



# Особенности разработки под графические ускорители на примере OpenCL: от архитектуры вычислителя до особенностей оптимизации

Семён Григорьев

Санкт-Петербургский Государственный Университет

27 января 2025

- Доцент кафедры системного программирования Санкт-Петербургского Государственного Университета
- Области интересов
  - ▶ **Высокопроизводительный анализ графов**
    - ★ В том числе, с использованием **графических ускорителей**
  - ▶ **Высокопроизводительная линейная алгебра** для анализа графов
    - ★ **Обобщённая**: матрицы и вектора параметризованы типом элемента, операции над ними могут быть заданы пользователем
    - ★ **Разреженная**: специализированные структуры для хранения матриц и векторов, специализированные алгоритмы для их обработки
    - ★ В том числе, с использованием **графических ускорителей**
  - ▶ Высокоуровневые средства разработки для графических ускорителей



- Email: [rsdpisuy@gmail.com](mailto:rsdpisuy@gmail.com)
- GitHub: [gsvgit](#)
- Google Scholar: [Semyon Grigorev](#)
- DBLP: [Semyon V. Grigorev](#)

Цель лекции: познакомить с графическими ускорителями, показать, на что стоит обращать внимание при первых попытках программировать под них

- 1 Параллельность: что, зачем, какая бывает, какая нас будет интересовать
- 2 CPU и GPGPU: сходства, различия, взаимосвязь
- 3 Основные принципы архитектуры GPGPU
- 4 Основы OpenCL
- 5 Немного про оптимизации для GPGPU

SIMT

Полностью независимые задачи

MIMD

SPMD

...

Bulk synchronous parallelism

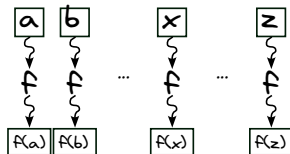
Irregular parallelism

Message-passing parallelism

Параллелизм по данным

# Параллелизм по данным

- Одна и та же функция применяется **независимо** ко всем элементам набора данных<sup>1</sup>



- Single Instruction Multiple Data (SIMD)

- ▶ Не функция, но **инструкция**: поток управления функции должен быть максимально простым
  - ⇒ Много ветвлений — **плохо**<sup>2</sup>
  - ⇒ Много однотипных вычислений — **хорошо**
  - ⇒ Данные должны быть **однородными**
- ▶ Типично для обработки изображений, линейной алгебры
  - ★ Вектора, матрицы — хорошие массивы данных
  - ★ Много вычислительно интенсивной, но «однотипной» работы

<sup>1</sup> Можно представить map, iter и прочие аналоги на коллекциях

<sup>2</sup> Современные GPGPU стараются бороться с этим ограничением

# CPU

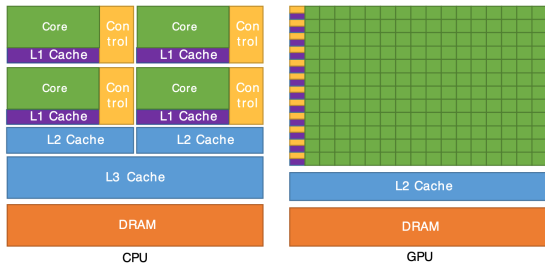
- Процессор общего назначения
- Независимые задачи
- Много сложных кэшей
- Сложный поток управления

# GPGPU

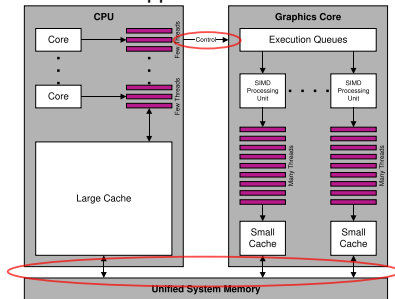
- Специализированный (co)процессор
- Параллелизм по данным
- Высокая пропускная способность
- Большое количество простых ALU

Типичный «взгляд издалека»

Сравнение архитектур<sup>3</sup>



Взаимодействие в SoC<sup>4</sup>



<sup>3</sup><https://cvw.cac.cornell.edu/gpu-architecture/gpu-characteristics/design>

<sup>4</sup>PowerVR Hardware. Architecture Overview for Developers

- Масса вендоров:

- ▶ Nvidia
- ▶ AMD
- ▶ Intel
- ▶ Imagination Technologies
- ▶ Qualcomm
- ▶ ...

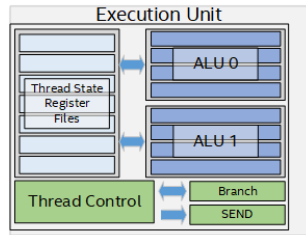
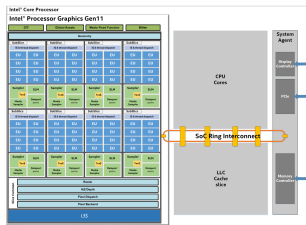
- Микроархитектура различных графических ускорителей сильно различается
- Архитектура различных графических ускорителей сильно различается
- Терминология сильно различается

Nvidia	AMD	Intel	Im.Tech.
Streaming Processor	SIMD lane	Processing element	ALU Group
SIMT unit	SIMD unit	Vector engine	ALU Pipeline
Streaming Multiprocessor	Computing Unit	Execution Unit	Unified Shading Cluster
GPU processing clusters	Compute Engine	Xe-slice	Core

---

<sup>5</sup> как всегда, в деталях

# Интегрированная графика от Intel<sup>6,7</sup>



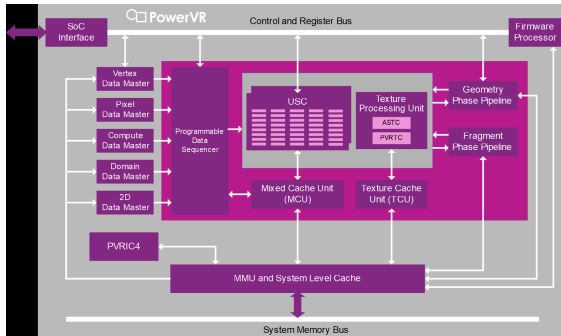
- Много специфики по обработке изображений
- Иерархическая группировка элементов
- Общий с процессором кэш

<sup>6</sup> Древняя, Gen11

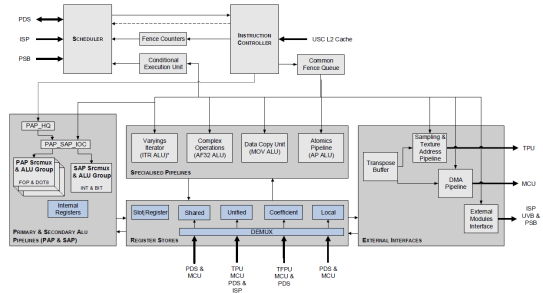
<sup>7</sup> Изображения из «Intel® Processor Graphics Gen11 Architecture»



# Графический ускоритель от Imagination Technologies (PowerVR)<sup>8,9</sup>



Общий вид

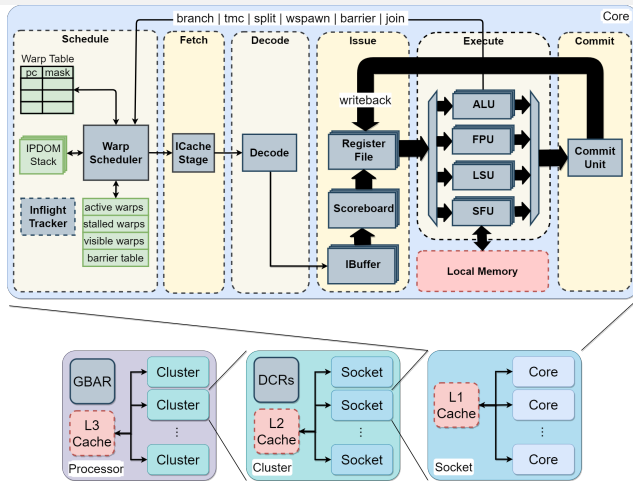


Детали Unified Shading Cluster

<sup>8</sup>Примерно то, с чем некоторые столкнутся на школе

<sup>9</sup>Изображения из документации: <https://docs.imgtec.com/performance-guides/compute-recommendations/html/topics/architecture/overview.html>

# Микроархитектура RISC-V GPGPU Vortex<sup>10,11</sup>



<sup>10</sup>Примерно то, с чем другие некоторые столкнутся на школе

<sup>11</sup>Изображение со страницы проекта:

<https://github.com/vortexgpgpu/vortex/blob/master/docs/microarchitecture.md>

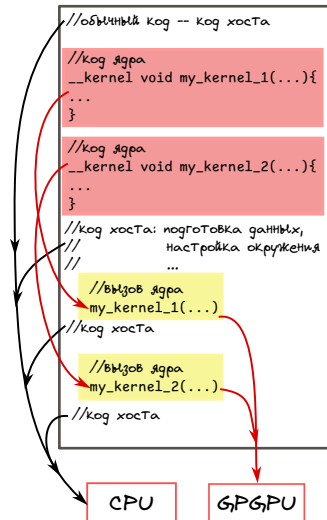
- Khronos<sup>12</sup>: «Open Standard for **Parallel** Programming of **Heterogeneous** Systems»
  - ▶ Многоядерные CPU
  - ▶ **Графические ускорители, GPGPU**
  - ▶ Программируемые логические интегральные схемы, ПЛИС/FPGA
  - ▶ Процессоры для цифровой обработки сигналов, DSP
  - ▶ ...
- Активно разрабатывается и используется
  - ▶ Intel, Google, Nvidia, AMD, Imagination Technologies, ...
  - ▶ OpenCV, FFmpeg, ViennaCL, GNU Octave, Matlab, ...
- Версии стандарта
  - ▶ 1.2: 2011 год, обязательная база для 3.0
  - ▶ 2.x: 2013 год, «откатали» в 3.0
  - ▶ 3.0: 2020 год, 3.0.17 в октябре 2024

---

<sup>12</sup>Консорциум для работы над открытыми стандартами в области параллельных вычислений, машинного обучения, компьютерной графики: больше 180 организаций, 21 стандарт.  
<https://www.khronos.org/>.

# OpenCL: структура программы

- Часть кода выполняется на центральном процессоре, часть на ускорителе
- Центральный процессор («хост», host)
  - ▶ Взаимодействие с «внешним миром», подготовка данных, запуск кода на ускорителе
  - ▶ Хочется писать код на «типичном» языке программирования
- Ускоритель («устройство», device)
  - ▶ Обработка данных «по указке» центрального процессора
  - ▶ Придётся пользоваться специализированным языком
  - ▶ Процедура, предназначенная для выполнения на ускорителе — «ядро», **kernel**<sup>13</sup>



<sup>13</sup>OpenCL kernel, CUDA kernel, ...

# OpenCL: взаимодействие хоста и устройства

- Обнаружение, подключение, конфигурация устройств<sup>14</sup>
- Управление памятью<sup>15</sup>
  - ▶ Выделение/освобождение
  - ▶ Контроль доступа и синхронизации между хостом и устройствами
- Выполнение команд
  - ▶ Передача данных, запуск ядер, синхронизация, ...
  - ▶ Основной интерфейс взаимодействия — **очередь команд**<sup>16</sup>
    - ★ Неблокирующая: центральный процессор ставит задачи, но не дожидается их исполнения
    - ★ Нужны дополнительные действия чтобы узнать, когда закончилась определённая задача
    - ★ In-order: команды выполняются строго друг за другом
    - ★ Может быть out-of-order
    - ★ Синхронизация между несколькими очередями возможна, но требует отелльной работы

---

<sup>14</sup>ICD-Loader: <https://github.com/KhronosGroup/OpenCL-ICD-Loader>

<sup>15</sup>Основной примитив — буфер:

[https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_API.html#\\_buffer\\_objects](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_API.html#_buffer_objects)

<sup>16</sup>Command queue:

[https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_API.html#\\_command\\_queues](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_API.html#_command_queues)

# Язык для написания ядер: OpenCL C/C++

- C++ for OpenCL language<sup>17</sup>
- OpenCL C language<sup>18</sup>
  - ▶ Диалект C99: специальные типы данных<sup>19</sup>, модификаторы памяти<sup>20</sup>
  - ▶ JIT-компиляция<sup>21</sup>, но бывают и AOT решения
  - ▶ Дополнительные библиотеки времени выполнения
  - ▶ Расширения<sup>22</sup>
    - ★ Предусмотренные стандартом<sup>23</sup>: `cl_khr_fp64`, `cl_khr_int64_extended_atomics`, ...
    - ★ Расширения от вендоров<sup>24</sup>: `cl_nv_pragma_unroll`, `cl_amd_fp64`, ...

---

<sup>17</sup> В 2.x был OpenCL C++ Kernel Language

<sup>18</sup> Спецификация: [https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_C.html](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_C.html)

<sup>19</sup> [https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_C.html#supported-data-types](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_C.html#supported-data-types)

<sup>20</sup> [https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_C.html#address-space-qualifiers](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_C.html#address-space-qualifiers)

<sup>21</sup> Для переносимости: вендору целевого устройства «лучше знать», как компилировать

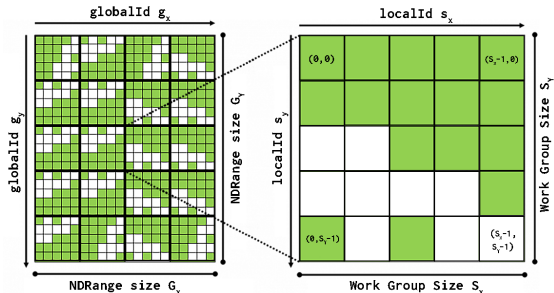
<sup>22</sup> Знакомая ситуация, не правда ли?

<sup>23</sup> Поискать «Khronos Extension Specifications» на <https://registry.khronos.org/OpenCL/>

<sup>24</sup> Поискать «Vendor and Multi-Vendor Extension Specifications» на <https://registry.khronos.org/OpenCL/>

# OpenCL C: параллельность

- Виртуальная решётка (NDRange)  
«потоков» (**work-items**)
  - одно-, двух-, трёхмерная
  - Группировка узлов в рабочие группы (**work-group**) фиксированного размера
  - Возможность получить глобальные и локальные координаты «потока»
- Параметры (размеры) решётки связаны с «размерами» обрабатываемых данных



Пример двухмерной виртуальной решётки<sup>25</sup>

- $globalId$  — глобальные координаты
- $localId$  — локальные координаты

<sup>25</sup>Изображение из поста: <http://thebeardsage.com/opencl-kernel-space/>

## И снова проблемы<sup>29</sup>

- ? Как отобразить абстракции OpenCL в реальное устройство?<sup>26</sup>
- ? Как писать код, чтобы он показывал хорошую<sup>27</sup> производительность?
- ☹ Придётся писать много вендору/платформу/device-специфичного кода.
  - 😊 Кое-что полезное об устройстве можно узнать во время выполнения<sup>28</sup>
  - ☹ От перебора вороха вариантов в коде всё равно не уйти
- ! Есть некоторые общие принципы

---

<sup>26</sup> При таком-то зоопарке устройств/микроархитектур

<sup>27</sup> Хотелось бы, конечно, максимальную но ...

<sup>28</sup> Стандартный API для получения данных об устройстве: [https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL\\_API.html#platform-querying-devices](https://registry.khronos.org/OpenCL/specs/3.0-unified/html/OpenCL_API.html#platform-querying-devices)

<sup>29</sup> «Шо, опять?»



# Общие принципы работы с очередью команд

Помни!

- Постановка задачи в очередь — накладные расходы
- Разные задачи в очереди задействуют разные ресурсы
- Одна задача может не утилизировать какой-то ресурс полностью
- Ожидание завершения задачи может привести к потере вычислительного времени

Потому

- Не дробь ядра
- Используй разные очереди или out-of-order<sup>30</sup>
- Ставь в очередь как можно больше задач<sup>31</sup>
- Пользуйся тем, что очередь неблокирующая
- Уменьшай количество точек синхронизации

---

<sup>30</sup> Там, где это возможно

<sup>31</sup> В разумных пределах

# Общие принципы работы с памятью

Помни!

- Передача данных<sup>32</sup> — накладные расходы
- Память бывает разная: объём и время доступа различаются на порядки
- Последовательный доступ — хорошо, произвольный — плохо

Потому

- Помни что где лежит и клади нужное поближе<sup>33</sup>
- Думай над тем, как представить данные в памяти
- Структурируй код ядер так, чтобы минимизировать непоследовательный доступ к памяти

---

<sup>32</sup>Как между хостом и устройством, так и между устройствами

<sup>33</sup>Изучи модификаторы памяти в OpenCL C и настройки буферов

# Общие принципы написания ядер

Помни!

- Ветвления<sup>34</sup> в ядрах — зло<sup>35</sup>
  - ▶ Соответствующее неприятное явление называется **thread divergence** или **branch divergence**
- На аппаратном уровне вполне определённая<sup>36</sup> группа «потоков» гарантированно выполняет одну инструкцию и рабочая группа как правило отображается в неё

Потому

- Старайся писать как можно более линейный код
  - ▶ Иногда ради этого приходится писать несколько ядер
- Подбирай размер рабочей группы так, чтобы максимально использовать аппаратные ресурсы

---

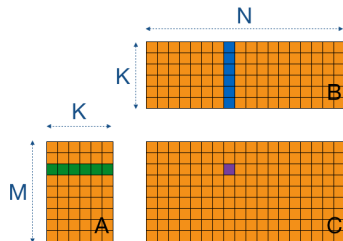
<sup>34</sup>И речь не только об if-then-else

<sup>35</sup>К сожалению, иногда — необходимое зло

<sup>36</sup>В зависимости от вендора: Task, Thread Group, Warp, Wavefront, ...

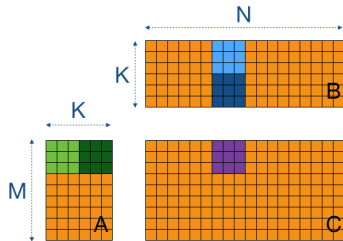
# Классический пример<sup>37</sup>

Наивная реализация



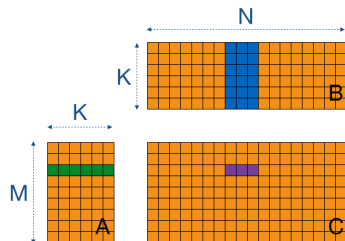
139 GFLOPS

Локальная память



373 GFLOPS

Регистры



689 GFLOPS

$$\begin{matrix} \text{purple} \\ \text{block} \end{matrix} = \begin{matrix} \text{green} \\ \text{block} \end{matrix} \times \begin{matrix} \text{light blue} \\ \text{block} \end{matrix} + \begin{matrix} \text{dark green} \\ \text{block} \end{matrix} \times \begin{matrix} \text{dark blue} \\ \text{block} \end{matrix}$$

- cuBLAS в тех же условиях: 3056 GFLOPS

<sup>37</sup>Материалы из tutorials <https://cnugteren.github.io/tutorial/pages/page1.html>.  
Соответствующий код: <https://github.com/cnugteren/myGEMM>

- Khronos OpenCL
  - ▶ <https://www.khronos.org/opencvl/>
- Курс по разработке под GPGPU от EuroCC National Competence Center Sweden
  - ▶ <https://enccs.github.io/gpu-programming/>
- OpenCL для ускорителей от Imagination Technologies
  - ▶ <https://blog.imaginationtech.com/a-quick-guide-to-writing-openccl-kernels-for-rogue/>
- Рекомендации разработчику под ускорители от Imagination Technologies
  - ▶ <https://docs.imgtec.com/performance-guides/compute-recommendations/html/index.html>
- Cornell University, Understanding GPU Architecture
  - ▶ <https://cvw.cac.cornell.edu/gpu-architecture>
- Антонюк В. А. OpenCL Открытый язык для параллельный программ, учебное пособие, 2017.
  - ▶ <https://istina.msu.ru/publications/book/76785028/>