

## OpenCL&Cuda

### 1. 环境配置

#### 1.1 cuda的安装和配置

##### 1.1.1 安装CUDA

##### 1.1.2 配置环境变量

##### 1.1.3 检查cuda是否安装成功

#### 1.2 VS中配置cuda

#### 1.3 VS中配置opencl

### 2. OpenCL&Cuda编程

#### 2.1 编程基础

#### 2.2 编程案例

##### 2.2.1 查看配置信息

##### 2.2.2. 向量运算

实验目的

实验步骤

opencl 代码分析

cuda 代码分析

实验结果

结果分析

##### 2.2.3 Pi计算

实验目的

实验步骤

opencl代码分析

cuda 代码分析

实验结果

结果分析

### 3. 总结

### 4. 附录

cuda 完整代码

openCL 完整代码

姓名：张煜玮      班级：计科1905      学号:19291255

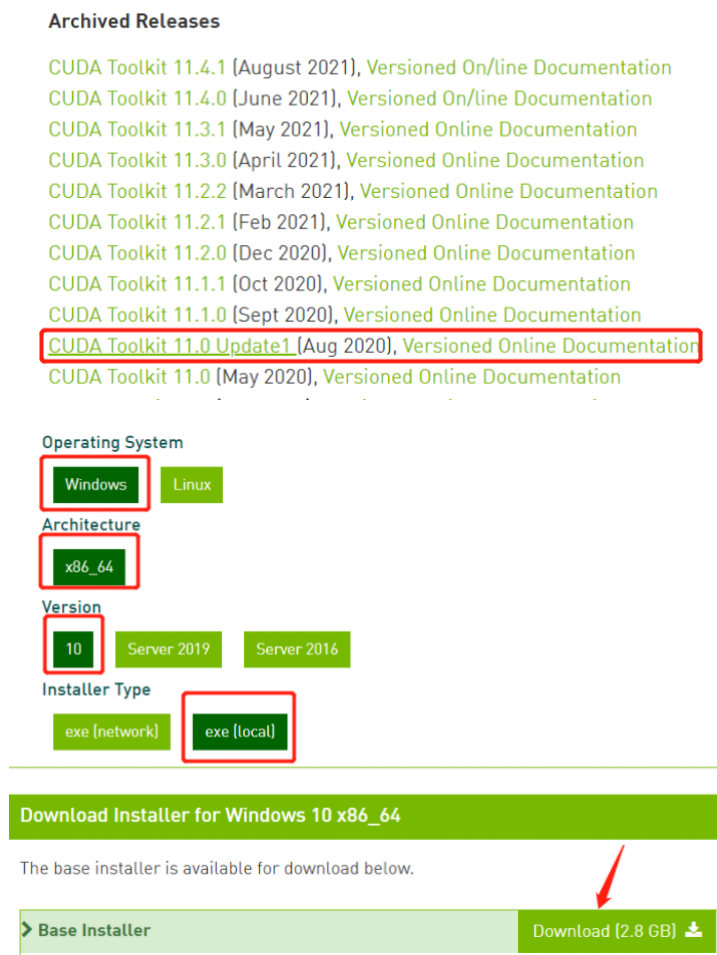
# OpenCL&Cuda

## 1. 环境配置

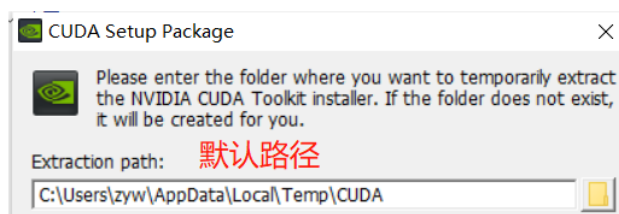
### 1.1 cuda的安装和配置

#### 1.1.1 安装CUDA

1. 从[CUDA Toolkit 官网下载地址](#)上选择合适的版本进行下载



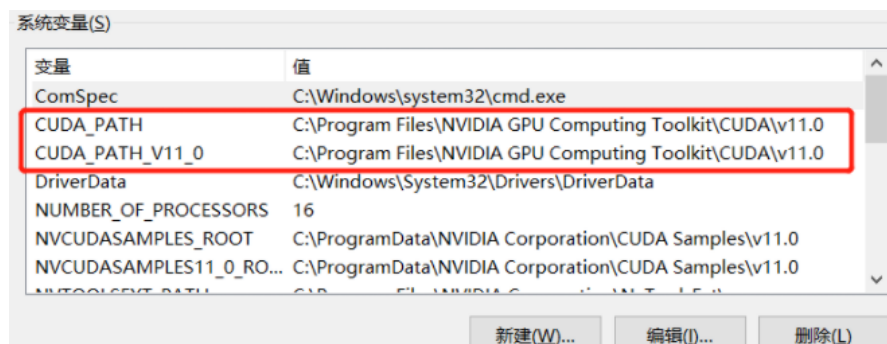
2. 选择默认的路径进行安装



#### 1.1.2 配置环境变量

1. 检查是否安装成功

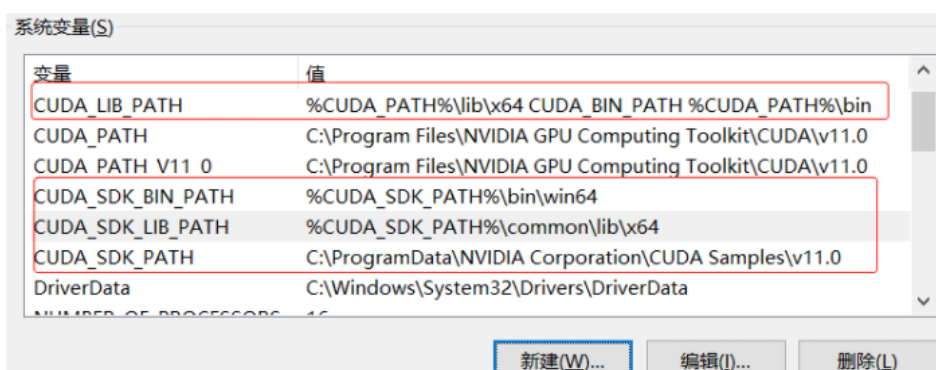
安装成功后在系统的环境变量中会自动添加两个变量，如图：“CUDA\_PATH”，和“CUDA\_PATH\_V11\_0”



## 2. 添加系统变量

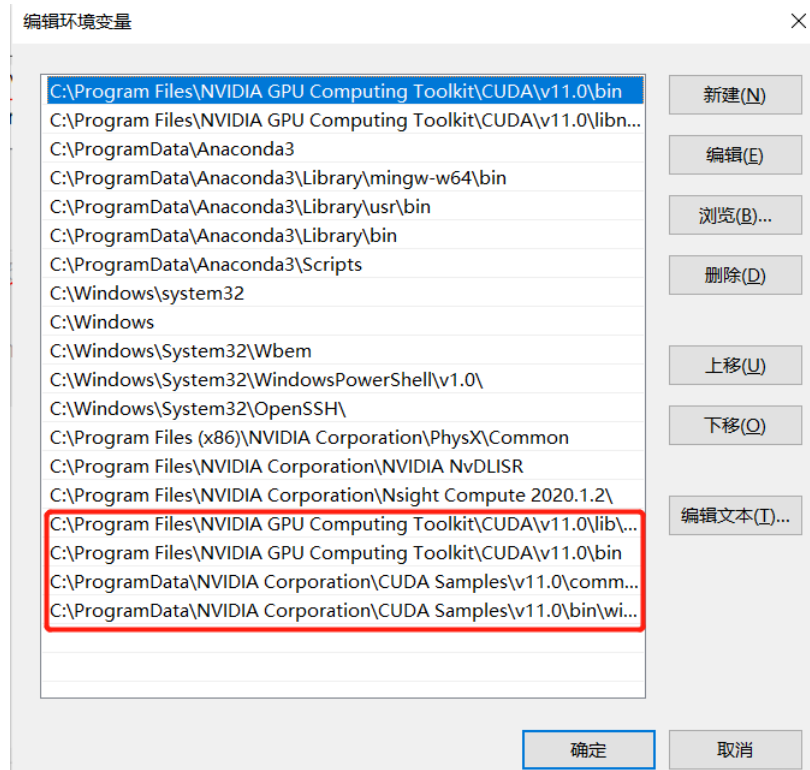
点击“新建”，然后依次添加下述5个系统变量，然后点击“确认”

1	变量	值
2	CUDA_SDK_PATH	C:\ProgramData\NVIDIA Corporation\CUDA Samples\v11.0
3	CUDA_LIB_PATH	%CUDA_PATH%\lib\x64 CUDA_BIN_PATH %CUDA_PATH%\bin
4	CUDA_SDK_BIN_PATH	%CUDA_SDK_PATH%\bin\win64
5	CUDA_SDK_LIB_PATH	%CUDA_SDK_PATH%\common\lib\x64



然后在系统变量Path里添加变量，双击Path添加如下变量

1	C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\lib\x64
2	C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\bin
3	C:\ProgramData\NVIDIA Corporation\CUDA Samples\v11.0\common\lib\x64
4	C:\ProgramData\NVIDIA Corporation\CUDA Samples\v11.0\bin\win64



### 1.1.3 检查cuda是否安装成功

1. 打开cmd命令窗口，切换路径

cd C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\extras\demo\_suite

2. 分别运行如下命令:

bandwidthTest.exe

如下图，Result=Pass

```
C:\Users\zyw>cd C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\extras\demo_suite
C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\extras\demo_suite>bandwidthTest.exe
[CUDA Bandwidth Test] - Starting...
Running on...

Device 0: GeForce RTX 3050 Ti Laptop GPU
Quick Mode

Host to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
  Transfer Size (Bytes)      Bandwidth(MB/s)
  33554432                  6376.5

Device to Host Bandwidth, 1 Device(s)
PINNED Memory Transfers
  Transfer Size (Bytes)      Bandwidth(MB/s)
  33554432                  6310.4

Device to Device Bandwidth, 1 Device(s)
PINNED Memory Transfers
  Transfer Size (Bytes)      Bandwidth(MB/s)
  33554432                  169626.3

Result = PASS
```

键入命令

pass: 安装成功

deviceQuery.exe

如下图，Result=Pass，安装成功。

```
deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 11.2, CUDA Runtime Version = 11.0, NumDevs = 1, Device0 = GeForce RTX 3050 Ti Laptop GPU
Result = PASS
```

## 1.2 VS中配置cuda

1. 在 C:\ProgramData\NVIDIA GPU Computing Toolkit\v10.1\extras\visual\_studio\_integration\CudaProjectVswizards 中找到 Nvda.Vsip.Cudawizards.dll.pkgdef 并创建快捷方式
2. 在 C:\ProgramData\NVIDIA GPU Computing Toolkit\v11.0\extras\visual\_studio\_integration\CudaProjectVswizards\2019 中找到 extension.vsixmanifest 并创建快捷方式
3. 在 D:\Program Files\Microsoft Visual Studio\2019\Professional\Common7\IDE\Extensions 下创建新的文件夹 NVIDIA -> CUDA 11.0 wizards -> 11.0。

(D:) > Program Files > Microsoft Visual Studio > 2019 > Professional > Common7 > IDE > Extensions > NVIDIA > CUDA 11.0 Wizards > 11.0

名称	修改日期	类型	大小
extension.vsixmanifest - 快捷方式	2021/12/16 9:43	快捷方式	2 KB
Nvda.Vsip.Cudawizards.dll.pkgdef - ...	2021/12/16 9:47	快捷方式	2 KB



管理员: x64 Native Tools Command Prompt for VS 2019 - devenv.com /setup /nosetupvstemplates

```
*****
** Visual Studio 2019 Developer Command Prompt v16.11.3
** Copyright (c) 2021 Microsoft Corporation
*****
[vcvarsall.bat] Environment initialized for: 'x64'
C:\Windows\System32>devenv.com /setup /nosetupvstemplates
```

4. 将 C:\Program Files\NVIDIA GPU Computing Toolkit\CUDA\v11.0\extras\visual\_studio\_integration\MSBuildExtensions 下的文件复制，并粘贴到 D:\Program Files\Microsoft Visual Studio\2019\Professional\MSBuild\Microsoft\VC\v160\BuildCustomizations 文件下。

ws-SSD (C:) > Program Files > NVIDIA GPU Computing Toolkit > CUDA > v11.0 > extras > visual\_studio\_integration > MSBuildExtensions

名称	修改日期	类型	大小
CUDA 11.0.props	2020/7/23 18:59	Project Property...	14 KB
CUDA 11.0.targets	2020/7/23 18:59	Project Targets ...	49 KB
CUDA 11.0.xml	2020/7/23 18:59	XML 文档	30 KB
Nvda.Build.CudaTasks.v11.0.dll	2020/7/23 18:59	应用程序扩展	260 KB

D:\ > Program Files > Microsoft Visual Studio > 2019 > Professional > MSBuild > Microsoft > VC > v160 > BuildCustomizations

名称	修改日期	类型	大小
CUDA 11.0.props	2020/7/23 18:59	Project Property...	14 KB
CUDA 11.0.targets	2020/7/23 18:59	Project Targets ...	49 KB
CUDA 11.0.xml	2020/7/23 18:59	XML 文档	30 KB
ImageContentTask.props	2021/9/16 12:00	Project Property...	1 KB
ImageContentTask.targets	2021/9/16 12:00	Project Targets ...	4 KB
ImageContentTask.xml	2021/9/16 12:00	XML 文档	5 KB
Ic.props	2021/9/16 12:00	Project Property...	2 KB
Ic.targets	2021/9/16 12:00	Project Targets ...	4 KB
Ic.xml	2021/9/16 12:00	XML 文档	7 KB
marmasm.props	2021/9/16 12:00	Project Property...	2 KB
marmasm.targets	2021/9/16 12:00	Project Targets ...	4 KB
marmasm.xml	2021/9/16 12:00	XML 文档	9 KB
masm.props	2021/9/16 12:00	Project Property...	2 KB
masm.targets	2021/9/16 12:00	Project Targets ...	5 KB
masm.xml	2021/9/16 12:00	XML 文档	15 KB
MeshContentTask.props	2021/9/16 12:00	Project Property...	1 KB
MeshContentTask.targets	2021/9/16 12:00	Project Targets ...	4 KB
MeshContentTask.xml	2021/9/16 12:00	XML 文档	2 KB
Nvda.Build.CudaTasks.v11.0.dll	2020/7/23 18:59	应用程序扩展	260 KB
ShaderGraphContentTask.props	2021/9/16 12:00	Project Property...	1 KB
ShaderGraphContentTask.targets	2021/9/16 12:00	Project Targets ...	4 KB
ShaderGraphContentTask.xml	2021/9/16 12:00	XML 文档	2 KB

## 5. 新建一个cuda runtime 的文件



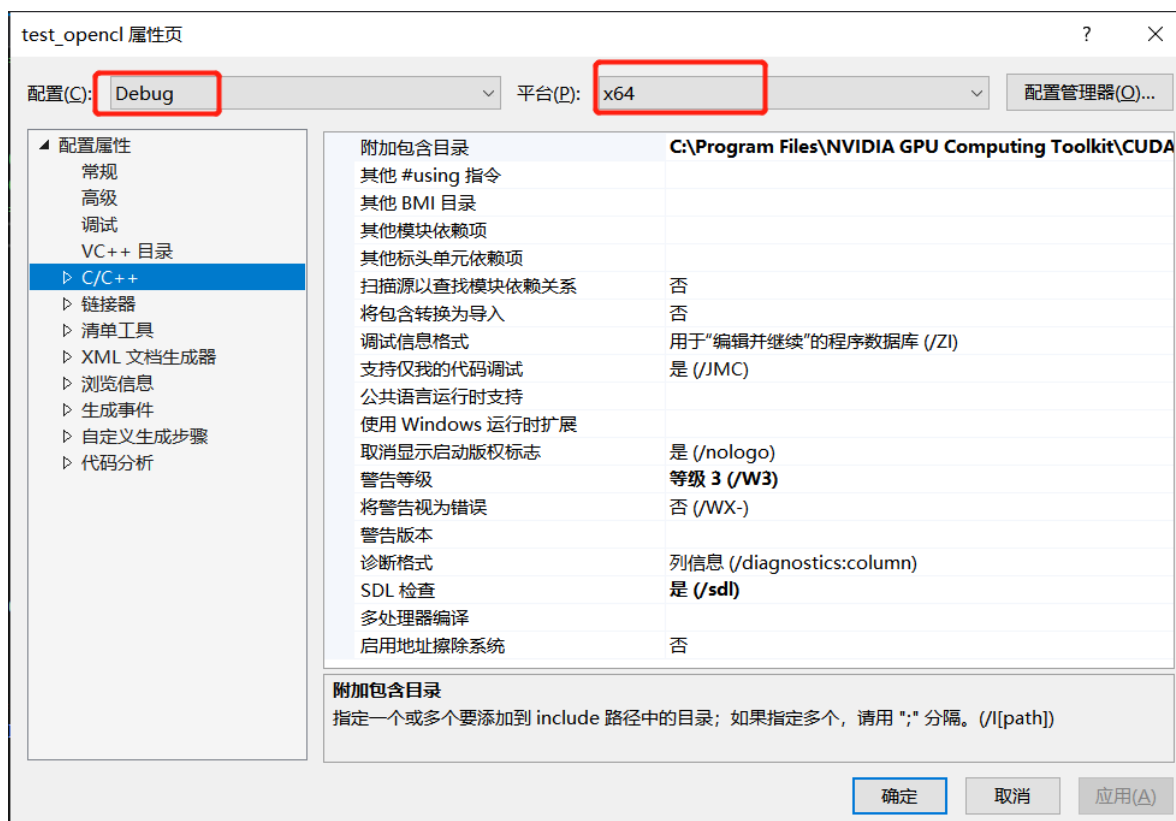
## 6. 测试运行

```
Microsoft Visual Studio 调试控制台
{1, 2, 3, 4, 5} + {10, 20, 30, 40, 50} = {11, 22, 33, 44, 55}
C:\Users\zyw\source\repos\CUDA 11.0 Runtime3\x64\Debug\CUDA 11.0 Runtime3.exe (进程 14436) 已退出, 代码为 0。
要在调试停止时自动关闭控制台, 请启用“工具”->“选项”->“调试”->“调试停止时自动关闭控制台”。
按任意键关闭窗口...
```

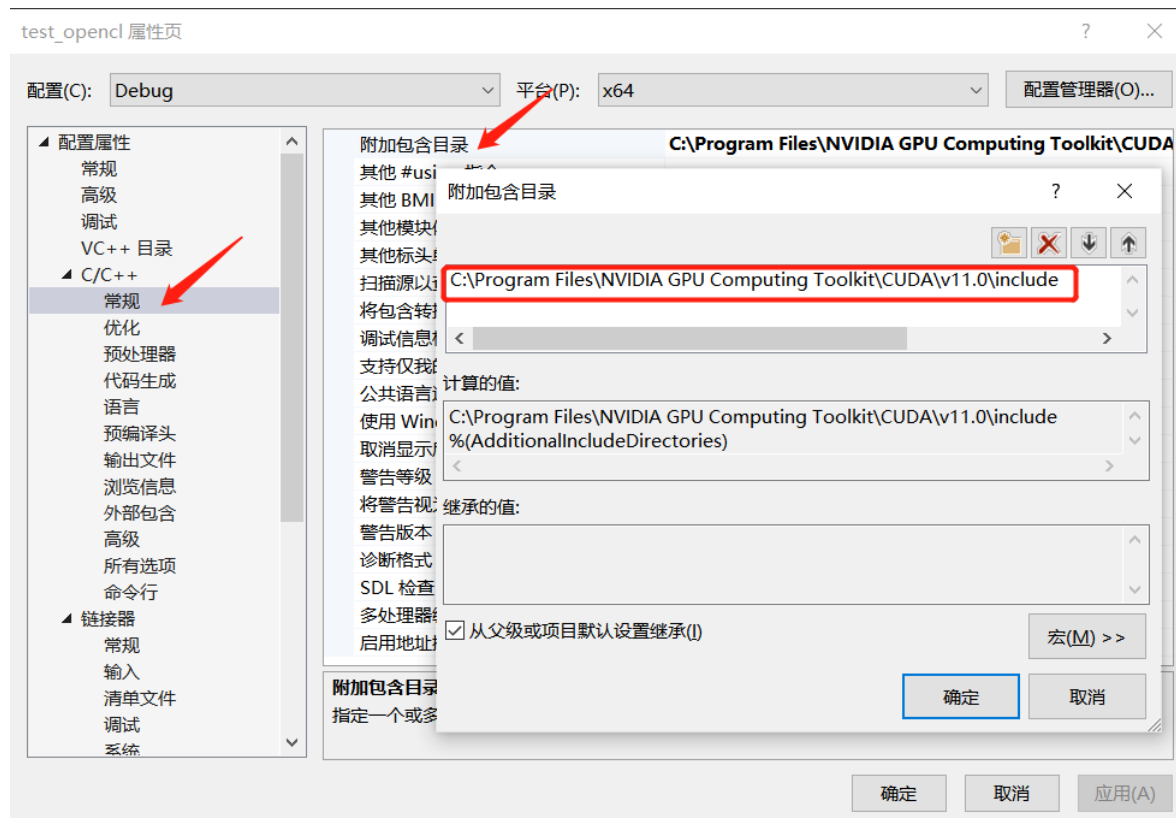
## 1.3 VS中配置opengl

前提: 下载好了cuda

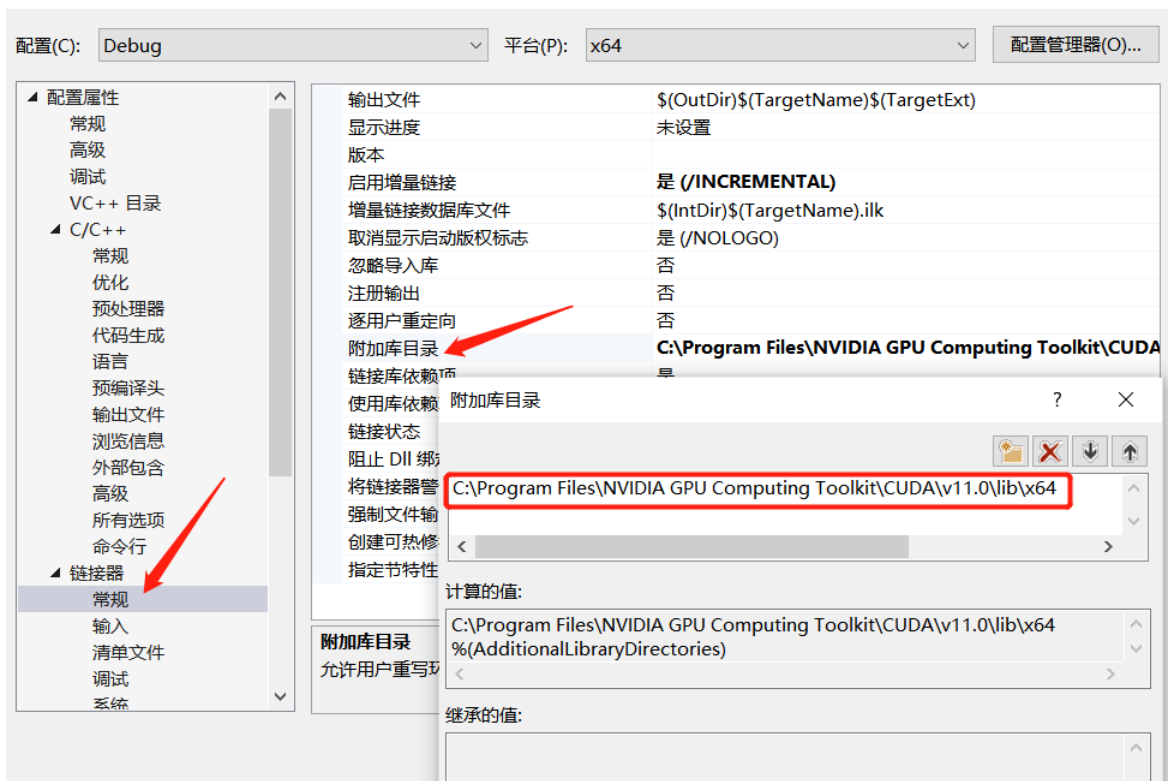
1. 在vs中新建一个项目，然后点击项目属性，选择配置debug,平台x64



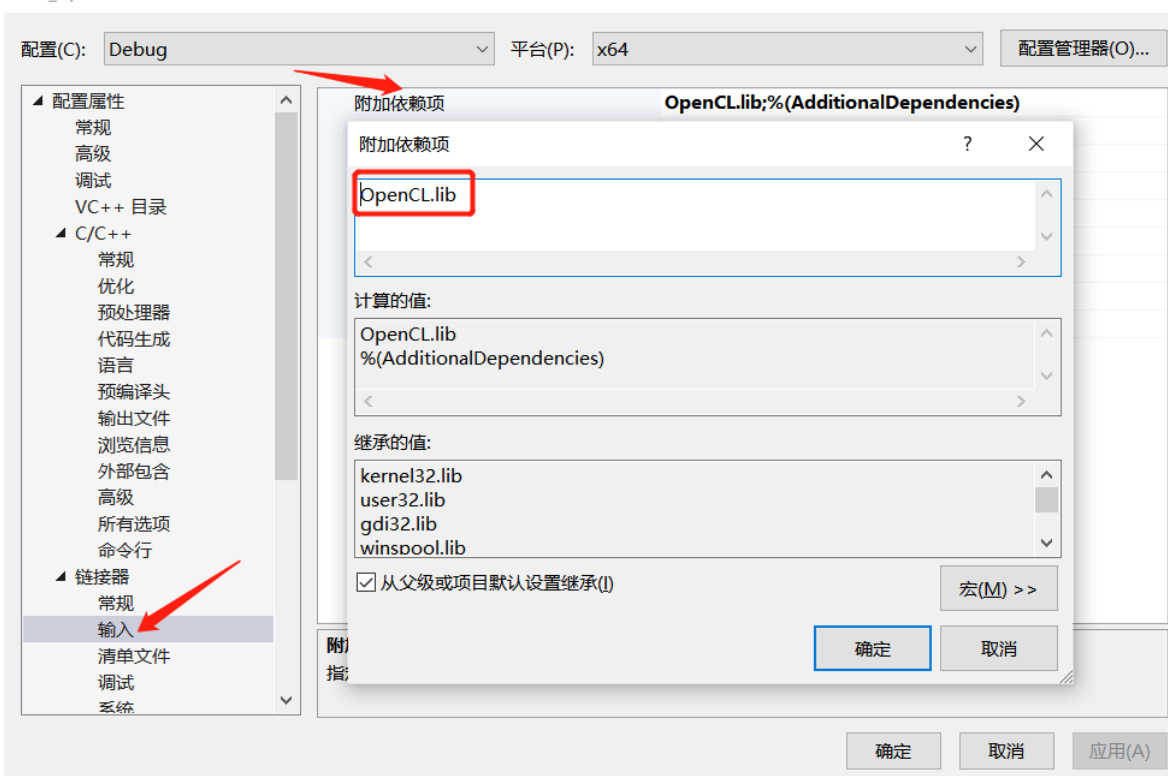
2. 添加附加包含目录



3. 添加附加库目录



#### 4. 添加依赖项



## 2. OpenCL&Cuda编程

### 2.1 编程基础

opengl编程的基本流程

- 选择OpenCL平台并创建一个上下文



平台 (Platform) 是指主机和OpenCL管理框架下的若干个设备构成的可以运行OpenCL程序的完整硬件系统, 这个是跑OpenCL程序的基础, 所以第一步要选择一个可用的OpenCL平台。一台机器上可以不止一个这样的平台, 一个平台也可以不止一个GPU。

`clGetPlatformIDs()`: 用于获取可用的平台

`clCreateContextFromType()`: 创建一个openCL运行时的上下文环境

- **选择设备并创建命令队列**

选择平台并创建好OpenCL上下文环境之后, 要做的事选择运行时用到的设备, 还要创建一个命令队列, 命令队列里定义了设备要完成的操作, 以及各个操作的运行次序。

`clCreateCommandQueue()`: 用于创建一个指定设备上的上下文环境, 第二个参数定义了选择的设备。

- **创建和构建程序对象**

程序对象用来存储与上下文相关联的设备的已编译可执行代码, 同时也完成内核源代码的加载编译工作。

`clCreateProgramWithSource()`: 这个函数会创建一个程序对象, 在创建的同时, 把已经转化成字符串形式的内核源代码加载到该程序对象中。

`clBuildProgram()`: 用于编译指定程序对象中的内核源代码, 编译成功之后, 再把编译代码存储在程序对象中。

- **创建内核和内存对象**

要执行程序对象中的已编译成功的内核运算, 需要在内存中创建内核并分配内核函数的参数, 在GPU上定义内存对象并分配存储空间。

`clCreateKernel`: 创建内核

`clCreateBuffer()`: 分配内存对象的存储空间, 这些对象可以由内核函数直接访问。

- **设置内核数据并执行内核**

创建内核和内存对象之后, 接下来要设置核函数的数据, 并将要执行的内核排队。

`clEnqueueNDRangeKernel()` 用于设置内核函数的所有参与运算的数据。 利用命令队列对要在设备上执行的内核排队。需要注意的是, 执行内核排队之后并不意味着这个内核一定会立即执行, 只是排队到了执行队列中。

- **读取执行结果并释放OpenCL资源**

内核执行完成之后, 需要把数据从GPU拷贝到CPU中, 供主机进一步处理, 所有者写工作完成之后需要释放所有的OpenCL资源。

`clEnqueueReadBuffer()`: 读取设备内存数据到主机内存

`clReleasexxx()`: 释放OpenCL的资源

## 2.2 编程案例

### 2.2.1 查看配置信息

运行CUDA的安装默认路径下的deviceQuery工程 `C:\ProgramData\NVIDIA Corporation\CUDA Samples\v11.0\1_Uutilities\deviceQuery`, 执行结果如下图:

```
C:\ProgramData\NVIDIA Corporation\CUDA Samples\v11.0\bin\win64\Debug\deviceQuery.exe Starting...

CUDA Device Query (Runtime API) version (CUDA static linking)

Detected 1 CUDA Capable device(s)

Device 0: "NVIDIA GeForce RTX 3050 Ti Laptop GPU"
  CUDA Driver Version / Runtime Version      11.4 / 11.0
  CUDA Capability Major/Minor version number: 8.6
  Total amount of global memory:             4096 MBytes (4294967296 bytes)
MapSMtoCores for SM 8.6 is undefined. Default to use 64 Cores/SM
MapSMtoCores for SM 8.6 is undefined. Default to use 64 Cores/SM
  (20) Multiprocessors, ( 64) CUDA Cores/MP: 1280 CUDA Cores
  GPU Max Clock rate:                       1485 MHz (1.49 GHz)
  Memory Clock rate:                        6001 Mhz
  Memory Bus Width:                         128-bit
  L2 Cache Size:                           2097152 bytes
  Maximum Texture Dimension Size (x,y,z)    1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)
  Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048 layers
  Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768), 2048 layers
  Total amount of constant memory:          65536 bytes
  Total amount of shared memory per block:   49152 bytes
  Total number of registers available per block: 65536
  Warp size:                               32
  Maximum number of threads per multiprocessor: 1536
  Maximum number of threads per block:      1024
  Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
  Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
  Maximum memory pitch:                     2147483647 bytes
  Texture alignment:                        512 bytes
  Concurrent copy and kernel execution:     Yes with 5 copy engine(s)
  Run time limit on kernels:                 Yes
  Integrated GPU sharing Host Memory:        No
  Support host page-locked memory mapping:   Yes
  Alignment requirement for Surfaces:        Yes
  Device has ECC support:                    Disabled
  CUDA Device Driver Mode (TCC or WDDM):     WDDM (Windows Display Driver Model)
  Device supports Unified Addressing (UVA):   Yes
  Device supports Managed Memory:            Yes
  Device supports Compute Preemption:        Yes
  Supports Cooperative Kernel Launch:        Yes
  Supports MultiDevice Co-op Kernel Launch:  No
  Device PCI Domain ID / Bus ID / location ID: 0 / 1 / 0
  Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 11.4, CUDA Runtime Version = 11.0, NumDevs = 1
Result = PASS
```

序号	名称	值	解释
1	Detected 1 CUDA Capable device(s)	1	检测到1个可用的NVIDIA显卡设备
2	Device 0: "NVIDIA GeForce RTX 3050 Ti Laptop GPU"	NVIDIA GeForce RTX 3050 Ti Laptop GPU	当前显卡型号为"NVIDIA GeForce RTX 3050 Ti Laptop GPU"
3	CUDA Driver Version / Runtime Version	11.4 / 11.0	CUDA驱动版本
4	<b>CUDA Capability Major/Minor version number</b>	8.6	CUDA设备支持的计算架构版本，即计算能力，该值越大越好
5	<b>Total amount of global memory</b>	4096Mbytes	Global memory全局存储器的大小。使用CUDA RUNTIME API调用函数cudaMalloc后，会消耗GPU设备上的存储空间，合理分配和释放空间避免程序出现crash
6	<b>(3) Multiprocessors, (128) CUDA Cores/MP</b>	1280 CUDA Cores	10个流多处理器（即SM），每个多处理器中包含128个流处理器，共1280个CUDA核
7	GPU Max Clock rate	1485 MHz	GPU最大频率
8	Memory Clock rate	6001 MHz	显存的频率
9	Memory Bus Width	128-bit	内存总线带宽
10	L2 Cache Size	2097152 bytes	第二级缓冲区大小
11	Maximum Texture Dimension Size (x, y, z)	1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)	
12	Maximum Layered 1D Texture Size, (num) layers	1D=(32768), 2048 layers	
13	Maximum Layered 2D Texture Size, (num) layers	2D=(32768, 32768), 2048 layers	
14	<b>Total amount of constant memory</b>	65536 bytes	常量存储器的大小
15	<b>Total amount of shared memory per block</b>	49152 bytes	共享存储器的大小，共享存储器速度比全局存储器快；多处理器上的所有线程块可以同时共享这些存储器
16	Total number of registers available per block	65536	
17	<b>Warp Size</b>	32	Warp，线程束，是SM运行的最基本单位，一个线程束含有32个线程
18	<b>Maximum number of threads per multiprocessor</b>	1536	一个SM中最多有2048个线程，即一个SM中可以有2048/32=64个线程束Warp
19	<b>Maximum number of threads per block</b>	1024	一个线程块最多可用的线程数目
20	<b>Max dimension size of a thread block (x, y, z)</b>	(1024,1024,64)	ThreadIdx.x<=1024,ThreadIdx.y<=1024,ThreadIdx.z<=64Block内三维中各维度的最大值
21	<b>Max dimension size of a grid size (x, y, z)</b>	(2147483647,65535,65535)	Grid内三维中各维度的最大值
22	Maximum memory Pitch	2147483647 bytes	显存访问时对齐时的pitch的最大值
23	Texture alignment	512 bytes	纹理单元访问时对其参数的最大值
24	Concurrent copy and kernel execution	Yes with 1 copy engine(s)	并发复制和内核执行
25	Run time limit on kernels	Yes	内核上的运行时间限制

序号	名称	值	解释
26	Integrated GPU sharing Host Memory	No	集成GPU共享主机内存
27	Support host page-locked memory mapping	Yes	支持主机页面锁定内存映射
28	Alignment requirement for Surfaces	Yes	表面对中要求
29	Device has ECC support	Disabled	设备是否有ECC支持
30	CUDA Device Driver Mode (TCC or WDDM)	WDDM (Windows Display Driver Model)	CUDA设备驱动模式(TCC或WDDM)
31	Device supports Unified Addressing (UVA)	Yes	设备支持统一寻址(UVA)
32	Device supports Managed Memory:	Yes	设备支持托管内存
33	Device supports Compute Preemption:	Yes	设备支持计算抢占
34	Supports Cooperative Kernel Launch:	Yes	支持协作内核启动
35	Supports MultiDevice Co-op Kernel Launch:	No	不支持多设备合作内核启动
36	Device PCI Domain ID / Bus ID / location ID:	0 / 1 / 0	Device PCI Domain ID / Bus ID / location ID

## 2.2.2. 向量运算

### 实验目的

利用OpenCL平台模型进行规模大的向量加法运算

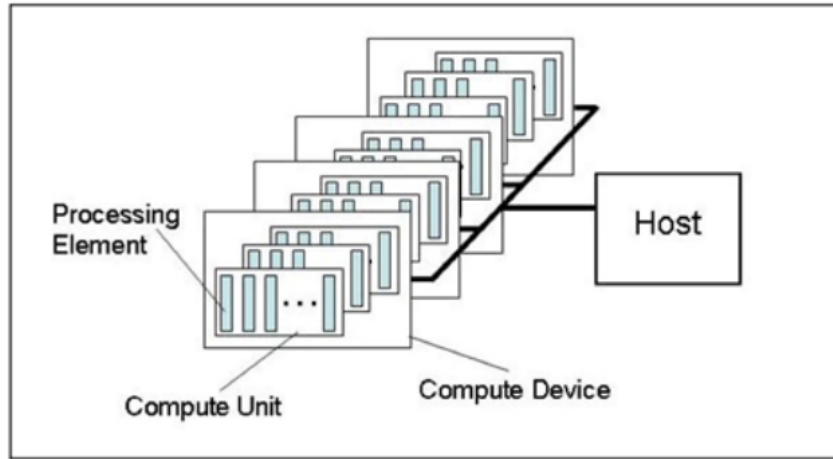
### 实验步骤

#### opengl 代码分析

基于openCL的实验代码分析如下：

- 首先进行定义变量，并初始化。

OpenCL 平台通常包括一个主机 (Host) 和多个 OpenCL 设备 (device)，每个 OpenCL 设备包括一个或多个 CU(compute units)，每个 CU 包括又一个或多个 PE (process element)。每个 PE 都有自己的程序计数器 (PC)。



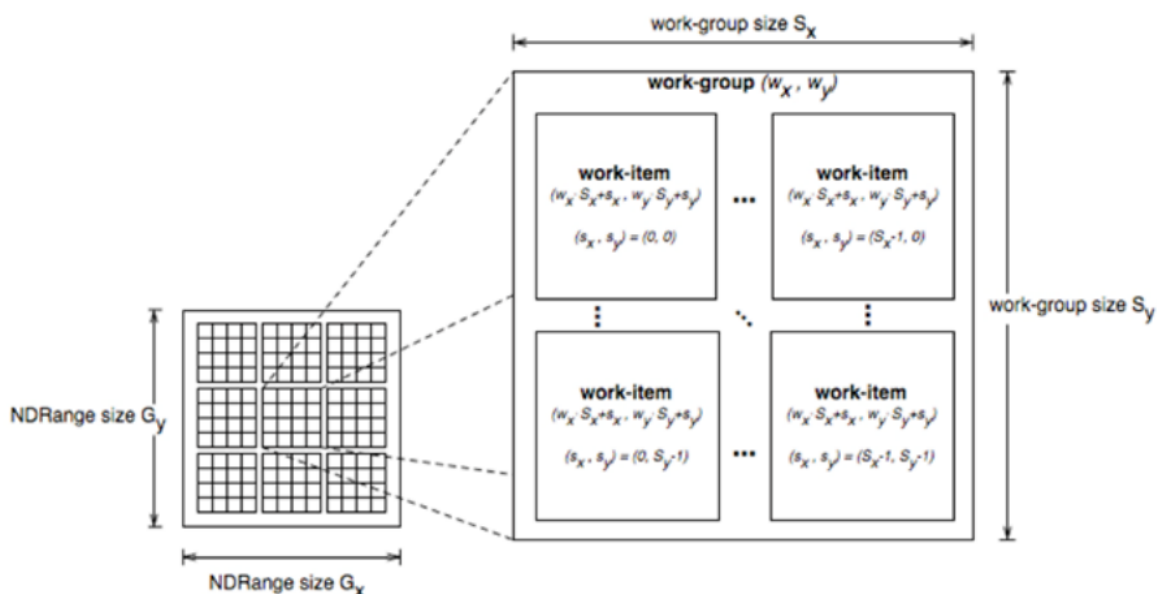
```

1  int i=0;
2  size_t globalSize, localSize;
3  cl_int err;
4  double sum = 0;
5  // 向量的维度大小n
6  int n = 100000;
7
8  //在主机中指向待相加的向量a,b的指针 h_a,h_b
9  double *h_a;
10 double *h_b;
11 //在主机中指向存放求和结果c的指针h_c
12 double *h_c;
13
14 //设备中存放待相加数据的缓冲区
15 cl_mem d_a;
16 cl_mem d_b;
17 //设备中存放求和结果数据的缓冲区
18 cl_mem d_c;
19
20 cl_platform_id platform;           // OpenCL 平台
21 cl_device_id device_id;           // 设备的ID
22 cl_context context;               // 上下文
23 cl_command_queue queue;           // 命令队列
24 cl_program program;               // 程序
25 cl_kernel kernel;                 // 核函数
26
27 // 按照字节的方式统计一个向量占用缓冲区大小
28 size_t bytes = n*sizeof(double);
29
30 // 在主机当中为abc向量开辟内存空间，并用指针h_a,h_b,h_c指向该内存空间
31 h_a = (double*)malloc(bytes);
32 h_b = (double*)malloc(bytes);
33 h_c = (double*)malloc(bytes);
34
35 //初始化主机中的向量数据，使得分别为sin^2和cos^2，目的是求和时结果为1
36 for( i = 0; i < n; i++){
37     h_a[i] = sinf(i)*sinf(i);
38     h_b[i] = cosf(i)*cosf(i);
39 }

```

- 定义globalSize 和 localSize 的daxiao

OpenCL 中的线程结构是可缩放的，Kernel 的每个运行实例称作 WorkItem(也就是线程)，WorkItem 组织在一起称作 WorkGroup，OpenCL 中，每个 workgroup 之间都是相互独立的。通过一个 global id(在索引空间，它是唯一的) 或者 一个 work group id 和一个 work group 内的 local id，就能标定一个 workitem。



```
1 //规定每个work-group 中work-item的数量为64
2 localSize = 64;
3
4 // 计算所占用的work group块的所有空间大小
5 globalSize = ceil(n/(float)localSize)*localSize;// work-group 的数量用ceil函数
   向上取整求得n/localSize，然后再乘上localSize就时globalSize
```

- 选择一个 OpenCL 平台

```
1 err = clGetPlatformIDs(1, &platform, NULL); //第一次调用，得到系统中可使用的平台
   数目，为（Platform）平台对象分配空间
2 err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NULL);
   //第二次调用，是查询所有的平台，选择自己需要的 OpenCL 平台
```

通常调用 clGetPlatformIDs 来获得平台，这个函数要调用 2 次，第一次得到系统中可使用的平台数目，然后为 (Platform) 平台对象分配空间，第二次 调用就是查询所有的平台，选择自己需要的 OpenCL 平台。将函数的返回值赋给err,如果函数执行正确，err的值是 0。

- 设备与context关联

Context 是指管理 OpenCL 对象和资源的上下文环境。为了管 理 OpenCL 程序，下面的一些对象都要和 Context 关联起来：

- 设备 (Devices)：执行 Kernel 程序对象。
- 程序对象 (Program objects)：kernel 程序源代码
- Kernels: 运行在 OpenCL 设备上的函数
- 内存对象 (Memory objects)：设备上存放数据
- 命令队列 (Command queues)：设备的交互机制
- 内存命令 (Memory commands) (用于在主机内存和设备内存之间拷贝数据)
- Kernel 执行 (Kernel execution)
- 同步 (Synchronization)

```
1 context = clCreateContext(0, 1, &device_id, NULL, NULL, &err);
```

`clCreateContext` 函数指定了和 context 关联的一个或多个设备对象，传入的参数 `device_id` 表示绑定的设备编号，如参数值为 `NULL`，则厂商选择的缺省值被使用。

- 命令队列

在 OpenCL 中，命令队列就是主机的请求，在设备上执行的一种机制。命令队列在 device 和 context 之间建立了一个连接。对于不同的设备，它们都有自己的独立的命令队列；命令队列中的命令 (kernel 函数) 可能是同步的，也可能是异步的，它们的执行顺序可以是有序的，也可以是乱序的。

```
1 | queue = clCreateCommandQueue(context, device_id, 0, &err);
```

- OpenCL 程序对象

程序对象就是通过读入 Kernel 函数源代码或二进制文件，然后在指定的设备上编译而产生的 OpenCL 对象。

```
1 | program = clCreateProgramWithSource(context, 1,  
2 |                                     (const char **) & kernelSource, NULL, &err);
```

`clCreateProgramWithSource` 个函数通过源代码 (`kernelSource`)，创建一个程序对象，其中 `counts=1` 指定源代码串的数量，`lengths` 指定源代码串的长度（为 `NULL` 结束的串时，可以省略）此外还需要用 `clBuildProgram` 编写一个从文件中读取源代码串的函数。

```
1 | clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
```

对 context 中的每个设备，这个函数编译、连接源代码对象，产生 device 可以执行的文件

- Kernel 对象

Kernel 就是在程序代码中的一个函数，这个函数能在 OpenCL 设备上执行。一个 Kernel 对象就是 kernel 函数以及其相关的输入参数。Kernel 对象通过程序对象以及指定的函数名字创建。

```
1 | kernel = clCreateKernel(program, "vecAdd", &err);
```

然后为 Kernel 对象设置参数。我们可以在 Kernel 运行后，重新设置参数再次运行。

```
1 | err = clSetKernelArg(kernel, 0, sizeof(cl_mem), &d_a);  
2 | err |= clSetKernelArg(kernel, 1, sizeof(cl_mem), &d_b);  
3 | err |= clSetKernelArg(kernel, 2, sizeof(cl_mem), &d_c);  
4 | err |= clSetKernelArg(kernel, 3, sizeof(unsigned int), &n);
```

调用函数 `clSetKernelArg`，该函数的第二个参数指定该参数为 Kernel 函数中的第几个参数（比如第一个参数为 0，第二个为 1,...）。内存对象和单个的值都可以作为 Kernel 参数。

- OpenCL 内存对象

OpenCL 内存对象就是一些 OpenCL 数据，这些数据一般在设备内存中，能够被拷入也能够被拷出。将 `d_a`, `d_b` 指定为只读的，`d_c` 指定为只写的。函数创建 OpenCL buffer 对象后，会把对应 host buffer 的内容拷贝到 OpenCL buffer 中。

```

1 | d_a = clCreateBuffer(context, CL_MEM_READ_ONLY, bytes, NULL, NULL);
2 | d_b = clCreateBuffer(context, CL_MEM_READ_ONLY, bytes, NULL, NULL);
3 | d_c = clCreateBuffer(context, CL_MEM_WRITE_ONLY, bytes, NULL, NULL);

```

在 Kernel 执行之前，host 中原始输入数据必须显式的传到 device 中，Kernel 执行完后，结果也要从 device 内存中传回到 host 内存中。我们主要通过函数 `clEnqueue{Read/Write}Buffer/Image` 来实现这两种操作。从 host 到 device，我们用 `clEnqueueWrite`，从 device 到 host，我们用 `clEnqueueRead`。`clEnqueueWrite` 命令包括初始化内存对象以及把 host 数据传到 device 内存这两种操作。

```

1 | err = clEnqueueWriteBuffer(queue, d_a, CL_TRUE, 0,
2 |                               bytes, h_a, 0, NULL, NULL);
3 | err |= clEnqueueWriteBuffer(queue, d_b, CL_TRUE, 0,
4 |                               bytes, h_b, 0, NULL, NULL);

```

通过 `clEnqueueWriteBuffer` 函数将主机中的 `h_a` 和 `h_b` 分别传到设备内存中的 `d_a`、`d_b`。

```

1 | clEnqueueReadBuffer(queue, d_c, CL_TRUE, 0,
2 |                      bytes, h_c, 0, NULL, NULL );

```

通过 `clEnqueueWriteBuffer` 函数将设备中计算结果 `d_c` 送回到内存中的 `h_c`。

- kernel 执行

通过函数 `clEnqueueNDRangeKernel` 把 Kernel 放在一个队列里，但不保证它马上执行，OpenCL driver 会管理 队列，调度 Kernel 的执行。注意：每个线程执行的代码都是相同的，但是它们执行数据却是不同的。

```

1 | err = clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &globalSize,
2 |                               &localSize, 0, NULL, NULL);

```

- 主机中计算求和

在本地计算对 `h_c` 指向的缓冲区的值进行累加求得 `sum` 并输出

```

1 | clFinish(queue);
2 | for(i=0; i<n; i++)
3 |     sum += h_c[i];
4 | printf("final result: %f\n", sum/n);

```

- 释放资源

OpenCL 资源都是指针，不使用的时候需要用

`clReleaseMemObject`、`clReleaseProgram`、`clReleaseKernel`、`clReleaseCommandQueue`、`clReleaseContext` 释放掉在设备中的内存空间。而在主机中的内存空间是用 `malloc` 开辟的，因此对应的要用 `free` 函数释放。



```

1 //释放OPenc1的空间
2 clReleaseMemObject(d_a);
3 clReleaseMemObject(d_b);
4 clReleaseMemObject(d_c);
5 clReleaseProgram(program);
6 clReleaseKernel(kernel);
7 clReleaseCommandQueue(queue);
8 clReleaseContext(context);
9
10 //释放主机的内存空间
11 free(h_a);
12 free(h_b);
13 free(h_c);

```

## cuda 代码分析

```

1 #include <stdio.h>
2 #include <cuda.h>
3
4 // 运行在cuda设备上的核函数
5 __global__ void square_array(float *a, int N)
6 {
7     int idx = blockIdx.x * blockDim.x + threadIdx.x; //id编号为块数*每个块的大小
8     if (idx < N) a[idx] = a[idx] * a[idx]; //计算每个id对应的平方和
9 }
10
11 int main()
12 {
13     float *a_h, *a_d; // a_h指向主机数组内存首地址, a_d指向设备数组内存首地址
14     const int N = 10; // 数组中共10个元素
15     size_t size = N * sizeof(float); //数组缓冲区的大小
16     a_h = (float *)malloc(size); // 分配主机上的数组内存
17     cudaMalloc((void **) &a_d, size); // 分配设备上的数组
18     // 初始化主机上的数组
19     for (int i=0; i<N; i++) a_h[i] = (float)i;
20     //把主机上的数组拷贝到设备的内存中
21     cudaMemcpy(a_d, a_h, size, cudaMemcpyHostToDevice);
22     int block_size = 32;
23     int n_blocks = N/block_size + (N%block_size == 0 ? 0:1); //这一步主要是为了向上取整
24     //对于大小为block_size 的n_blocks个块计算平方和
25     square_array <<< n_blocks, block_size >>> (a_d, N);
26     // 将设备上计算的结果拷贝到主机中
27     cudaMemcpy(a_h, a_d, sizeof(float)*N, cudaMemcpyDeviceToHost);
28     // 打印结果
29     for (int i=0; i<N; i++) printf("%d %f\n", i, a_h[i]);
30     // 释放主机和设备的资源
31     free(a_h); cudaFree(a_d);
32     return 0;
33 }

```

## 实验结果

opengl 代码运行结果:

```
Microsoft Visual Studio 调试控制台  
final result: 1.000000
```

cuda 代码运行结果:

```
Microsoft Visual Studio 调试控制台  
0 0.000000  
1 1.000000  
2 4.000000  
3 9.000000  
4 16.000000  
5 25.000000  
6 36.000000  
7 49.000000  
8 64.000000  
9 81.000000
```

## 结果分析

opengl的代码输出结果是 $\sin^2 + \cos^2$ ,其对应的结果应该是1。

cuda的代码计算的是对等差数组 (d=1) 中的每个数求平方, 故其结果分别是0, 1, 4, 9...

## 2.2.3 Pi计算

### 实验目的

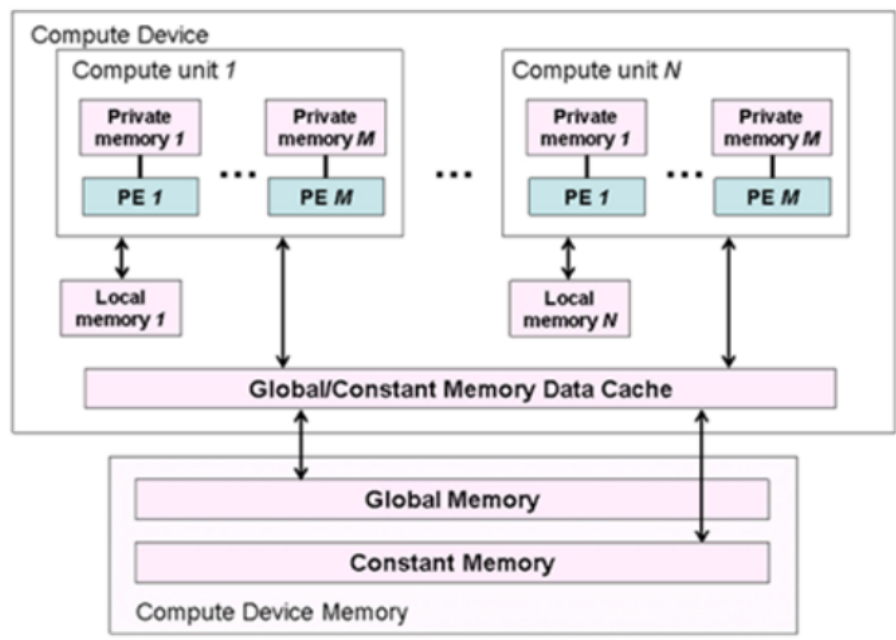
利用OpenCL平台模型进行pi值的计算

### 实验步骤

#### opengl代码分析

- 核函数

每个线程 (workitem) 都有一个 kenerl 函数的实例。每个 Kernel 函数都必须以 `_kernel` 开始, 而且必须返回 `void`。每个输入参数都必须声明使用的内存类型。通过一些 API, 比如 `get_global_id` 之类的得到线程 id。WorkGroup 被映射到硬件的 CU 上执行。



响应的内存对象的地址空间标识符有以下几种：

```

1 | __global //主机中的全局内存空间
2 | __constant //主机中一块只读内存，不允许修改
3 | __local // work_group 共享的内存空间
4 | __private //每个work_item 私有的内存空间

```

- 核函数的定义

```

1 | const char* kernelSource = KERNEL(
2 |     __kernel void Pi(__global float* workGroupBuffer, // 0..NumWorkGroups-1
3 |         __local float* insideworkGroup, // 0..workGroupSize-1
4 |         const uint n, // Total iterations
5 |         const uint chunk) // Chunk size
6 | {
7 |     const uint lid = get_local_id(0); // work_group内每个work_item的编号
8 |     const uint gid = get_global_id(0); //work_group编号
9 |
10 |    const float step = (1.0 / (float)n); //步长
11 |    float partial_sum = 0.0; // 每个小部分的和
12 |
13 |    // 每一个work_item 迭代chunks 次数（即计算chunks个小长方形）
14 |    for (uint i = gid * chunk; i < (gid * chunk) + chunk; i++) {
15 |        float x = step * ((float)i - 0.5);
16 |        partial_sum += 4.0 / (1.0 + x * x);
17 |    }
18 |
19 |    // insideworkGroup[lid]存放着work_item 编号为lid的分组计算的值
20 |    insideworkGroup[lid] = partial_sum;
21 |
22 |    //设置一个barrier，使得所有的work_item 都完成运算
23 |    barrier(CLK_LOCAL_MEM_FENCE);
24 |
25 |    float local_pi = 0;
26 |
27 |    // 然后对work_group 中的所有work_item 计算的值累加
28 |    if (lid == 0) {

```

```

29     const uint length = lid + get_local_size(0); //work item的个数
30     for (uint i = lid; i < length; i++) {
31         local_pi += insideworkGroup[i];
32     }
33     //最终对所有块的值进行规约
34     workGroupBuffer[get_group_id(0)] = local_pi;
35 }
36 }
37

```

- 主函数变量初始化

```

1     int i = 0;
2     float pi; //最终的pi值
3     float* pi_partial; // 每个部分的pi
4     size_t maxworkGroupSize; //work group的最大值
5     cl_int err; //
6     cl_mem memObjects; //
7     int niter, chunks, workGroups; //
8     size_t globalworkSize; //work group的总大小
9     size_t localworkSize; // 一个work group 的大小
10
11     cl_platform_id platform; // OpenCL 平台
12     cl_device_id device_id; // 设备 ID
13     cl_context context; // 上下文
14     cl_command_queue queue; // 命令队列
15     cl_program program; // 程序
16     cl_kernel kernel; // 核
17
18     niter = 262144; //总计算次数（分割成的小长方形数）
19     chunks = 64;

```

- 主函数程序运行逻辑

```

1     //获取平台的ID
2     err = clGetPlatformIDs(1, &platform, NULL);
3     //获取设备的ID
4     err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id, NULL);
5     //获取设备信息
6     clGetDeviceInfo(device_id, CL_DEVICE_MAX_WORK_GROUP_SIZE,
sizeof(size_t),
7         &maxworkGroupSize, NULL);
8     //work group的数量=（总迭代次数/chunk数）/每个work group的最大值
9     workGroups = ceil((float)(niter / maxworkGroupSize / chunks));
10    //为每个部分计算pi的中间结果分配内存空间
11    pi_partial = (float*)malloc(sizeof(float) * workGroups);
12    // 创建上下文
13    context = clCreateContext(0, 1, &device_id, NULL, NULL, &err);
14    // 创建命令队列
15    queue = clCreateCommandQueue(context, device_id, 0, &err);
16    //创建程序
17    program = clCreateProgramWithSource(context, 1,
&kernelSource, NULL, &err);
18    //建立可执行程序
19    err = clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
20

```

```

21     localWorkSize = maxWorkGroupSize;//每个work group 中work item的数量
22     globalWorkSize = niter / chunks; //所有的work group 的大小
23
24     // 创建核函数
25     kernel = clCreateKernel(program, "pi", &err);
26
27     // 创建 OpenCL 设备缓冲 (buffer)
28     memObjects = clCreateBuffer(context, CL_MEM_READ_WRITE,
29         sizeof(float) * workGroups, NULL, &err);
30     //创建 kernel 后, 运行 kernel 之前, 为 kernel 对象设置参数
31     err |= clSetKernelArg(kernel, 0, sizeof(cl_mem), &memObjects);
32     err = clSetKernelArg(kernel, 1, sizeof(float) * maxWorkGroupSize, NULL);
33     err |= clSetKernelArg(kernel, 2, sizeof(unsigned int), &niter);
34     err |= clSetKernelArg(kernel, 3, sizeof(unsigned int), &chunks);
35     //把 kernel 放在一个队列里
36     err = clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &globalWorkSize,
&localWorkSize,
37         0, NULL, NULL);
38     //结束命令队列
39     clFinish(queue);
40     //从设备中读取每个pi_partial 的值
41     err = clEnqueueReadBuffer(queue, memObjects, CL_TRUE, 0,
42         sizeof(float) * workGroups, pi_partial, 0, NULL, NULL);
43     pi = 0;
44     //cpu将每个pi_partial 的值累加在一起
45     for (i = 0; i < workGroups; i++) {
46         pi += pi_partial[i];
47     }
48     //计算最终的pi值
49     pi *= (1.0 / (float)niter);
50     printf("final result: %f\n", pi);
51     //释放OpenCL资源
52     clReleaseMemObject(memObjects);
53     clReleaseProgram(program);
54     clReleaseKernel(kernel);
55     clReleaseCommandQueue(queue);
56     clReleaseContext(context);
57     //释放主机资源
58     free(pi_partial);
59     return 0;
60 }

```

## cuda 代码分析

- 核函数的定义

```

1  __global__ void cal_pi(double *sum, long long nbin, float step, long long
   nthreads, long long nblocks) {
2      long long i;
3      float x;
4      long long idx = blockIdx.x*blockDim.x+threadIdx.x; //每一个work item的编号
   =块编号*块的维度+块内编号
5      for (i=idx; i< nbin; i+=nthreads*nblocks) { //每个线程每隔
6          x = (i+0.5)*step;
7          sum[idx] = sum[idx]+4.0/(1.+x*x);
8      }
9  }

```

- 主函数的编写

```

1  int main()
2  {
3      long long tid;
4      double pi = 0;
5      long long num_steps = 100000000; //总共迭代的次数
6      float step = 1./(float)num_steps; //步长
7      long long size = NUM_THREAD*NUM_BLOCK*sizeof(double); //总线程数
8      clock_t before, after; //时钟标志
9      double *sumHost, *sumDev;
10     sumHost = (double *)malloc(size); //为主机分配内存空间
11     cudaMalloc((void **)&sumDev, size); //为设备分配内存空间
12     // 初始化为全0
13     cudaMemset(sumDev, 0, size);
14     before = clock();
15     // Do calculation on device
16     printf("Before Compute \n\n");
17     dim3 numBlocks(NUM_BLOCK,1,1);
18     dim3 threadsPerBlock(NUM_THREAD,1,1);
19     cal_pi <<<numBlocks, threadsPerBlock>>> (sumDev, (int)num_steps, step,
   NUM_THREAD, NUM_BLOCK); //调用kernel 函数实现pi值的计算
20
21     printf("After Compute \n\n");
22     //将设备的计算结果复制到主机中
23     cudaMemcpy(sumHost, sumDev, size, cudaMemcpyDeviceToHost);
24     printf("After Copy \n\n");
25     for(tid=0; tid<NUM_THREAD*NUM_BLOCK; tid++){ //对每一个线程的计算结果累加
26         pi = pi+sumHost[tid];
27     }
28     pi = pi*step; //最后乘以步长, 得到最终的pi值
29     after = clock();
30     printf("The value of PI is %15.12f\n",pi);
31     printf("The time to calculate PI was %f seconds\n",((float)(after -
   before)/1000.0));
32     free(sumHost); //释放主机内存空间
33     cudaFree(sumDev); //释放设备的内存空间
34     return 0;
35 }

```

## 实验结果

opencl代码结果

```
Microsoft Visual Studio 调试控制台  
final result: 3.141600
```

cuda代码结果

```
Microsoft Visual Studio 调试控制台  
Before Compute  
After Compute  
After Copy  
The value of PI is 3.141592641428  
The time to calculate PI was 0.516000 seconds
```

## 结果分析

opencl和cuda版本都可以正确的计算出pi值的大小，但是cuda版本设置的迭代次数要更多，因此更加准确。

## 3. 总结

openCL 和cuda编程的核心思想是一样的，多核 cpu 就适合基于任务的并行编程，而 GPU 更适应于数据并行编程，即 GPU 上执行一个任务，都是把任务中的数据分配到各个独立的 core 中执行。

无论是计算矩阵求和，或者是计算pi，是把本身在CPU上串行进行的代码改成在GPU上并行计算的。GPU 上，都是轻量级的线程，创建、调度线程的开销比较小，所以我们可以做到把循环完全展开，一个线程处理一个数据。最后把每个线程的运算结果从设备复制到主机中，再在CPU上进行后续的处理,如累加（计算pi）,或者输出数组中的每个值（计算数组和）

## 4. 附录

### cuda 完整代码

#### 1. 计算pi

```
1  #include <stdio.h>  
2  #include <cuda.h>  
3  #include <math.h>  
4  #define NUM_THREAD 1024  
5  #define NUM_BLOCK 1  
6  
7  __global__ void cal_pi(double *sum, long long nbin, float step, long long  
   nthreads, long long nblocks) {  
8  
9     long long i;  
10    float x;  
11    long long idx = blockIdx.x*blockDim.x+threadIdx.x;  
12
```

```

13     for (i=idx; i< nbin; i+=nthreads*nblocks) {
14         x = (i+0.5)*step;
15         sum[idx] = sum[idx]+4.0/(1.+x*x);
16     }
17
18 }
19
20 int main()
21 {
22     long long tid;
23     double pi = 0;
24     long long num_steps = 100000000;
25
26     float step = 1./(float)num_steps;
27     long long size = NUM_THREAD*NUM_BLOCK*sizeof(double);
28     clock_t before, after;
29     double *sumHost, *sumDev;
30     sumHost = (double *)malloc(size);
31     cudaMalloc((void **)&sumDev, size);
32     // Initialize array in device to 0
33     cudaMemset(sumDev, 0, size);
34     before = clock();
35     // Do calculation on device
36     printf("Before Compute \n\n");
37     dim3 numBlocks(NUM_BLOCK,1,1);
38     dim3 threadsPerBlock(NUM_THREAD,1,1);
39     cal_pi <<<numBlocks, threadsPerBlock>>> (sumDev, (int)num_steps, step,
NUM_THREAD, NUM_BLOCK); // call CUDA kernel
40
41     printf("After Compute \n\n");
42     // Retrieve result from device and store it in host array
43     cudaMemcpy(sumHost, sumDev, size, cudaMemcpyDeviceToHost);
44     printf("After Copy \n\n");
45     for(tid=0; tid<NUM_THREAD*NUM_BLOCK; tid++){
46         pi = pi+sumHost[tid];
47     }
48     pi = pi*step;
49     after = clock();
50     printf("The value of PI is %15.12f\n",pi);
51     printf("The time to calculate PI was %f seconds\n",((float)(after -
before)/1000.0));
52     free(sumHost);
53     cudaFree(sumDev);
54     return 0;
55 }

```

## 2. 向量运算

```

1  #include <stdio.h>
2  #include <cuda.h>
3
4  // kernel that executes on the CUDA device
5  __global__ void square_array(float *a, int N)
6  {
7      int idx = blockIdx.x * blockDim.x + threadIdx.x;

```



```

8   if (idx<N) a[idx] = a[idx] * a[idx];
9   }
10
11  int main()
12  {
13      float *a_h, *a_d; // Pointer to host & device arrays
14      const int N = 10; // Number of elements in arrays
15      size_t size = N * sizeof(float);
16      a_h = (float *)malloc(size); // Allocate array on host
17      cudaMalloc((void **) &a_d, size); // Allocate array on device
18      // Initialize host array and copy it to CUDA device
19      for (int i=0; i<N; i++) a_h[i] = (float)i;
20      cudaMemcpy(a_d, a_h, size, cudaMemcpyHostToDevice);
21      // Do calculation on device:
22      int block_size = 32;
23      int n_blocks = N/block_size + (N%block_size == 0 ? 0:1);
24      square_array <<< n_blocks, block_size >>> (a_d, N);
25      // Retrieve result from device and store it in host array
26      cudaMemcpy(a_h, a_d, sizeof(float)*N, cudaMemcpyDeviceToHost);
27      // Print results
28      for (int i=0; i<N; i++) printf("%d %f\n", i, a_h[i]);
29      // Cleanup
30      free(a_h); cudaFree(a_d);
31      return 0;
32  }

```

## openCL 完整代码

### 1. 计算Pi

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <math.h>
5  #include <CL/opencl.h>
6
7
8  // OpenCL kernel. many workGroups compute n iterations
9  #define KERNEL(...) #__VA_ARGS__
10 const char* kernelSource = KERNEL(
11     __kernel void Pi(__global float* workGroupBuffer, // 0..NumWorkGroups-1
12         __local float* insideworkGroup, // 0..workGroupSize-1
13         const uint n, // Total iterations
14         const uint chunk) // Chunk size
15 {
16     const uint lid = get_local_id(0);
17     const uint gid = get_global_id(0);
18
19     const float step = (1.0 / (float)n);
20     float partial_sum = 0.0;
21
22     // Each work-item computes chunk iterations
23     for (uint i = gid * chunk; i < (gid * chunk) + chunk; i++) {
24         float x = step * ((float)i - 0.5);
25         partial_sum += 4.0 / (1.0 + x * x);

```

```

26     }
27
28     // Each work-item stores its partial sum in the workgroup array
29     insideworkGroup[lid] = partial_sum;
30
31     // Synchronize all threads within the workgroup
32     barrier(CLK_LOCAL_MEM_FENCE);
33
34     float local_pi = 0;
35
36     // Only work-item 0 of each workgroup perform the reduction
37     // of that workgroup
38     if (lid == 0) {
39         const uint length = lid + get_local_size(0);
40         for (uint i = lid; i < length; i++) {
41             local_pi += insideworkGroup[i];
42         }
43         // It store the workgroup sum
44         // Final reduction, between block, is done out by CPU
45         workGroupBuffer[get_group_id(0)] = local_pi;
46     }
47 }
48 );
49
50 int main(int argc, char* argv[])
51 {
52     int i = 0;
53     float pi;
54     float* pi_partial;
55     size_t maxWorkGroupSize;
56     cl_int err;
57     cl_mem memObjects;
58     int niter, chunks, workGroups;
59     size_t globalworkSize;
60     size_t localworkSize;
61
62
63     cl_platform_id platform;           // OpenCL platform
64     cl_device_id device_id;           // device ID
65     cl_context context;               // context
66     cl_command_queue queue;           // command queue
67     cl_program program;               // program
68     cl_kernel kernel;                 // kernel
69
70     niter = 262144;
71     chunks = 64;
72
73
74     err = clGetPlatformIDs(1, &platform, NULL);
75
76     // Get ID for the device
77     err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id,
78 NULL);
79     clGetDeviceInfo(device_id, CL_DEVICE_MAX_WORK_GROUP_SIZE,
80 sizeof(size_t),
81 &maxWorkGroupSize, NULL);
82
83     workGroups = ceil((float)(niter / maxWorkGroupSize / chunks));

```

```

82
83     pi_partial = (float*)malloc(sizeof(float) * workGroups);
84
85     // Create a context
86     context = clCreateContext(0, 1, &device_id, NULL, NULL, &err);
87
88     // Create a command queue
89     queue = clCreateCommandQueue(context, device_id, 0, &err);
90
91     // Create the compute program from the source buffer
92
93     program = clCreateProgramWithSource(context, 1,
94         &kernelSource, NULL, &err);
95     // Build the program executable
96     err = clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
97     localWorkSize = maxWorkGroupSize;
98     globalWorkSize = niter / chunks;
99
100    // Create the compute kernel in the program we wish to run
101    kernel = clCreateKernel(program, "pi", &err);
102
103    // Create the input and output arrays in device memory for our
    calculation
104    memObjects = clCreateBuffer(context, CL_MEM_READ_WRITE,
105        sizeof(float) * workGroups, NULL, &err);
106
107    err |= clSetKernelArg(kernel, 0, sizeof(cl_mem), &memObjects);
108    err = clSetKernelArg(kernel, 1, sizeof(float) * maxWorkGroupSize,
    NULL);
109    err |= clSetKernelArg(kernel, 2, sizeof(unsigned int), &niter);
110    err |= clSetKernelArg(kernel, 3, sizeof(unsigned int), &chunks);
111
112    err = clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &globalWorkSize,
    &localWorkSize,
113        0, NULL, NULL);
114    clFinish(queue);
115    err = clEnqueueReadBuffer(queue, memObjects, CL_TRUE, 0,
116        sizeof(float) * workGroups, pi_partial, 0, NULL, NULL);
117    pi = 0;
118
119    for (i = 0; i < workGroups; i++) {
120        pi += pi_partial[i];
121    }
122    pi *= (1.0 / (float)niter);
123    printf("final result: %f\n", pi);
124
125    // release OpenCL resources
126    clReleaseMemObject(memObjects);
127    clReleaseProgram(program);
128    clReleaseKernel(kernel);
129    clReleaseCommandQueue(queue);
130    clReleaseContext(context);
131
132    //release host memory
133    free(pi_partial);
134
135    return 0;
136 }

```

## 2. 向量运算

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <math.h>
4  #include <CL/opencl.h>
5
6  // OpenCL kernel. Each work item takes care of one element of c
7  const char *kernelSource = "\n" \
8  "#pragma OPENCL EXTENSION cl_khr_fp64 : enable\n" \
9  "__kernel void vecAdd( __global double *a,\n" \
10 "                        __global double *b,\n" \
11 "                        __global double *c,\n" \
12 "                        const unsigned int n)\n" \
13 "{\n" \
14 "    //Get our global thread ID\n" \
15 "    int id = get_global_id(0);\n" \
16 "\n" \
17 "    //Make sure we do not go out of bounds\n" \
18 "    if (id < n)\n" \
19 "        c[id] = a[id] + b[id];\n" \
20 "}\n" \
21 "\n" ;
22
23 int main( int argc, char* argv[] )
24 {
25     int i=0;
26     size_t globalSize, localSize;
27     cl_int err;
28     double sum = 0;
29
30     // Length of vectors
31     // unsigned int n = 100000;
32     int n = 100000;
33     // Host input vectors
34     double *h_a;
35     double *h_b;
36     // Host output vector
37     double *h_c;
38
39     // Device input buffers
40     cl_mem d_a;
41     cl_mem d_b;
42     // Device output buffer
43     cl_mem d_c;
44
45     cl_platform_id platform;    // OpenCL platform
46     cl_device_id device_id;    // device ID
47     cl_context context;        // context
48     cl_command_queue queue;    // command queue
49     cl_program program;        // program
50     cl_kernel kernel;          // kernel
51
52     // Size, in bytes, of each vector
```

```

53     size_t bytes = n*sizeof(double);
54
55     // Allocate memory for each vector on host
56     h_a = (double*)malloc(bytes);
57     h_b = (double*)malloc(bytes);
58     h_c = (double*)malloc(bytes);
59
60     // Initialize vectors on host
61
62     for( i = 0; i < n; i++ ){
63         h_a[i] = sinf(i)*sinf(i);
64         h_b[i] = cosf(i)*cosf(i);
65     }
66
67     // size_t globalSize, localSize;
68
69     //cl_int err;
70
71     // Number of work items in each local work group
72     localSize = 64;
73
74     // Number of total work items - localSize must be divisor
75     globalSize = ceil(n/(float)localSize)*localSize;
76
77     // Bind to platform
78     err = clGetPlatformIDs(1, &platform, NULL);
79
80     // Get ID for the device
81     err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device_id,
82 NULL);
83
84     // Create a context
85     context = clCreateContext(0, 1, &device_id, NULL, NULL, &err);
86
87     // Create a command queue
88     queue = clCreateCommandQueue(context, device_id, 0, &err);
89
90     // Create the compute program from the source buffer
91     program = clCreateProgramWithSource(context, 1,
92 (const char **) & kernelSource, NULL, &err);
93
94     // Build the program executable
95     clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
96
97     // Create the compute kernel in the program we wish to run
98     kernel = clCreateKernel(program, "vecAdd", &err);
99
100    // Create the input and output arrays in device memory for our
101    calculation
102    d_a = clCreateBuffer(context, CL_MEM_READ_ONLY, bytes, NULL, NULL);
103    d_b = clCreateBuffer(context, CL_MEM_READ_ONLY, bytes, NULL, NULL);
104    d_c = clCreateBuffer(context, CL_MEM_WRITE_ONLY, bytes, NULL, NULL);
105
106    // Write our data set into the input array in device memory
107    err = clEnqueueWriteBuffer(queue, d_a, CL_TRUE, 0,
108 bytes, h_a, 0, NULL, NULL);
109    err |= clEnqueueWriteBuffer(queue, d_b, CL_TRUE, 0,
110 bytes, h_b, 0, NULL, NULL);

```

```

109
110 // Set the arguments to our compute kernel
111 err = clSetKernelArg(kernel, 0, sizeof(cl_mem), &d_a);
112 err |= clSetKernelArg(kernel, 1, sizeof(cl_mem), &d_b);
113 err |= clSetKernelArg(kernel, 2, sizeof(cl_mem), &d_c);
114 err |= clSetKernelArg(kernel, 3, sizeof(unsigned int), &n);
115
116 // Execute the kernel over the entire range of the data set
117 err = clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &globalSize,
118 &localSize,
119                                     0, NULL,
120 NULL);
121
122 // wait for the command queue to get serviced before reading back
123 results
124 clFinish(queue);
125
126 // Read the results from the device
127 clEnqueueReadBuffer(queue, d_c, CL_TRUE, 0,
128                     bytes, h_c, 0, NULL, NULL );
129
130 //Sum up vector c and print result divided by n, this should equal 1
131 within error
132 //double sum = 0;
133 for(i=0; i<n; i++)
134     sum += h_c[i];
135 printf("final result: %f\n", sum/n);
136
137 // release OpenCL resources
138 clReleaseMemObject(d_a);
139 clReleaseMemObject(d_b);
140 clReleaseMemObject(d_c);
141 clReleaseProgram(program);
142 clReleaseKernel(kernel);
143 clReleaseCommandQueue(queue);
144 clReleaseContext(context);
145
146 //release host memory
147 free(h_a);
148 free(h_b);
149 free(h_c);
150
151 return 0;
152 }

```