM / 2;
329.        float fy = y - DIM / 2;
330.        float d = sqrtf(fx * fx + fy * fy);
331.        unsigned char grey = (unsigned char)(128.0f + 127.0f * cos(d/10.0f - ticks/7.0f) / (d/
332.        ptr[offset].x = grey;
333.        ptr[offset].y = grey;
334.        ptr[offset].z = grey;
335.        ptr[offset].w = 255;
336.    }
337.
338.    __global__ void Heat_blend_kernel_opengl(float *dst, bool dstOut)
339.    {
340.        //map from threadIdx/BlockIdx to pixel position
341.        int x = threadIdx.x + blockIdx.x * blockDim.x;
```

关闭

```
342.        int y = threadIdx.y + blockIdx.y * blockDim.y;
343.        int offset = x + y * blockDim.x * gridDim.x;
344.
345.        int left = offset - 1;
346.        int right = offset + 1;
347.        if (x == 0) left++;
348.        if (x == DIM-1) right--;
349.
350.        int top = offset - DIM;
351.        int bottom = offset + DIM;
352.        if (y == 0) top += DIM;
353.        if (y == DIM-1) bottom -= DIM;
354.
355.        float t, l, c, r, b;
356.        if (dstOut) {
357.            t = tex1Dfetch(texIn, top);
358.            l = tex1Dfetch(texIn, left);
359.            c = tex1Dfetch(texIn, offset);
360.            r = tex1Dfetch(texIn, right);
361.            b = tex1Dfetch(texIn, bottom);
362.
363.        } else {
364.            t = tex1Dfetch(texOut, top);
365.            l = tex1Dfetch(texOut, left);
366.            c = tex1Dfetch(texOut, offset);
367.            r = tex1Dfetch(texOut, right);
368.            b = tex1Dfetch(texOut, bottom);
369.        }
370.        dst[offset] = c + SPEED * (t + b + r + l - 4 * c);
371.    }
372.
373.    __global__ void Heat_copy_const_kernel_opengl(float *iptr)
374.    {
375.        int x = threadIdx.x + blockIdx.x * blockDim.x;
376.        int y = threadIdx.y + blockIdx.y * blockDim.y;
377.        int offset = x + y * blockDim.x * gridDim.x;
378.
379.        float c = tex1Dfetch(texConstSrc, offset);
380.        if (c != 0)
381.            iptr[offset] = c;
382.    }
383.
384.    __global__ void histo_kernel(unsigned char *buffer, long size, unsigned int *histo)
385.    {
386.        //clear out the accumulation buffer called temp since we are launched with 256 threads
387.        //it is easy to clear that memory with one write per thread
388.        __shared__ unsigned int temp[256]; //共享内存缓冲区
389.        temp[threadIdx.x] = 0;
390.        __syncthreads();
391.
392.        //calculate the starting index and the offset to the next block that each thread will
393.        int i = threadIdx.x + blockIdx.x * blockDim.x;
394.        int stride = blockDim.x * gridDim.x;
395.        while (i < size) {
396.            atomicAdd(&temp[buffer[i]], 1);
397.            i += stride;
398.        }
399.
400.        //sync the data from the above writes to shared memory then add the shared memory valu
401.        //the other thread blocks using global memory atomic adds same as before, since we hav
402.        //updating the global histogram is just one write per thread!
403.        __syncthreads();
404.        atomicAdd(&(histo[threadIdx.x]), temp[threadIdx.x]);
405.    }
406.
407.    __global__ void singlestream_kernel(int *a, int *b, i                              关闭
408.    {
409.        int idx = threadIdx.x + blockIdx.x * blockDim.x;
410.        if (idx < NUM) {
411.            int idx1 = (idx + 1) % 256;
412.            int idx2 = (idx + 2) % 256;
413.            float as = (a[idx] + a[idx1] + a[idx2]) / 3.0f;
414.            float bs = (b[idx] + b[idx1] + b[idx2]) / 3.0f;
415.            c[idx] = (as + bs) / 2;
416.        }
417.    }
418.
419.    __global__ void dot_kernel(int size, float *a, float *b, float *c)
420.    {
```

```
421.        __shared__ float cache[threadsPerBlock];
422.        int tid = threadIdx.x + blockIdx.x * blockDim.x;
423.        int cacheIndex = threadIdx.x;
424.
425.        float temp = 0;
426.        while (tid < size) {
427.            temp += a[tid] * b[tid];
428.            tid += blockDim.x * gridDim.x;
429.        }
430.
431.        //set the cache values
432.        cache[cacheIndex] = temp;
433.
434.        //synchronize threads in this block
435.        __syncthreads();
436.
437.        //for reductions(归
约), threadsPerBlock must be a power of 2 because of the following code
438.        int i = blockDim.x / 2;
439.        while (i != 0) {
440.            if (cacheIndex < i)
441.                cache[cacheIndex] += cache[cacheIndex + i];
442.            __syncthreads();
443.            i /= 2;
444.        }
445.
446.        if (cacheIndex == 0)
447.            c[blockIdx.x] = cache[0];
448.    }
```

以上来自于对《GPU高性能编程CUDA实战》书中内容整理。

**GitHub**：https://github.com/fengbingchun/CUDA_Test

<div align="center">

顶　　踩

2　　　　0

</div>

上一篇　Git基础(常用命令)介绍

下一篇　Linux下常用的C/C++开源Socket库

我的同类文章

**CUDA（10）**

| | | | |
|---|---|---|---|
| • CUDA基础介绍 | 2017-01-23 阅读 0 | • windows7 64位机上配置支… | 2016-12-27 阅读 111 |
| • windows7 64位机上安装配… | 2016-12-27 阅读 506 | • Ubuntu14.04 64位机上安装… | 2016-12-23 阅读 186 |
| • Ubuntu14.04 64位机上安装… | 2016-12-23 阅读 433 | • CUDA Runtime API 汇总 | 2015-04-26 阅读 2495 |
| • windows7 64位机上CUDA7.… | 2015-04-09 阅读 17677 | • GPU及GPU通用计算编程模… | 2014-02-21 阅读 4017 |
| • Windows7 32位机上，Open… | 2013-08-08 阅读 17385 | • OpenCV中G… | |

关闭

猜你在找

ArcGIS for javascript 项目实战（环境监测系统）　　　　《GPU高性能编程CUDA实战》

从零开始学习mac系统视频录制编辑软件Camtasia For　《GPU高性能编程CUDA实战》学习笔记十一

查看评论

1楼 luyyuang 2015-10-08 09:44发表

请问你用的CUDA SDK是那个版本的？？

Re: fengbingchun 2015-10-08 11:27发表

回复luyyuang：7.0

您还没有登录,请[登录]或[注册]

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

核心技术类目

| 全部主题 | Hadoop | AWS | 移动游戏 | Java | Android | iOS | Swift | 智能硬件 | Docker | Op |
| VPN | Spark | ERP | IE10 | Eclipse | CRM | JavaScript | 数据库 | Ubuntu | NFC | WAP | jQuery |
| BI | HTML5 | Spring | Apache | .NET | API | HTML | SDK | IIS | Fedora | XML | LBS | Unity |
| Splashtop | UML | components | Windows Mobile | Rails | QEMU | KDE | Cassandra | CloudStack | FTC |
| coremail | OPhone | CouchBase | 云计算 | iOS6 | Rackspace | Web App | SpringSide | Maemo |
| Compuware | 大数据 | aptech | Perl | Tornado | Ruby | Hibernate | ThinkPHP | HBase | Pure | Solr |
| Angular | Cloud Foundry | Redis | Scala | Django | Bootstrap |

关闭