



南開大學

Nankai University

计算机学院  
并行程序设计实验报告

基于 Intel DevCloud 平台学习 SYCL  
课程学习报告

姓名：蔡沅宏

学号：2213897

专业：计算机科学与技术

2024 年 6 月 15 日

# 目录

<b>1 DevCloud 平台学习基本要求</b>	<b>2</b>
1.1 oneAPI_Intro . . . . .	2
1.1.1 学习目标 . . . . .	2
1.1.2 学习内容 . . . . .	2
1.1.3 实践 . . . . .	3
1.2 SYCL_Program_Structure . . . . .	5
1.2.1 学习目标 . . . . .	5
1.2.2 学习内容 . . . . .	5
1.2.3 实践 . . . . .	8
1.3 unified_Shared_Memory . . . . .	10
1.3.1 学习目标 . . . . .	10
1.3.2 学习内容 . . . . .	10
1.3.3 实践 . . . . .	11
<b>2 DevCloud 平台学习进阶要求</b>	<b>14</b>
2.1 Sub_Groups . . . . .	14
2.1.1 学习目标 . . . . .	14
2.1.2 学习内容 . . . . .	14
2.1.3 实践 . . . . .	15
2.2 Intel_Advisor . . . . .	17
2.2.1 学习目标 . . . . .	17
2.2.2 学习内容 . . . . .	18

# 1 DevCloud 平台学习基本要求

## 1.1 oneAPI\_Intro

### 1.1.1 学习目标

- 解释 oneAPI 编程模型如何解决异构环境中的编程挑战
- 使用 oneAPI 项目启用 workflow
- 了解 SYCL 语言和编程模型
- 熟悉在整个课程中使用 Jupyter 笔记本进行培训

### 1.1.2 学习内容

目前，在数据中心领域，专用工作负载的增长显著。每种类型的数据中心硬件通常需要使用不同的语言和库进行编程，因为没有通用的编程语言或 API，这需要维护独立的代码库。开发者必须学习一整套不同的工具，因为各平台的工具支持不一致。

为解决以上问题，oneAPI 旨在提供统一的编程模型，以简化跨多种架构的开发，如图1.1所示。它包括统一且简化的语言和库，用于表达并行性。

**oneAPI** oneAPI 编程模型提供了一个全面而统一的开发人员工具组合，可以用于跨硬件目标的开发，它包括跨多个工作负载的一系列性能库。这些库包括为每个目标架构自定义编码的函数，因此相同函数调用可在支持的体系结构中提供优化的性能。DPC++ 基于行业标准和开放规范，以鼓励生态系统协作和创新。

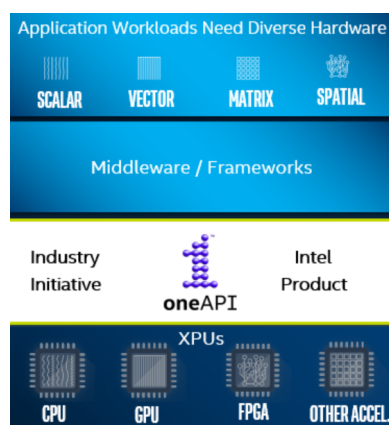


图 1.1: OneAPI 实现跨架构开发

接下来，平台用一个小程序为我们展示 SYCL 编程，并以此为基础帮助我们在 DevCloud 环境中入门，同时介绍了用于编辑和保存代码的 Jupyter 笔记本环境。这部分将在实践模块中细讲。

然后平台介绍了 SYCL，SYCL 构建在 OpenCL 之上，提供跨平台抽象层，使得异构处理器代码可以用 C++ “单源” 风格编写。而与 OpenCL 不同，SYCL 包含模板和 lambda 函数，使高层次应用软件可以干净地编码，同时优化内核代码的加速。

## oneAPI 编程模型

1. 平台模型：oneAPI 的平台模型基于 SYCL 平台模型，主要涉及主机和设备之间的协调与控制，如图2(a)所示。主机通常是基于 CPU 的系统，负责控制一个或多个设备执行的计算工作。设备则是加速器，包含专用的计算资源，能够高效执行操作，通常比主机的 CPU 更有效。每个设备包含一个或多个计算单元，这些单元能够并行执行多个操作，每个计算单元内部又包含一个或多个处理元素，充当独立的计算引擎。这种关系可以通过图示直观地描述出来，其中一个主机可以与一台或多台设备通信，每个设备可以包含一个或多个计算单元，每个计算单元又可以包含一个或多个处理元素。
2. 执行模型：执行模型基于 SYCL 执行模型，主要定义了代码（内核）在设备上执行以及与主机交互的方式。主机执行模型通过命令组协调主机和设备之间的执行和数据管理，其中命令组是内核调用和访问器等命令的分组，被提交到队列中执行。访问器是内存模型的一部分，同时传达执行的排序要求。程序使用队列来声明和实例化，并可以按照程序可控制的顺序或无序策略执行。设备执行模型指定了在加速器上完成计算的方式，将计算范围分布在 ND 范围、工作组、子组和工作项的层次结构中。实际内核代码表示为一个工作项执行的工作，其并行性由外部控制，工作的数量和分布由 ND 范围和工作组的大小规范控制，如图2(b)。
3. 内存型号：oneAPI 的内存模型基于 SYCL 内存模型，定义了主机和设备之间的内存交互和管理，如图2(c)所示。内存可以驻留在主机或设备上，通过缓冲区和图像两种类型的内存对象进行指定和访问。访问器指定了访问位置和模式。例如，可以通过创建缓冲区对象使主机分配的内存能够与设备通信，访问器对象定义了对缓冲区的具体访问方式。
4. 内核编程模型：oneAPI 的内核编程模型基于 SYCL 内核编程模型，支持显式并行性，程序员可以明确决定在主机和设备上执行哪些代码。内核代码在加速器上执行。采用 oneAPI 编程模型的程序支持单源代码，即主机代码和设备代码可以位于同一源文件中。不过，主机代码和设备代码在语言一致性和功能上有所不同，SYCL 规范详细定义了主机代码和设备代码所需的语言功能。

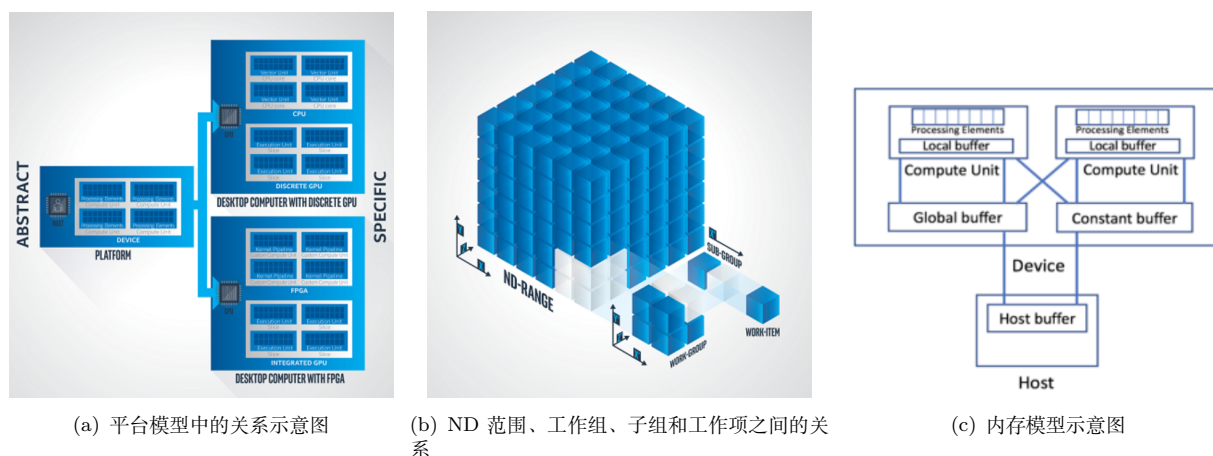


图 1.2: oneAPI 编程模型相关示意图

### 1.1.3 实践

本部分课程的实践只有一个 `simple.cpp`，并且不要求进行修改或编写，代码主要是为我们展示 SYCL 编程的格式，让我们熟悉平台的使用。

## simple.cpp

```

1  #include <sycl/sycl.hpp>
2  using namespace sycl;
3  static const int N = 16;
4  int main(){
5      ///  

6      queue q;
7      std::cout << "Device: " << q.get_device().get_info<info::device::name>() <<
8          "\n";
9
10     ///  

11     ///  

12     ///  

13     ///  

14     ///  

15     ///  

16     ///  

17     ///  

18     ///  

19     ///  

20     ///  

21     ///  

22     ///  

23     ///  

24     ///  

25     }

```

我们观察这段代码，代码实际上是在设备上并行计算与内存管理操作。首先程序创建一个 SYCL 队列，选择默认计算设备。用统一共享内存（USM）为大小为 16 的整数数组分配内存，然后在设备上并行地将数组中的每个元素乘以 2。再输出计算后的数组元素并释放内存。

运行后可以得到如下图1.3结果：

### 2.1.1.构建和运行

选择下面的单元格，然后单击“运行”以编译并执行上面的代码：

```

[2]: ! chmod 755 q; chmod 755 run_simple.sh;if [ -x "$(command -v qsub)" ]; then ./q run_simple.sh; else ./run_simple.sh; fi
## ua08be3491b6b99332287772d458427 is compiling SYCL_Essentials Module1 -- oneAPI Intro sample - 1 of 1 simple.cpp
Device: Intel(R) Data Center GPU Max 1100
0
2
4
6
8
10
12
14
16
18
20
22
24
26
28
30

```

图 1.3: simple.cpp 程序运行截图

## 1.2 SYCL\_Program\_Structure

### 1.2.1 学习目标

- 解释 SYCL 基础类
- 使用设备选择卸载内核工作负载
- 决定何时使用基本并行内核和 ND 范围内核
- 在 SYCL 程序中使用统一共享内存或缓冲区访问器内存模型
- 通过动手实验练习构建示例 SYCL 应用程序

### 1.2.2 学习内容

本章的课程内容主要介绍 SYCL 的相关编程结构知识。

SYCL 当中有一些重要的类，来方便我们在代码中进行调用使用，以下对其进行总结整理，并未展示出课程中所示所有类，总结整理如下：

1. Device Selector（设备选择器）：这些类允许根据用户提供的启发式方法来选择特定的设备来执行内核。以下是标准设备选择器的使用示例：

```
1 queue q(gpu_selector_v); // GPU 选择器
2 // queue q(cpu_selector_v); // CPU 选择器
3 // queue q(accelerator_selector_v); // 加速器选择器
4 // queue q(default_selector_v); // 默认选择器
5 // queue q; // 默认选择器（默认选择器等同于 default_selector_v）
```

输出选择的设备名称：

```
1 std::cout << "Device: " << q.get_device().get_info<info::device::name>() << "\n";
```

2. Queue（队列）：队列用于提交命令组以在 SYCL 运行时执行。每个队列与一个设备相关联，可以在异构系统中使用多个队列。

```
1 q.submit([&](handler& h) {
2     // 命令组代码
3 });
```

3. Kernel（内核）：内核类封装了在设备上执行代码所需的方法和数据。内核对象不需要用户显式构造，而是在调用诸如 parallel\_for 的内核调度函数时构造。

```
1 q.submit([&](handler& h) {
2     h.parallel_for(range<1>(N), [=](id<1> i) {
3         A[i] = B[i] + C[i];
4     });
5 });
```





在以上所示的使用多种操作实现并行后，接下来自然会遇到同步的问题，教程中给出了使用主机访问器和缓冲区销毁进行同步的方法。以下是使用主机访问器的代码示例。

#### host\_accessor\_sample.cpp

```

1  %%writefile lab/host_accessor_sample.cpp
2  #include <sycl/sycl.hpp>
3  using namespace sycl;
4
5  int main() {
6      constexpr int N = 16;
7      auto R = range<1>(N);
8      std::vector<int> v(N, 10);
9      queue q;
10     // Buffer takes ownership of the data stored in vector.
11     buffer buf(v);
12     q.submit([&](handler& h) {
13         accessor a(buf, h);
14         h.parallel_for(R, [=](auto i) { a[i] -= 2; });
15     });
16     // Creating host accessor is a blocking call and will only return after all
17     // enqueued SYCL kernels that modify the same buffer in any queue completes
18     // execution and the data is available to the host via this host accessor.
19     host_accessor b(buf, read_only);
20     for (int i = 0; i < N; i++) std::cout << b[i] << " ";
21     return 0;
22 }
```

在运用 HostAccessor 实现数据同步的代码中，首先初始化了一个包含 16 个元素的向量，每个元素的值为 10。然后创建了一个 SYCL 队列和一个与向量关联的 buffer 对象。在提交的命令组中，代码通过 accessor 对象访问 buffer，并在并行 for 循环中将每个元素的值减去 2。接着，代码创建了一个 hostaccessor 对象来同步数据回主机。值得注意的是，Host accessor 的创建是一个阻塞调用，只有在所有修改同一 buffer 的 SYCL 内核执行完毕后，数据才会通过 host accessor 可用。最后，代码通过 host accessor 读取并打印 buffer 中的每个元素，展示了数据同步后的结果。

随后教程介绍了用缓冲区销毁 (BufferDestruction) 来实现数据同步 (相关数据与上述代码中的设置是相同的，只是实现数据同步的方式不同)。缓冲区的创建发生在一个单独的函数作用域内。当执行超出这个函数作用域时，缓冲区的析构函数被调用，释放数据的所有权并将数据复制回主机内存。

虽然以上实现数据同步的方式不同，但是两种方式产生的结果是相同的，如图1.7所示。

```

[4]: [1] cheod 755 q; cheod 755 run_host_accessor.sh; if [ -x "$i{command -v qsub}" ]; then ./q_run_host_accessor.sh; else ./run_host_accessor.sh; fi
## uae08e3493d0993322877726458427 is compiling SYCL_Essentials Module2 -- SYCL Program Structure sample - 3 of 7 host_accessor_sample.cpp
*****
```

(a) 主机访问器

```

[4]: [1] cheod 755 q; cheod 755 run_buffer_destruction.sh; if [ -x "$i{command -v qsub}" ]; then ./q_run_buffer_destruction.sh; else ./run_buffer_destruction.sh; fi
## uae08e3493d0993322877726458427 is compiling SYCL_Essentials Module2 -- SYCL Program Structure sample - 4 of 7 buffer_destruction.cpp
*****
```

(b) 缓冲区销毁

图 1.6: 使用主机访问器和缓冲区销毁同步两种方式的结果

同时，SYCL 中还允许开发者使用自定义设备选择器来指定程序应该在何种类型的设备上运行。自定义设备选择器的示例包括根据供应商名称选择设备、选择特定 GPU 设备名称以及基于设备类型设定优先级。



1. 根据供应商名称选择设备：可以编写自定义选择器函数，返回 1 以选择供应商名称为”Intel” 的设备，或返回 0 以允许其他设备作为备选选择。

2. 选择特定 GPU 设备名称：通过检查设备是否为 GPU 并且名称包含”Intel” 来选择特定的 GPU 设备。

3. 基于设备类型设置优先级：可以根据设备类型设定优先级，例如优先选择 Xeon 设备，其次是如何 GPU，再次是如何 CPU。

自定义设备选择器允许开发者根据应用程序的特定需求优化设备选择，从而实现更高效的计算资源利用。

### 1.2.3 实践

最后，在本章课程中给出了一个向量添加的实验练习，要求使用 SYCL Buffer 和 Accessor 概念完成以下编码练习：

- 该代码在主机上共初始化了三个向量
- 内核代码将向量 vector1 的每个元素增加 1
- 创建一个新的向量 vector2，并初始化为值 20
- 为上述第二个向量创建 SYCL 缓冲区
- 在内核代码中，为第二个向量缓冲区创建第二个访问器
- 将向量增量修改为向量相加，将 vector1 的值加到 vector2 上

以下是我补充完成的代码：

vector\_add.cpp

```
1 %%writefile lab/vector_add.cpp
2 #include <sycl/sycl.hpp>
3 using namespace sycl;
4 int main() {
5     const int N = 256;
6
7     ///  
Initialize a vector and print values  
8     std::vector<int> vector1(N, 10);  
9     std::cout<<"\nInput Vector1: ";  
10    for (int i = 0; i < N; i++) std::cout << vector1[i] << " ";  
11
12    ///  
STEP 1 : Create second vector, initialize to 20 and print values  
13    ///  
YOUR CODE GOES HERE  
14    std::vector<int> vector2(N, 20);  
15    std::cout << "\nInput Vector2: ";  
16    for (int i = 0; i < N; i++) std::cout << vector2[i] << " ";  
17
18    ///  
Create Buffer  
19    buffer vector1_buffer(vector1);  
20
21    ///  
STEP 2 : Create buffer for second vector
```

```

22 // # YOUR CODE GOES HERE
23 buffer vector2_buffer(vector2);
24
25 // # Submit task to add vector
26 queue q;
27 q.submit([&](handler &h) {
28     // # Create accessor for vector1_buffer
29     accessor vector1_accessor (vector1_buffer, h);
30
31     // # STEP 3 – add second accessor for second buffer
32     // # YOUR CODE GOES HERE
33     accessor vector2_accessor(vector2_buffer, h, read_only);
34
35     h.parallel_for(range<1>(N), [=](id<1> index) {
36
37         // # STEP 4 : Modify the code below to add the second vector to first one
38         vector1_accessor[index] += vector2_accessor[index];
39
40     });
41 });
42
43 // # Create a host accessor to copy data from device to host
44 host_accessor h_a(vector1_buffer, read_only);
45
46 // # Print Output values
47 std::cout<<"\nOutput Values: ";
48 for (int i = 0; i < N; i++) std::cout<< vector1[i] << " ";
49 std::cout<<"\n";
50
51 return 0;
52 }

```

程序运行结果如下图1.7:

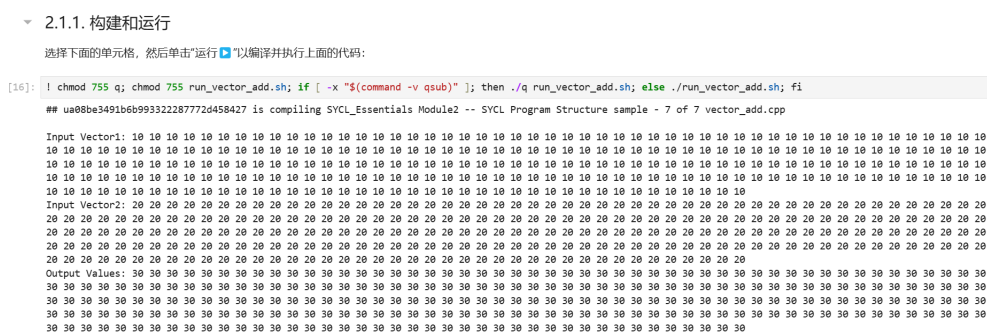


图 1.7: 向量添加实验练习运行结果

## 1.3 unified\_Shared\_Memory

### 1.3.1 学习目标

- 使用新的 SYCL2020 功能（如统一共享内存）来简化编程。
- 了解使用 USM 移动内存的隐式和显式方式。
- 以最优方式解决内核任务之间的数据依赖性。

### 1.3.2 学习内容

本小节主要介绍了有关 USM 模型的使用。Unified Shared Memory (USM) 是一种编程模型，允许在同一地址空间中共享内存，无论是在 CPU 还是 GPU 上。这简化了内存管理，使得数据在设备之间可以无需显式拷贝而共享。

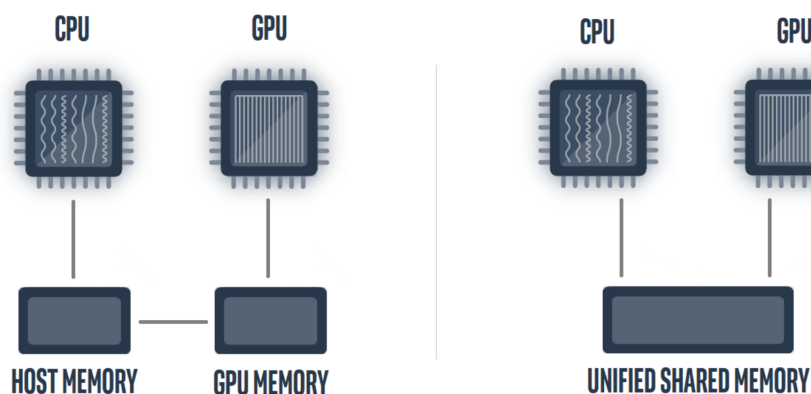


图 1.8: 不带 USM 和带 USM 的内存的开发人员视图

在 USM 编程时需要注意在 USM 中的数据依赖性。因为任务是异步执行的，多个任务可以同时执行，所以当使用 USM 时，任务之间的依赖关系必须通过事件来指定，因为任务是异步执行的，多个任务可以同时执行。

为解决数据依赖性，USM 提供了以下几种解决数据依赖性的不同选项：

1. `wait()` 方法：在任务之间使用 `q.wait()` 来等待依赖完成，但会阻塞主机执行。
2. `in_order` 队列属性：使用 `property::queue::in_order()` 属性创建的队列可以序列化所有内核任务，确保执行顺序不重叠。
3. `depends_on` 方法：在命令组中使用 `h.depends_on(e)` 方法来指定依赖事件 `e`，确保在任务开始之前等待事件完成。

此外，还可以使用简化的方式指定依赖性，通过在 `parallel_for` 中传递额外参数（事件对象）来简化指定依赖性，例如 `q.parallel_for(range < 1 > (N), e, [=](id < 1 > i) data[i] += 3;);`。

这些方法和属性使得在使用 USM 时可以有效地管理和控制任务之间的数据依赖关系，以避免非期望的结果。



- 将 data1 复制回主机并验证结果

根据以上实验要求，我补充完成了以下代码：

#### usm\_lab.cpp

```

1 %%writefile lab/usm_lab.cpp
2 #include <sycl/sycl.hpp>
3 using namespace sycl;
4 static const int N = 1024;
5 int main() {
6     queue q;
7     std::cout << "Device : " << q.get_device().get_info<info::device::name>() << "\n";
8
9     //intialize 2 arrays on host
10    int *data1 = static_cast<int *>(malloc(N * sizeof(int)));
11    int *data2 = static_cast<int *>(malloc(N * sizeof(int)));
12    for (int i = 0; i < N; i++) {
13        data1[i] = 25;
14        data2[i] = 49;
15    }
16
17    //STEP 1 : Create USM device allocation for data1 and data2
18    //YOUR CODE GOES HERE
19    int *d_data1 = malloc_device<int>(N, q);
20    int *d_data2 = malloc_device<int>(N, q);
21
22    //STEP 2 : Copy data1 and data2 to USM device allocation
23    //YOUR CODE GOES HERE
24    q.memcpy(d_data1, data1, N * sizeof(int)).wait();
25    q.memcpy(d_data2, data2, N * sizeof(int)).wait();
26
27    //STEP 3 : Write kernel code to update data1 on device with square of its value
28    auto e1 = q.parallel_for(N, [=](auto i) {
29
30        //YOUR CODE GOES HERE
31        d_data1[i] = sqrt(d_data1[i]);
32    });
33
34    //STEP 4 : Write kernel code to update data2 on device with square of its value
35    auto e2 = q.parallel_for(N, [=](auto i) {
36
37        //YOUR CODE GOES HERE
38        d_data2[i] = sqrt(d_data2[i]);
39    });
40
41    //STEP 5 : Write kernel code to add data2 on device to data1
42    auto e3 = q.parallel_for(N, {e1, e2}, [=](auto i) {
43
44        //YOUR CODE GOES HERE

```

```

45     d_data1[i] += d_data2[i];
46 });
47
48 // # STEP 6 : Copy data1 on device to host
49 // # YOUR CODE GOES HERE
50 q.memcpy(data1, d_data1, N * sizeof(int)).wait();
51
52 // # verify results
53 int fail = 0;
54 for (int i = 0; i < N; i++) if(data1[i] != 12) {fail = 1; break;}
55 if(fail == 1) std::cout << " FAIL"; else std::cout << " PASS";
56 std::cout << "\n";
57
58 // # STEP 7 : Free USM device allocations
59 // # YOUR CODE GOES HERE
60 free(d_data1, q);
61 free(d_data2, q);
62
63 // # STEP 8 : Add event based kernel dependency for the Steps 2 - 6
64 free(data1);
65 free(data2);
66
67 return 0;
68 }

```

从代码中我，我们可以看到，它创建了三个并行计算任务（内核）。第一个内核任务将 `d_data1` 中每个元素的值更新为其平方根，因为其本来值为 25，所以更新后为 5；第二个内核任务将 `d_data2` 中每个元素的值更新为其平方根，其值本来为 49，所以更新之后为 7；第三个内核任务将 `d_data2` 中每个元素的值加到 `d_data1` 中对应的元素上，则 `d_data1` 中所有元素都为 12，之后将 `d_data1` 从设备复制回主机后，那么 `data1` 数组中所有元素都为 12。那么最后，验证 `data1` 数组中的每个元素是否都等于 12 时是成立的，所以如果代码没有编写错误的话，输出结果应该为 PASS。

程序运行后得到如下图1.11结果：

1.10.0.1. 构建和运行

选择下面的单元格，然后单击“运行”以编译并执行代码：

```

[10]: ! chmod 755 run_usm_lab.sh; if [ -x "$(command -v qsub)" ]; then ./q run_usm_lab.sh; else ./run_usm_lab.sh; fi
      ## ua08be3491b6b99332287772d458427 is compiling SYCL_Essentials Module3 -- SYCL Unified Shared Memory - 5 of 5 usm_lab.cpp
      Device : Intel(R) Data Center GPU Max 1100
      PASS

```

图 1.11: 统一共享内存实验练习程序运行结果

如上图所示，实验结果符合预期，程序编写正确。

## 2 DevCloud 平台学习进阶要求

### 2.1 Sub\_Groups

#### 2.1.1 学习目标

- 了解在 SYCL 中使用子组的优势
- 利用子组算法提高性能和生产力
- 使用子组随机播放操作可避免显式内存操作

#### 2.1.2 学习内容

首先，在许多现代硬件平台上，工作组中工作项的子集是同时执行的，或者具有额外的计划保证。这些工作项子集称为子组。利用子组将有助于将执行映射到低级硬件，并可能有助于实现更高的性能。

也就是说，通过子组，我们可以将任务划分为多个子集将其映射到不同硬件上以获得更快的效率。

同时，课程中提到：使用 ND\_RANGE 内核进行并行执行有助于将工作项分组映射到硬件资源。这有助于为性能进行应用调优。

ND-range 内核的执行范围如图2.12所示，分为工作组、子组和工作项。

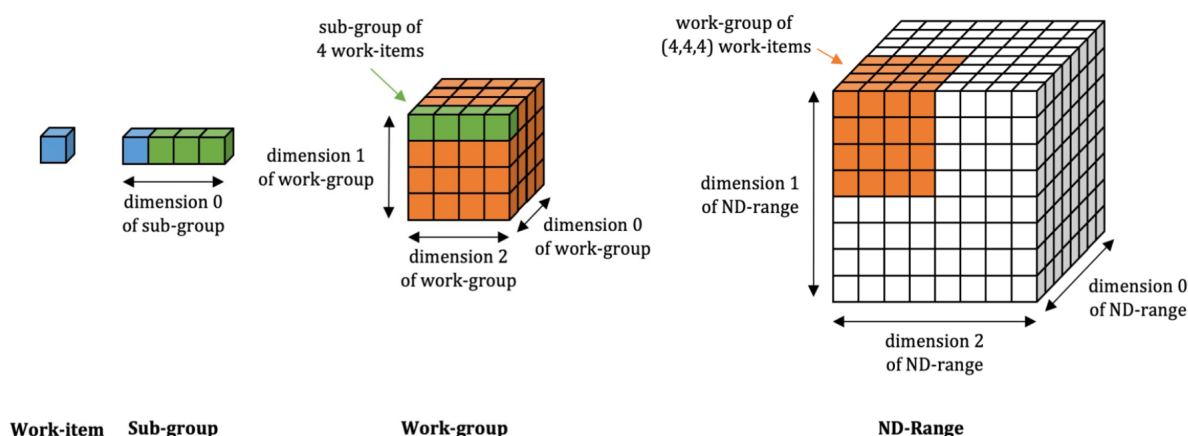


图 2.12: ND-range 内核的执行范围

分工如下所示：

工作项 (Work-item)	表示核函数的各个实例。
工作组 (Work-group)	整个迭代空间被划分为更小的组，称为工作组，工作组内的工作项在硬件上的单个计算单元上调度执行。
子组 (Subgroup)	工作组内的工作项的一个子集，同时执行，可能映射到向量硬件上。

接下来，教程介绍使用子组时可以用到的相关函数和成员变量，总结如下：

#### 1. 子组信息：

- 可以通过查询子组句柄获取详细信息，如子组中工作项的索引 (`get_local_id()`)、子组的大小 (`get_local_range()`)、子组在父工作组中的索引 (`get_group_id()`) 以及父工作组内子组的数量 (`get_group_range()`)。



## 2. 子组大小:

- 为了优化性能, 可以设置特定的子组大小, 如 8、16 或 32。编译器通常会选择最佳的大小, 但也可以使用 `reqd_sub_group_size(S)` 属性显式设置。

## 3. 子组函数和算法:

- 这些函数和算法作用于子组内的工作项, 例如 `select_by_group`、`shift_group_left`、`shift_group_right`、`permute_group_by_xor` 等。它们可以执行洗牌、减少、扫描和投票等操作, 用于优化并行计算。

## 4. 子组洗牌 (Shuffle):

- 允许子组内的工作项直接通信, 而无需显式内存操作。例如, 使用 `permute_group_by_xor` 可以交换工作项的值, 从而提高内核执行效率, 减少对全局内存的访问。

为了让我们更好地理解 Shuffle 的内涵, 教程中给出了如下的代码示例:

## 子组随机操作示例

```
1  h.parallel_for(nd_range<1>(N,B), [=](nd_item<1> item){
2      auto sg = item.get_sub_group();
3      auto i = item.get_global_id(0);
4      /* Shuffles */
5      //data[i] = select_by_group(sg, data[i], 2);
6      //data[i] = shift_group_left(sg, data[i], 1);
7      //data[i] = shift_group_right(sg, data[i], 1);
8      data[i] = permute_group_by_xor(sg, data[i], 1);
9  });
```

以上代码中利用 `permute_group_by_xor` 交换两个工作项的值, 如图2.13所示:

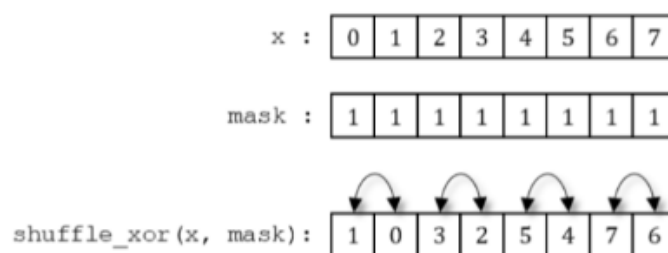


图 2.13: 两个工作项值的交换示意图

## 2.1.3 实践

本章课程实践部分需要我们通过使用子组中的 `reduce` 方法实现从 1 到 1024 的求和。

题目要求我们将把内核任务卸载到计算每个子组中所有项的总和, 并保存在新数组 `sg_data` 中。子组大小要设置为  $S=32$ , 这也就使得 `sg_data` 数组的大小为  $N/32$ 。然后我们需要运用第三章的内容创建 USM 共享分配, 用于 `data` 和 `sg_data`。同时创建一个 `nd-range` 内核任务, 固定子组大小为 32。

在内核任务中，我需要按要求使用 `reduce_over_group` 函数计算子组总和。并将每个子组的总和保存到 `sg_data` 数组中。最后将 `sg_data` 的所有元素相加以获得最终总和。

实现代码如下：

#### sub\_group\_lab.cpp

```

1  #include <CL/sycl.hpp>
2  using namespace cl::sycl;
3
4  static constexpr size_t N = 1024; // global size
5  static constexpr size_t B = 256;  // work-group size
6  static constexpr size_t S = 32;   // sub-group size
7
8  int main() {
9      queue q;
10
11     std::cout << "Device: " << q.get_device().get_info<info::device::name>() << "\n";
12
13     int* data = malloc_shared<int>(N, q);
14     int* sg_data = malloc_shared<int>(N/S, q);
15
16     for (int i = 0; i < N; i++)
17         data[i] = i;
18
19     for (int i = 0; i < N; i++)
20         std::cout << data[i] << " ";
21
22     std::cout << "\n\n";
23
24     q.parallel_for(nd_range<1>(N, B), [=](nd_item<1> item)
25         [[intel::reqd_sub_group_size(S)]] {
26         auto sg = item.get_sub_group();
27         auto i = item.get_global_id(0);
28         int result = reduce_over_group(sg, data[i], plus<>());
29         if (sg.leader()) {
30             sg_data[item.get_group(0) * (B / S) + sg.get_group_id()[0]] = result;
31         }
32     }).wait();
33
34     for (int i = 0; i < N/S; i++)
35         std::cout << sg_data[i] << " ";
36
37     std::cout << "\n";
38
39     int sum = 0;
40     for (int i = 0; i < N/S; i++) {
41         sum += sg_data[i];
42     }
43

```

```

44     std::cout << "\nSum = " << sum << "\n";
45
46     free(data, q);
47     free(sg_data, q);
48
49     return 0;
50 }

```

根据以上代码和预期完成功能，我们可以预测程序运行结果。如果程序逻辑正确，那么应该能正常分出 32 个子组，并下放至硬件进行计算，保存至 sg\_data[] 中，那么执行代码后，应该会输出 sg\_data 数组的内容，即每个子组中元素的总和。最后输出计算得到的所有子组和的总和 Sum。

运行结果如下图2.14所示：

### 1.11.0.1.构建和运行

选择下面的单元格，然后单击“运行”以编译并执行代码：

```

! chmod 755 run_sub_group_lab.sh; if [ -x "$(command -v qsub)" ]; then ./q run_sub_group_lab.sh; else ./run_sub_group_lab.sh; fi

## ua08be3491b6b993322287772d458427 is compiling SYCL_Essentials Module4 -- SYCL Sub Groups - 7 of 7 sub_group_lab.cpp
Device: Intel(R) Data Center GPU Max 1100
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66
67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122
123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 17
0 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217
218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 26
5 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312
313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 36
0 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407
408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 45
5 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502
503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 55
0 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597
598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 64
5 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692
693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 74
0 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787
788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 83
5 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882
883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 93
0 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977
978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 10
20 1021 1022 1023

496 1520 2544 3568 4592 5616 6640 7664 8688 9712 10736 11760 12784 13808 14832 15856 16880 17904 18928 19952 20976 22000 23024 24048 25072 26096 27120 28144 29168 30192 31216 32240

Sum = 523776

```

图 2.14: 程序运行截图

我们知道，从 1 加到 1023 的结果确实为 523776，且每一个子组的和与前后相邻子组的和应该相差 1024，这些数据在图中均可以对应上，说明程序编写是正确的。

## 2.2 Intel\_Advisor

### 2.2.1 学习目标

#### 1. 卸载顾问 (Offload Advisor)

- 运行卸载顾问并生成 HTML 报告
- 阅读和了解报告中的指标
- 在目标硬件上获取应用程序的性能估计
- 确定哪些循环适合卸载

#### 2. 车顶线分析 (Roofline Analysis)

- 说明英特尔顾问如何执行 GPU 屋顶线分析。
- 使用命令行语法运行 GPU Roofline Analysis。
- 使用 GPU 屋顶线分析确定有效的优化策略。

### 2.2.2 学习内容

**Intel Offload Advisor** Intel Offload Advisor 可以用分析应用程序，特别是循环，确定哪些部分适合加速或迁移到特定硬件上。它可以进行各种分析，包括时间、浮点运算、内存传输和数据依赖，以评估循环的性能。

课程中提到了可以使用这个 Advisor 的方法，即通过命令行脚本 `collect.py` 和 `analyze.py` 运行，当然，两个脚本各有不同的任务：`collect.py` 负责收集时间、浮点运算、循环次数等数据；`analyze.py` 则基于收集的数据生成详细报告。

详细使用步骤如下：

1. 从 GitHub 克隆 Intel oneAPI 示例存储库。
2. 使用 CMake 和 Make 构建和编译应用程序。
3. 运行 `advixe-python` 使用 `collect.py` 和 `analyze.py` 收集数据并生成报告。

那么生成的报告中有什么相关信息，课程中提到：生成 HTML 报告中会包括，比如在特定硬件上的预期加速（例如 Gen9），或者是发现潜在的瓶颈，以及对于建议适合迁移和不适合迁移的循环部分的分析。同时，还会有一些附加输出，比如包括 CSV 文件（`report.csv`, `whole_app_metric.csv`）、图形表示（`program_tree.dot`, `program_tree.pdf`）以及用于调试和详细分析的 JSON/LOG 文件。

那么这些信息都能帮助开发人员优化代码，可以针对特定硬件加速器通过提供性能瓶颈的洞察和适合迁移的区域。

Offload Advisor 报告示例如下图2.15所示：

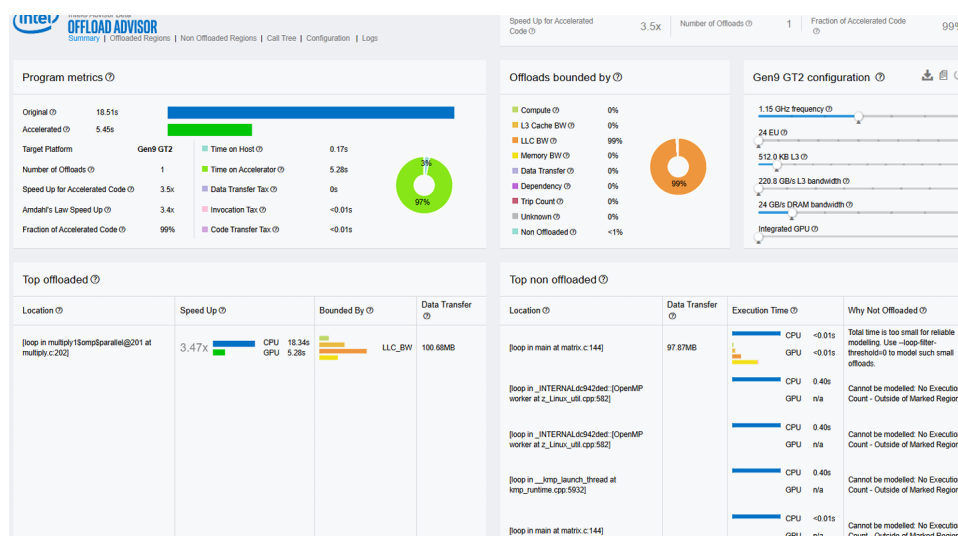


图 2.15: Offload Advisor 报告示例

根据以上总结分析，我们可以看到，Intel Offload Advisor 能够通过分析和报告应用程序的适合迁移到特定硬件的部分，帮助开发人员优化代码性能和效率。这无疑能帮助开发人员通过详细的分析和报告做出基于信息的决策，以实现最佳性能。

**Intel Roofline Analysis** 车顶线模型与 Offload Advisor 不同的是，它是一种图形化分析工具，用于评估和优化计算机程序的性能。它帮助确定程序在特定硬件上的最大理论性能极限。但是我们可

以注意到，两者针对的是计算机程序运行在硬件上能获得的最大收益的不同方面，所以两者实际上可以结合分析一个程序，以获得多方面的信息对程序进行优化和改进。

因为车顶线模型是一种图形化工具，所以它可以通过图表的形式清晰地展示程序在计算性能和内存带宽之间的关系。下图2.16是在平台上生成的一个车顶线分析图表：

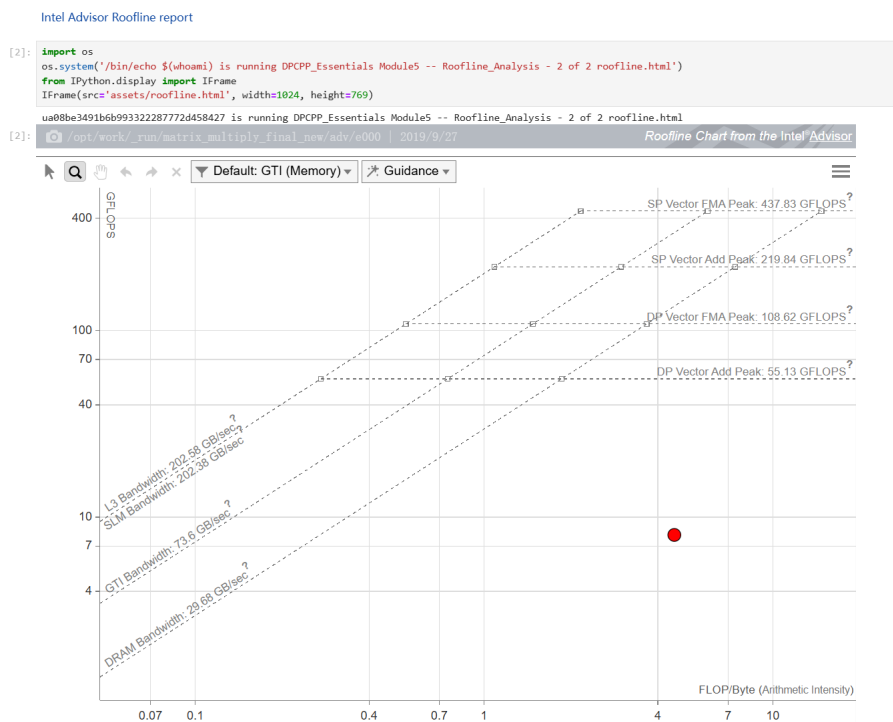


图 2.16: Roofline Analysis 分析图表示例

那么，这样的信息报告就能为开发人员提供指导，帮助确定优化策略，开发人员就能够针对硬件特性最大化程序性能，特别是在选择合适的硬件加速器或优化内存使用方面进行优化改进，从而达到最佳的加速效果。