

Министерство образования и науки Российской Федерации

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«САРАТОВСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ ИМЕНИ Н. Г. ЧЕРНЫШЕВСКОГО»

Кафедра математической
кибернетики и компьютерных наук

ВЫЧИСЛЕНИЯ НА ВИДЕОКАРТАХ

КУРСОВАЯ РАБОТА

студента 3 курса 351 группы
направления 09.03.04 — Программная инженерия
факультета КНиИТ
Григорьева Алексея Александровича

Научный руководитель
доцент

М. С. Семенов

Заведующий кафедрой
к. ф.-м. н.

С. В. Миронов

Саратов 2019

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	3
1 Краткая теория	4
1.1 Типовая модель видеокарты	4
1.2 Основные понятия OpenCL	5
2 Алгоритмы на видеокарте	8
2.1 Требования к алгоритмам	8
2.2 Настройка среды разработки	12
2.3 Инициализация OpenCL программы	13
2.4 Задачи на одномерных массивах	13
2.4.1 Вычисление суммы ряда	13
2.5 Задачи на двумерных массивах	14
ЗАКЛЮЧЕНИЕ	15
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	16
Приложение А Листинг программы	17
Приложение Б Листинг сборочных файлов CMake	26

ВВЕДЕНИЕ

Будет добавлено позже. Немного об истории, сравнение характеристик процессора и видеокарты

Цели курсовой работы:

- ознакомиться с теорией, необходимой для написания эффективных алгоритмов, исполняющихся на видеокарте, с использованием OpenCL;
- понять общие свойства архитектуры видеокарты и тем самым научиться оптимизировать алгоритмы;
- получить практический опыт разработки программ на видеокартах с помощью OpenCL;

1 Краткая теория

Составление эффективных алгоритмов вычисления на видеокарте в значительной степени отличается от привычных алгоритмов, исполняющихся на процессоре. При составлении программного кода необходимо учитывать как и общие особенности видеокарт, так и, возможно, характеристики конкретного устройства, для которого программируется алгоритм.

В данном разделе будет рассмотрена типовая модель видеокарты и основные понятия OpenCL, с которыми будем оперировать в данной работе.

1.1 Типовая модель видеокарты

Рассмотрим следующую архитектуру вычислительного устройства, используемого в видеокартах Nvidia 1.

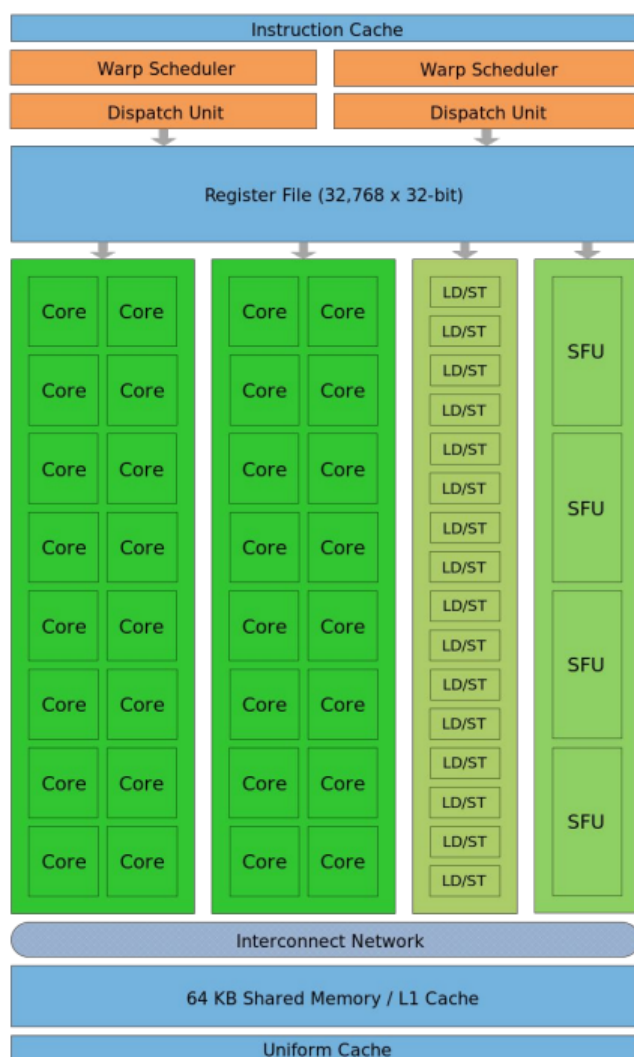


Рисунок 1 – Архитектура потокового мультипроцессора Fermi.

Вычислительное устройство в архитектуре Nvidia имеет 32 **ядра** (CUDA

cores), каждое из которых в состоянии работы является **поток**ом. В отличие от процессора, ядра выполняют более узкий набор задач, что позволяет с меньшими затратами увеличить их количество в устройстве [1]. Для управления ими существует **warp scheduler**, выполняющий роль указателя на инструкции соответствуя архитектуре SIMD. Данные для вычислений потоки берут из **локальной памяти** (shared memory), общей для всех ядер. Достигается это с использованием **устройств загрузки и хранения** (load-store units), соответственно способных загружать, а также сохранять данные в локальную память. Между всеми 32 ядрами вычислительного устройства динамически распределяются **регистры**, самая быстрая память, доступная им. У мультипроцессора в наличии намного больше регистров, чем могло быть нужно для выполнения программы. Это сделано для сокрытия времени на загрузку памяти и быстрого переключения контекста, подробнее - в разделе 2.1.

Количество таких устройств в видеокарте определяется следующим образом:

Количество ядер в видеокарте / 32, в случае Nvidia

Количество ядер в видеокарте / 64, в случае AMD

В терминологии Nvidia, поток из всех (32) активных ядер вычислительного устройства образует **warp**, Например, видеокарта Nvidia Geforce GTX 1050 Ti имеет 768 ядер CUDA, и, соответственно, 24 warp.

1.2 Основные понятия OpenCL

OpenCL — открытый для свободного пользования программный интерфейс для создания параллельных приложений, использующих многоядерные структуры как и центрального процессора (CPU), так и графического (GPU). Использование API необходимо для обеспечения совместимости программы с различными устройствами [2].

При построении задач, определяется **рабочее пространство** (NDRange), представляющее собой все возможные в рамках задачи значения индексов потоков. Размер рабочего пространства определяется программистом на этапе инициализации OpenCL программы. Рабочее пространство может представлять:

- одномерный массив длиной N элементов;
- двумерную сетку размерности NxM;

- трехмерное пространство размерностью $N \times M \times P$.

Код, выполняющийся параллельно на ядрах процессора, называется **kernel**. Копия kernel выполняется для каждого индекса рабочего пространства и называется **work-item** с глобальным ID, соответствующим некоторому ID рабочего пространства. Kernel для всех work-item в рабочем пространстве имеют одинаковый код и входные параметры, но может иметь различный путь выполнения программы соответственно своему глобальному индексу - индекс в рабочем пространстве, полученному с использованием функции `get_global_id()`. Kernel в отличие от остальной программы полностью выполняется на видеокарте [3].

Группа work-item называется **work-group**, и за каждой группой закреплен собственный warp (см. предыдущий раздел), в рамках которого work-item могут синхронизироваться. Для каждой рабочей группы существует ее индекс в рабочем пространстве, и каждый work-item может узнать свой индекс внутри рабочей группы. Нетрудно заметить следующее соотношение:

$$\text{global ID} = \text{group ID} * \text{размер группы} + \text{local ID}$$

Размер рабочей группы аналогично рабочему пространству определяется программистом.

Каждое ядро, выполняя заданный kernel, является work-item в некоторой рабочей группой, на которые разделено рабочее пространство `NDRange`.

Рассмотрим на примере следующей схемы 2 другие виды сущностей, с которыми будет взаимодействие в OpenCL.

- Платформа — драйвер, модель взаимодействия OpenCL и устройства. Распространены платформы от следующих производителей: Nvidia, Intel, AMD.
- Программа — хостовая часть, организующая подготовку к вычислениям и набор kernel-подпрограмм.
- Kernel — программа, исполняющаяся на видеокарте в каждом ядре.
- Контекст — окружение, в котором исполняется kernel.
- Объект памяти — создаваемый в контексте объект.
- Буфер — произвольный массив данных.

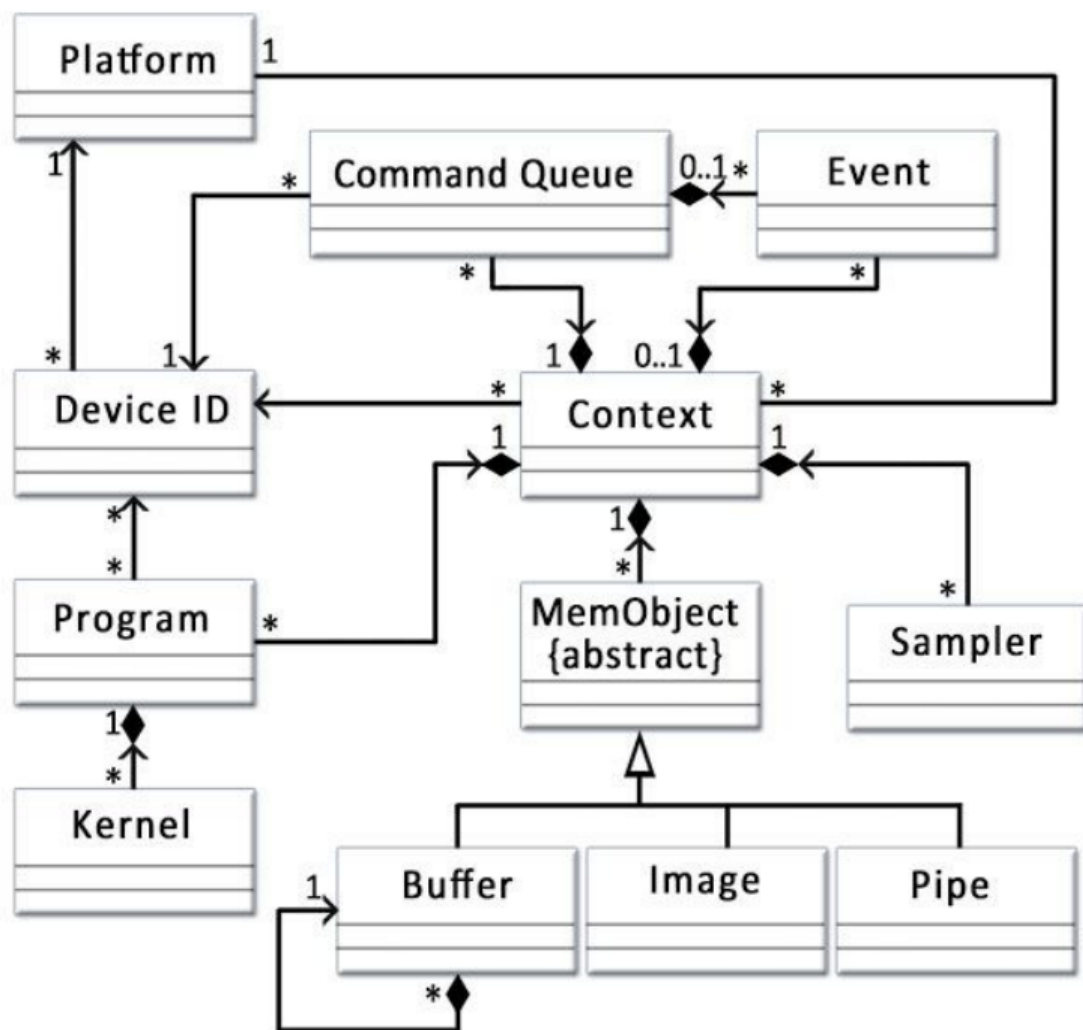


Рисунок 2 – Основные сущности в OpenCL.

2 Алгоритмы на видеокарте

В данном разделе будет рассмотрена анализ и практическая реализация алгоритмов на видеокарте, включая:

- описание общих требований к алгоритмам на основе доступа к памяти и параллельного исполнения;
- настройка среды разработки Microsoft Visual Studio 2017 под выполнение параллельных программ с использованием OpenCL;
- написание программ для задач, использующих входные данные разных размерностей.

2.1 Требования к алгоритмам

Любой алгоритм можно вычислить на видеокарте, но эффективность в сравнении с реализацией на центральном процессоре зависит от корректного построения алгоритма для видеокарты.

Основным требованием к составлению алгоритма на видеокарте считается наличие массового параллелизма. Он заключается в том что задачу можно разбить на рабочие группы так, что не будет требоваться постоянная синхронизация между work-item из разных рабочих групп.

Следует вспомнить, что все потоки в warp выполняют одинаковые инструкции в любой момент времени. Какая инструкция будет выполняться следующей определяется с помощью warp scheduler, единого для всех потоков в warp. Рассмотрим следующий фрагмент кода:

```
if (predicate) {  
    value = x[i];  
}  
else {  
    value = y[i];  
}
```

Учитывая сказанное выше, все потоки при срабатывания if-части должны выполнить внутреннюю часть, однако это не совсем так, и если у потока предикат — False, он будет спрятан от выполнения внутренней части, аналогично и с else-частью. Однако несмотря на то что результат выполнения конструкции if-else будет верным, часть потоков будет простаивать, ожидая выполнение маскированных для них частей.

Данная ситуация называется *code divergence*, и она может стать причиной низкой производительности программы. Этого можно избежать, если организовать код таким образом чтобы для всех потоков предикат возвращал одинаковое значение, тогда конструкция не соответствующая ему будет пропущена указателем на инструкции. Если это невозможно, то для эффективного выполнения алгоритма рекомендуется отказаться от многочисленных сложных ветвлений, так как сложность выполнения фрагмента алгоритма будет вычисляться как сумма *if*- и *else*- частей вместо максимума как в последовательных программах.

При выборе размера рабочих групп стоит учитывать особенности алгоритма, однако, есть некоторые общие правила, которых необходимо придерживаться.

1. Размер рабочей группы не должен быть меньше *warp*.
2. Размер рабочей группы должен быть кратен 32 (64 если используется AMD).

В противном случае, некоторые потоки будут простаивать, ожидая пока остальные завершат свою работу

Как известно, операции с памятью являются одними из самых долгих по времени выполнения, в связи с этим было решено сделать разбиение задач на рабочие группы, в результате у видеокарт появился аналог имеющегося у процессоров *hyper-threading*. Он заключается в использовании каждым вычислительным устройством регистров для переключения контекста при задержке, созданной обращением к памяти (*latency*) [4].

Другими словами, *warp* может быстро сохранить состояние выполнения в данной рабочей группе и пока выполняется долгая операция обращения к памяти, вычислительное устройство может переключиться на другой *warp* в рабочей группе, и если второй *warp* хочет выполнить операцию обращения к памяти, то происходит возвращение к первому *warp* если доступ к памяти завершился, либо активируется третий *warp* и так далее. Следствие — высокая вычислительная мощность и большая пропускная способность видеокарты [5].

Количество одновременно активных *warp* в рабочей группе определяется как минимум из:

- количества регистров / количество используемых в *warp* регистров;
- количества локальной памяти / количество используемой локальной па-

мяти;

- максимально допустимого количества warp (~10).

В соответствии с этим существует величина **occupancy**, определяемая соотношением

$$\text{среднее кол-во активных warp} / \text{максимальное кол-во активных warp}$$

Не всегда высокий occupancy означает что программа имеет высокую производительность. Например, если доступ к памяти в программе очень быстрый, только из регистров, то необходимости в сокращении задержки и переключения контекста нет, и occupancy будет низким.

Однако, низкий occupancy и высокая задержка при обращении к памяти может означать что программа написана не достаточно эффективно, и ей необходимы улучшения, если это возможно.

Чем больше на одном вычислителе warp — тем реже все warp оказываются в состоянии «ждем запрос памяти» и тем реже вычислитель будет простаивать, т.к. тем чаще у него находится рабочая группа в которой можно что-то посчитать [4].

Если потоки из одного warp делают запрос к памяти, то эти запросы склеются в столько запросов, сколькими кеш-линиями покрываются запрошенные данные.

Другими словами, если потоки запрашивают данные, которые в памяти лежат подряд, то достигнутая пропускная способность будет максимальная так как запросы «склеются». Размер кеш-линии обычно от 32 до 128 байт.

Если приложение использует OpenCL 1.x, то размеры NDRange должны нацело (без остатка) делиться на размеры рабочих групп. Там, где данные образуют NDRange с другим размером, необходимо самостоятельно изменить их чтобы выполнялось это условие, например, добавлением нулей или средних значений, которые не будут значимо влиять на результат вычислений [6].

В OpenCL 2.0 появилась новая возможность, в которой устранена данная проблема. Речь идет о так называемых неоднородных рабочих группах: выполняемый модуль OpenCL 2.0 может разделить NDRange на рабочие группы неоднородного размера по любому измерению. Если разработчик укажет размер рабочей группы, на который размер NDRange не делится нацело, выполняемый модуль разделит NDRange таким образом, чтобы создать как мож-

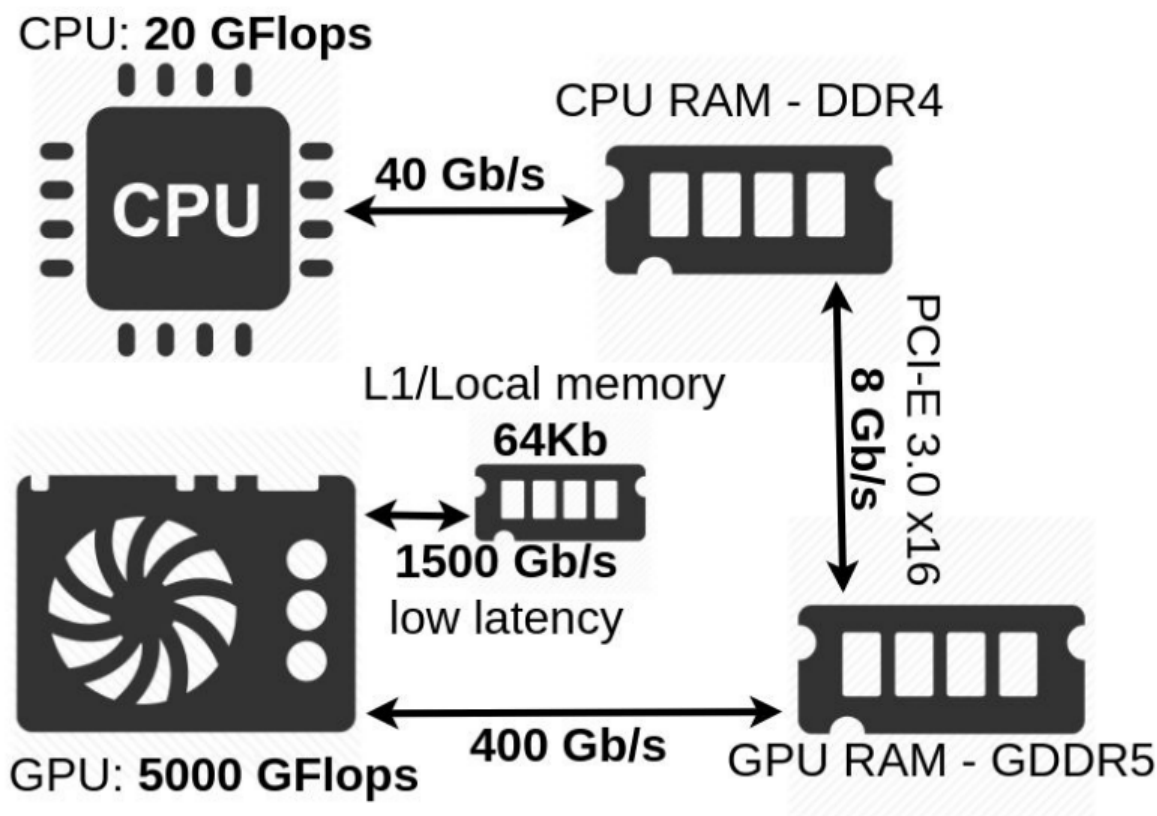


Рисунок 3 – Доступные ресурсы — память.

но больше рабочих групп с указанным размером, а остальные рабочие группы будут иметь другой размер. Например, для NDRange размером 1918x1078 рабочих элементов при размере рабочей группы 16x16 элементов среда выполнения OpenCL 2.0 разделит NDRange, как показано на приведенном ниже рисунке 4.

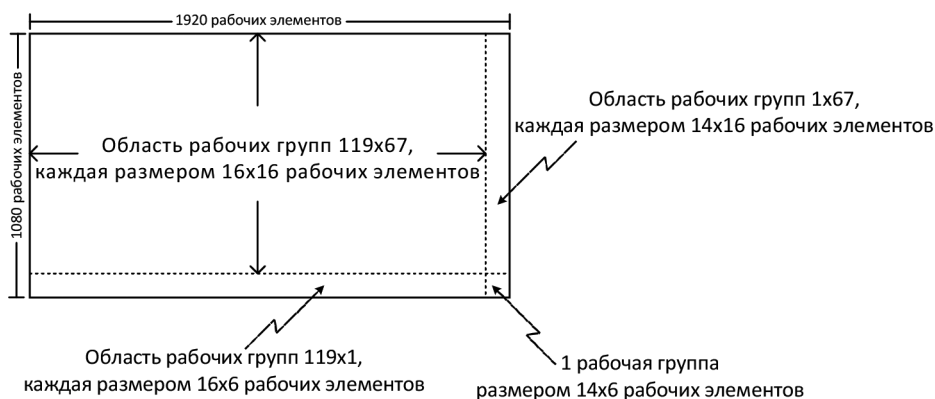


Рисунок 4 – Разделение NDRange на рабочие группы разных размеров.

2.2 Настройка среды разработки

В данном разделе будет рассмотрен процесс настройки среды разработки и создания первого OpenCL-проекта.

В качестве среды разработки для программирования с использованием OpenCL выбрана Microsoft Visual Studio, язык программирования — C++.

На компьютер была установлена реализация OpenCL от Nvidia — Nvidia GPU Computing SDK. А также программа CMake, являющаяся независимым от платформы инструментом для сборки проектов 5.

С помощью графического интерфейса выберем расположение файлов исходного кода и места сборки проекта. Директория с исходными файлами должны содержать текстовые файлы CMakeLists из приложения Б. Если OpenCL установлен корректно, то нажатие кнопки «Configure» выведет найденные на компьютеры файлы, связанные с OpenCL. Нажмем «Generate», и перейдем в папку с проектом, в котором можно увидеть созданный файл .sln проекта Microsoft Visual Studio, сконфигурированного под OpenCL.

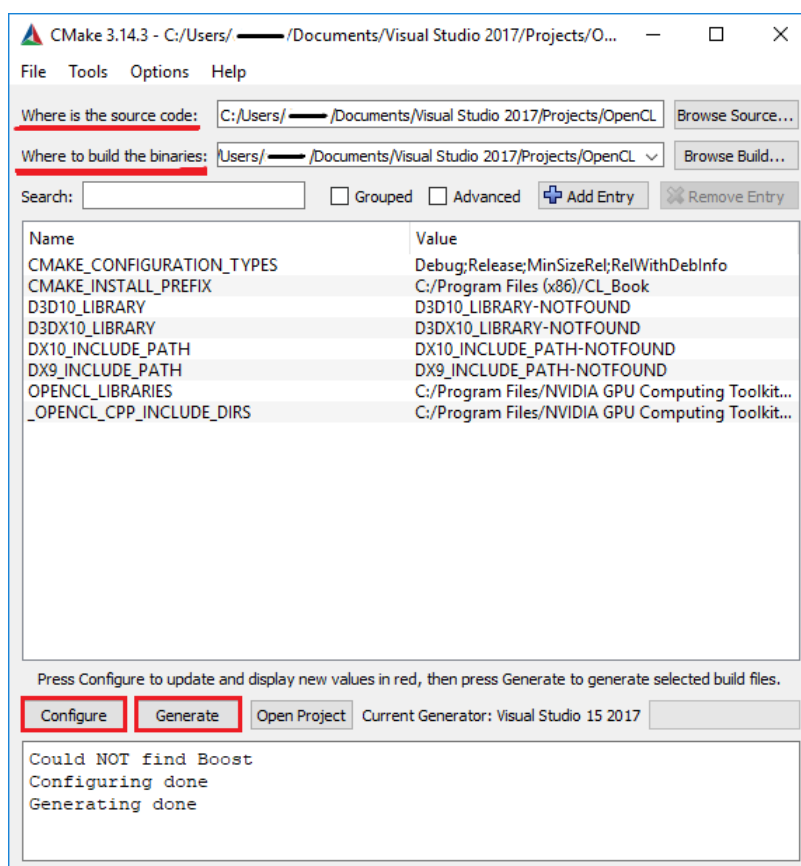


Рисунок 5 – Окно программы CMake-gui.

2.3 Инициализация OpenCL программы

В данном разделе будут рассмотрены базовые функции, необходимые для инициализации параллельной программы с использованием OpenCL. Данные функции будут предварять задачи из следующих разделов.

Рассмотрим пример, взятый из руководства по OpenCL [7]. С полным кодом, содержащим комментарии, переведенными на русский язык, можно ознакомиться в приложении А, файл HelloWorld.cpp. Обратим внимание на последовательность действий в функции main(). Многие понятия из данного раздела подробно описаны в 1.2.

Сначала с помощью функции CreateContext() создается контекст на основе первой найденной на компьютере платформы. Далее для первого доступного устройства в контексте создается командная очередь clCreateCommandQueue(), а в случае неудачи запускается функция очистки и программа завершается с кодом ошибки 1.

Из файла с исходным кодом kernel HelloWorld.cl создается OpenCL программа clCreateProgramWithSource(). После этого создается и сам kernel на основе созданной «программы» clCreateKernel().

После этого создаются объекты памяти для конкретной задачи clCreateBuffer(), и каждый из них поочередно загружается в kernel с помощью clSetKernelArg(). Затем kernel ставится в очередь на выполнение clEnqueueNDRangeKernel(), и после завершения работы, выводим в консоль буфер-результат, являющимся результатом выполнения данной OpenCL программы clEnqueueReadBuffer().

2.4 Задачи на одномерных массивах

(сумма ряда, какая-нибудь префикс-функция)

Решим типовые задачи на одномерных массивах. Самая тривиальная — посчитать сумму двух векторов. Даны векторы **a** и **b**, на основе их суммы должен получиться вектор **result**.

2.4.1 Вычисление суммы ряда

Да

```
1
2 __kernel void hello_kernel(__global const float *a,
3                               __global const float *b,
```

```

4                                     __global float *result)
5 {
6     int gid = get_global_id(0);
7
8     result[gid] = a[gid] + b[gid];
9 }

```

Рассмотрим kernel для решения данной задачи. Каждый work-item, исполняя копию kernel, узнает свой глобальный ID в рабочем пространстве с помощью `get_global_id(0)`. Таким образом он читает данные из локальной памяти соответствуя своему индексу. В основной программе выбран `NDRange` размерности 1024 и `local_work_size`, размер рабочей группы, равный 32. Эти параметры переданы в функцию `clEnqueueNDRangeKernel()`

Пусть значения компонент первого вектор соответствуют номеру компоненты начиная с 0, а значения компонент второго вектора соответствуют удвоенному номеру компоненты в векторе. Результат работы данной программы представлен на изображении 6.

```

C:\WINDOWS\system32\cmd.exe
86 1089 1092 1095 1098 1101 1104 1107 1110 1113 1116 1119 1122 1125 1128 1131 1134 1137 1140 1143 1146 1149 1152 1155 1158 1161 1164 1167 1170 1173 1176 1179 1182 1185 1188 1191 1194 1197 1200 1203 1206 1209 1212 1215 1218 1221 1224 1227 1230
30 1233 1236 1239 1242 1245 1248 1251 1254 1257 1260 1263 1266 1269 1272 1275 1278 1281 1284 1287 1290 1293 1296 1299 1302 1305 1308 1311 1314 1317 1320 1323 1326 1329 1332 1335 1338 1341 1344 1347 1350 1353 1356 1359 1362 1365 1368 1371 1374
74 1377 1380 1383 1386 1389 1392 1395 1398 1401 1404 1407 1410 1413 1416 1419 1422 1425 1428 1431 1434 1437 1440 1443 1446 1449 1452 1455 1458 1461 1464 1467 1470 1473 1476 1479 1482 1485 1488 1491 1494 1497 1500 1503 1506 1509 1512 1515 1518
46 1449 1452 1455 1458 1461 1464 1467 1470 1473 1476 1479 1482 1485 1488 1491 1494 1497 1500 1503 1506 1509 1512 1515 1518 1521 1524 1527 1530 1533 1536 1539 1542 1545 1548 1551 1554 1557 1560 1563 1566 1569 1572 1575 1578 1581 1584 1587 1590
18 1521 1524 1527 1530 1533 1536 1539 1542 1545 1548 1551 1554 1557 1560 1563 1566 1569 1572 1575 1578 1581 1584 1587 1590 1593 1596 1599 1602 1605 1608 1611 1614 1617 1620 1623 1626 1629 1632 1635 1638 1641 1644 1647 1650 1653 1656 1659 1662
62 1665 1668 1671 1674 1677 1680 1683 1686 1689 1692 1695 1698 1701 1704 1707 1710 1713 1716 1719 1722 1725 1728 1731 1734 1737 1740 1743 1746 1749 1752 1755 1758 1761 1764 1767 1770 1773 1776 1779 1782 1785 1788 1791 1794 1797 1800 1803 1806
06 1809 1812 1815 1818 1821 1824 1827 1830 1833 1836 1839 1842 1845 1848 1851 1854 1857 1860 1863 1866 1869 1872 1875 1878 1881 1884 1887 1890 1893 1896 1899 1902 1905 1908 1911 1914 1917 1920 1923 1926 1929 1932 1935 1938 1941 1944 1947 1950
78 1881 1884 1887 1890 1893 1896 1899 1902 1905 1908 1911 1914 1917 1920 1923 1926 1929 1932 1935 1938 1941 1944 1947 1950 1953 1956 1959 1962 1965 1968 1971 1974 1977 1980 1983 1986 1989 1992 1995 1998 2001 2004 2007 2010 2013 2016 2019 2022
22 2025 2028 2031 2034 2037 2040 2043 2046 2049 2052 2055 2058 2061 2064 2067 2070 2073 2076 2079 2082 2085 2088 2091 2094 2097 2100 2103 2106 2109 2112 2115 2118 2121 2124 2127 2130 2133 2136 2139 2142 2145 2148 2151 2154 2157 2160 2163 2166
94 2097 2100 2103 2106 2109 2112 2115 2118 2121 2124 2127 2130 2133 2136 2139 2142 2145 2148 2151 2154 2157 2160 2163 2166 2169 2172 2175 2178 2181 2184 2187 2190 2193 2196 2199 2202 2205 2208 2211 2214 2217 2220 2223 2226 2229 2232 2235 2238
66 2169 2172 2175 2178 2181 2184 2187 2190 2193 2196 2199 2202 2205 2208 2211 2214 2217 2220 2223 2226 2229 2232 2235 2238 2241 2244 2247 2250 2253 2256 2259 2262 2265 2268 2271 2274 2277 2280 2283 2286 2289 2292 2295 2298 2301 2304 2307 2310
38 2241 2244 2247 2250 2253 2256 2259 2262 2265 2268 2271 2274 2277 2280 2283 2286 2289 2292 2295 2298 2301 2304 2307 2310 2313 2316 2319 2322 2325 2328 2331 2334 2337 2340 2343 2346 2349 2352 2355 2358 2361 2364 2367 2370 2373 2376 2379 2382
10 2313 2316 2319 2322 2325 2328 2331 2334 2337 2340 2343 2346 2349 2352 2355 2358 2361 2364 2367 2370 2373 2376 2379 2382 2385 2388 2391 2394 2397 2400 2403 2406 2409 2412 2415 2418 2421 2424 2427 2430 2433 2436 2439 2442 2445 2448 2451 2454
82 2385 2388 2391 2394 2397 2400 2403 2406 2409 2412 2415 2418 2421 2424 2427 2430 2433 2436 2439 2442 2445 2448 2451 2454 2457 2460 2463 2466 2469 2472 2475 2478 2481 2484 2487 2490 2493 2496 2499 2502 2505 2508 2511 2514 2517 2520 2523 2526
54 2457 2460 2463 2466 2469 2472 2475 2478 2481 2484 2487 2490 2493 2496 2499 2502 2505 2508 2511 2514 2517 2520 2523 2526 2529 2532 2535 2538 2541 2544 2547 2550 2553 2556 2559 2562 2565 2568 2571 2574 2577 2580 2583 2586 2589 2592 2595 2598
26 2529 2532 2535 2538 2541 2544 2547 2550 2553 2556 2559 2562 2565 2568 2571 2574 2577 2580 2583 2586 2589 2592 2595 2598 2601 2604 2607 2610 2613 2616 2619 2622 2625 2628 2631 2634 2637 2640 2643 2646 2649 2652 2655 2658 2661 2664 2667 2670
98 2601 2604 2607 2610 2613 2616 2619 2622 2625 2628 2631 2634 2637 2640 2643 2646 2649 2652 2655 2658 2661 2664 2667 2670 2673 2676 2679 2682 2685 2688 2691 2694 2697 2700 2703 2706 2709 2712 2715 2718 2721 2724 2727 2730 2733 2736 2739 2742
70 2673 2676 2679 2682 2685 2688 2691 2694 2697 2700 2703 2706 2709 2712 2715 2718 2721 2724 2727 2730 2733 2736 2739 2742 2745 2748 2751 2754 2757 2760 2763 2766 2769 2772 2775 2778 2781 2784 2787 2790 2793 2796 2799 2802 2805 2808 2811 2814
42 2745 2748 2751 2754 2757 2760 2763 2766 2769 2772 2775 2778 2781 2784 2787 2790 2793 2796 2799 2802 2805 2808 2811 2814 2817 2820 2823 2826 2829 2832 2835 2838 2841 2844 2847 2850 2853 2856 2859 2862 2865 2868 2871 2874 2877 2880 2883 2886
14 2817 2820 2823 2826 2829 2832 2835 2838 2841 2844 2847 2850 2853 2856 2859 2862 2865 2868 2871 2874 2877 2880 2883 2886 2889 2892 2895 2898 2901 2904 2907 2910 2913 2916 2919 2922 2925 2928 2931 2934 2937 2940 2943 2946 2949 2952 2955 2958
86 2889 2892 2895 2898 2901 2904 2907 2910 2913 2916 2919 2922 2925 2928 2931 2934 2937 2940 2943 2946 2949 2952 2955 2958 2961 2964 2967 2970 2973 2976 2979 2982 2985 2988 2991 2994 2997 3000 3003 3006 3009 3012 3015 3018 3021 3024 3027 3030
58 2961 2964 2967 2970 2973 2976 2979 2982 2985 2988 2991 2994 2997 3000 3003 3006 3009 3012 3015 3018 3021 3024 3027 3030 3033 3036 3039 3042 3045 3048 3051 3054 3057 3060 3063 3066 3069
30 3033 3036 3039 3042 3045 3048 3051 3054 3057 3060 3063 3066 3069
Executed program successfully.
Для продолжения нажмите любую клавишу . . .

```

Рисунок 6 – Результат работы программы по суммированию двух векторов .

2.5 Задачи на двумерных массивах

(транспонирование, перемножение матриц)

ЗАКЛЮЧЕНИЕ

В настоящей работе

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 [NVIDIA's Next Generation CUDA Compute Architecture: Fermi]. — URL: https://www.nvidia.com/content/PDF/fermi_white_papers/NVIDIA_Fermi_Compute_Architecture_Whitepaper.pdf (Дата обращения 12.05.2019). Загл. с экр. Яз. англ.
- 2 [Введение | OpenCL]. — URL: <http://opencl.ru/node/8> (Дата обращения 12.05.2019). Загл. с экр. Яз. рус.
- 3 [The OpenCL Specification]. — URL: <https://www.khronos.org/registry/OpenCL/specs/> (Дата обращения 12.05.2019). Загл. с экр. Яз. англ.
- 4 [Введение в OpenCL. Архитектура видеокарты]. — URL: https://compscicenter.ru/courses/video_cards_computation/2018-autumn/classes/3980/ (Дата обращения 12.05.2019). Загл. с экр. Яз. рус.
- 5 [Achieved Occupancy]. — URL: <https://docs.nvidia.com/gameworks/content/developertools/desktop/analysis/report/cudaexperiments/kernellevel/achievedoccupancy.htm> (Дата обращения 12.05.2019). Загл. с экр. Яз. англ.
- 6 [Неоднородные рабочие группы OpenCL 2.0]. — URL: <https://software.intel.com/ru-ru/articles/opencl-20-non-uniform-work-groups> (Дата обращения 12.05.2019). Загл. с экр. Яз. рус.
- 7 *Aaftab Munshi Benedict R. Gaster, T. G. M. J. F. D. G. OpenCL Programming Guide / T. G. M. J. F. D. G. Aaftab Munshi, Benedict R. Gaster.* — Ann Arbor, Michigan: Edwards Brothers, 2012.

ПРИЛОЖЕНИЕ А

Листинг программы

Вычисление суммы двух массивов.

```
1 // HelloWorld.cpp
2 //
3 //     В данном примере продемонстрирована базовая установка и использование OpenCL
4 //
5
6 #include <iostream>
7 #include <fstream>
8 #include <sstream>
9
10 #ifdef __APPLE__
11 #include <OpenCL/cl.h>
12 #else
13 #include <CL/cl.h>
14 #endif
15
16 ///
17 //     Константы
18 //
19 const int ARRAY_SIZE = 1000;
20
21 ///
22 //     Создание OpenCL контекста на основе доступной платформы,
23 //     использующей GPU (в приоритете) или CPU
24 //
25 cl_context CreateContext()
26 {
27     cl_int errNum;
28     cl_uint numPlatforms;
29     cl_platform_id firstPlatformId;
30     cl_context context = NULL;
31
32     // Выберем OpenCL платформу, на которой будет запущен код.
33     // В данном примере выберем первую доступную платформу.
34     errNum = clGetPlatformIDs(1, &firstPlatformId, &numPlatforms);
35     if (errNum != CL_SUCCESS || numPlatforms <= 0)
36     {
37         std::cerr << "Failed to find any OpenCL platforms." << std::endl;
38         return NULL;
```

```

39     }
40
41     // Создадим OpenCL контекст на заданной платформе.
42     // Попробуем создать основанный на GPU контекст и в случае
43     // неудача попробуем создать основанный на CPU контекст
44     cl_context_properties contextProperties[] =
45     {
46         CL_CONTEXT_PLATFORM,
47         (cl_context_properties)firstPlatformId,
48         0
49     };
50     context = clCreateContextFromType(contextProperties, CL_DEVICE_TYPE_GPU,
51                                     NULL, NULL, &errNum);
52     if (errNum != CL_SUCCESS)
53     {
54         std::cout << "Could not create GPU context, trying CPU..." << std::endl;
55         context = clCreateContextFromType(contextProperties, CL_DEVICE_TYPE_CPU,
56                                     NULL, NULL, &errNum);
57         if (errNum != CL_SUCCESS)
58         {
59             std::cerr << "Failed to create an OpenCL GPU or CPU context." << std::endl;
60             return NULL;
61         }
62     }
63
64     return context;
65 }
66
67 ///
68 //  Создание командной очередь для первого доступного
69 //  устройства из контекста
70 //
71 cl_command_queue CreateCommandQueue(cl_context context, cl_device_id *device)
72 {
73     cl_int errNum;
74     cl_device_id *devices;
75     cl_command_queue commandQueue = NULL;
76     size_t deviceBufferSize = -1;
77
78     // Получить размер буфера устройства
79     errNum = clGetContextInfo(context, CL_CONTEXT_DEVICES, 0, NULL, &deviceBufferSize);

```

```

80     if (errNum != CL_SUCCESS)
81     {
82         std::cerr << "Failed call to clGetContextInfo(...,GL_CONTEXT_DEVICES,...)";
83         return NULL;
84     }
85
86     if (deviceBufferSize <= 0)
87     {
88         std::cerr << "No devices available.";
89         return NULL;
90     }
91
92     // Выделить память под буфер устройства
93     devices = new cl_device_id[deviceBufferSize / sizeof(cl_device_id)];
94     errNum = clGetContextInfo(context, CL_CONTEXT_DEVICES, deviceBufferSize, devices,
95     if (errNum != CL_SUCCESS)
96     {
97         delete [] devices;
98         std::cerr << "Failed to get device IDs";
99         return NULL;
100     }
101
102     // Выбор первого доступного устройства
103     commandQueue = clCreateCommandQueue(context, devices[0], 0, NULL);
104     if (commandQueue == NULL)
105     {
106         delete [] devices;
107         std::cerr << "Failed to create commandQueue for device 0";
108         return NULL;
109     }
110
111     *device = devices[0];
112     delete [] devices;
113     return commandQueue;
114 }
115
116 ///
117 //  Создание OpenCL программы из файла-kernel
118 //
119 cl_program CreateProgram(cl_context context, cl_device_id device, const char* fileName)
120 {

```

```

121     cl_int errNum;
122     cl_program program;
123
124     std::ifstream kernelFile(fileName, std::ios::in);
125     if (!kernelFile.is_open())
126     {
127         std::cerr << "Failed to open file for reading: " << fileName << std::endl;
128         return NULL;
129     }
130
131     std::ostringstream oss;
132     oss << kernelFile.rdbuf();
133
134     std::string srcStdStr = oss.str();
135     const char *srcStr = srcStdStr.c_str();
136     program = clCreateProgramWithSource(context, 1,
137                                         (const char**)&srcStr,
138                                         NULL, NULL);
139     if (program == NULL)
140     {
141         std::cerr << "Failed to create CL program from source." << std::endl;
142         return NULL;
143     }
144
145     errNum = clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
146     if (errNum != CL_SUCCESS)
147     {
148         char buildLog[16384];
149         clGetProgramBuildInfo(program, device, CL_PROGRAM_BUILD_LOG,
150                               sizeof(buildLog), buildLog, NULL);
151
152         std::cerr << "Error in kernel: " << std::endl;
153         std::cerr << buildLog;
154         clReleaseProgram(program);
155         return NULL;
156     }
157
158     return program;
159 }
160
161 ///

```

```

162 // Создание объектов-памяти для kernel
163 // Kernel принимает 3 аргумента: result (вывод), a (ввод),
164 // and b (ввод)
165 //
166 bool CreateMemObjects(cl_context context, cl_mem memObjects[3],
167                      float *a, float *b)
168 {
169     memObjects[0] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
170                                     sizeof(float) * ARRAY_SIZE, a, NULL);
171     memObjects[1] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
172                                     sizeof(float) * ARRAY_SIZE, b, NULL);
173     memObjects[2] = clCreateBuffer(context, CL_MEM_READ_WRITE,
174                                     sizeof(float) * ARRAY_SIZE, NULL, NULL);
175
176     if (memObjects[0] == NULL || memObjects[1] == NULL || memObjects[2] == NULL)
177     {
178         std::cerr << "Error creating memory objects." << std::endl;
179         return false;
180     }
181
182     return true;
183 }
184
185 ///
186 // Очистка от созданных OpenCL ресурсов
187 //
188 void Cleanup(cl_context context, cl_command_queue commandQueue,
189             cl_program program, cl_kernel kernel, cl_mem memObjects[3])
190 {
191     for (int i = 0; i < 3; i++)
192     {
193         if (memObjects[i] != 0)
194             clReleaseMemObject(memObjects[i]);
195     }
196     if (commandQueue != 0)
197         clReleaseCommandQueue(commandQueue);
198
199     if (kernel != 0)
200         clReleaseKernel(kernel);
201
202     if (program != 0)

```

```

203         clReleaseProgram(program);
204
205     if (context != 0)
206         clReleaseContext(context);
207
208 }
209
210 ///
211 //         main() для HelloWorld
212 //
213 int main(int argc, char** argv)
214 {
215     cl_context context = 0;
216     cl_command_queue commandQueue = 0;
217     cl_program program = 0;
218     cl_device_id device = 0;
219     cl_kernel kernel = 0;
220     cl_mem memObjects[3] = { 0, 0, 0 };
221     cl_int errNum;
222
223     // Создание OpenCL контекста для первой доступной платформы
224     context = CreateContext();
225     if (context == NULL)
226     {
227         std::cerr << "Failed to create OpenCL context." << std::endl;
228         return 1;
229     }
230
231     // Создание очереди команд для первого доступного устройства
232     // в заданном контексте
233     commandQueue = CreateCommandQueue(context, &device);
234     if (commandQueue == NULL)
235     {
236         Cleanup(context, commandQueue, program, kernel, memObjects);
237         return 1;
238     }
239
240     // Создание OpenCL программы из файла исходного кода HelloWorld.cl для kernel
241     program = CreateProgram(context, device, "HelloWorld.cl");
242     if (program == NULL)
243     {

```

```

244     Cleanup(context, commandQueue, program, kernel, memObjects);
245     return 1;
246 }
247
248 // Create OpenCL kernel
249 kernel = clCreateKernel(program, "hello_kernel", NULL);
250 if (kernel == NULL)
251 {
252     std::cerr << "Failed to create kernel" << std::endl;
253     Cleanup(context, commandQueue, program, kernel, memObjects);
254     return 1;
255 }
256
257 // Создание объектов памяти, используемых kernel.
258 // Сначала создаются объекты памяти, содержащие данные
259 // для аргументов kernel
260 float result[ARRAY_SIZE];
261 float a[ARRAY_SIZE];
262 float b[ARRAY_SIZE];
263 for (int i = 0; i < ARRAY_SIZE; i++)
264 {
265     a[i] = (float)i;
266     b[i] = (float)(i * 2);
267 }
268
269 if (!CreateMemObjects(context, memObjects, a, b))
270 {
271     Cleanup(context, commandQueue, program, kernel, memObjects);
272     return 1;
273 }
274
275 // Задание аргументов kernel (a, b, result).
276 errNum = clSetKernelArg(kernel, 0, sizeof(cl_mem), &memObjects[0]);
277 errNum |= clSetKernelArg(kernel, 1, sizeof(cl_mem), &memObjects[1]);
278 errNum |= clSetKernelArg(kernel, 2, sizeof(cl_mem), &memObjects[2]);
279 if (errNum != CL_SUCCESS)
280 {
281     std::cerr << "Error setting kernel arguments." << std::endl;
282     Cleanup(context, commandQueue, program, kernel, memObjects);
283     return 1;
284 }

```

```

285
286     size_t globalWorkSize[1] = { ARRAY_SIZE };
287     size_t localWorkSize[1] = { 1 };
288
289     // Поставить kernel в очередь на исполнение
290     errNum = clEnqueueNDRangeKernel(commandQueue, kernel, 1, NULL,
291                                     globalWorkSize, localWorkSize,
292                                     0, NULL, NULL);
293     if (errNum != CL_SUCCESS)
294     {
295         std::cerr << "Error queuing kernel for execution." << std::endl;
296         Cleanup(context, commandQueue, program, kernel, memObjects);
297         return 1;
298     }
299
300     // Считать выходной буфер в основную программу
301     errNum = clEnqueueReadBuffer(commandQueue, memObjects[2], CL_TRUE,
302                                  0, ARRAY_SIZE * sizeof(float), result,
303                                  0, NULL, NULL);
304     if (errNum != CL_SUCCESS)
305     {
306         std::cerr << "Error reading result buffer." << std::endl;
307         Cleanup(context, commandQueue, program, kernel, memObjects);
308         return 1;
309     }
310
311     // Вывод результирующего буфера
312     for (int i = 0; i < ARRAY_SIZE; i++)
313     {
314         std::cout << result[i] << " ";
315     }
316     std::cout << std::endl;
317     std::cout << "Executed program succesfully." << std::endl;
318     Cleanup(context, commandQueue, program, kernel, memObjects);
319
320     return 0;
321 }

```

Kernel для данной задачи

```

1
2 __kernel void hello_kernel(__global const float *a,

```



```
3                                     __global const float *b,
4                                     __global float *result)
5 {
6     int gid = get_global_id(0);
7
8     result[gid] = a[gid] + b[gid];
9 }
```

ПРИЛОЖЕНИЕ Б

Листинг сборочных файлов CMake

Код сборочного файла CMakeLists.txt

```
1 # This is an example project to show and test the usage of the FindOpenCL
2 # script.
3
4 cmake_minimum_required( VERSION 2.6 )
5 project( CL_Book )
6
7 set(CMAKE_MODULE_PATH "${CMAKE_SOURCE_DIR}/cmake")
8
9 find_package( OpenCL REQUIRED )
10
11 include_directories( ${OPENCL_INCLUDE_DIRS} )
12 include_directories( "${CMAKE_SOURCE_DIR}/khronos" )
13
14 SUBDIRS(src)
15
```