



Казанский  
федеральный  
университет

ВЫСШАЯ ШКОЛА  
информационных технологий  
и информационных систем

# Программная модель CUDA

# Компиляция

- ▶ nvcc – строковый компилятор для CUDA
- ▶ Исходные файлы – \*.cu
- ▶ Работа с компилятором nvcc идентична gcc (RTFM)
- ▶ Пример компиляции исходного файла:
  - `$nvcc -arch=sm_37 -O3 test.cu -o test`



# Программная модель

- ▶ Программная модель CUDA является гетерогенной – использует и GPU и CPU
- ▶ В CUDA выделяют
  - Host (хост) – обозначает CPU и оперативную память
  - Device – обозначает GPU и графическую память
  - Kernels (ядра) – функции исполняемые только на GPU



# Программная модель

- ▶ Типичная последовательность команд CUDA-программы
  - Объявить и выделить память на хосте и девайсе
  - Инициализировать данные на хосте
  - Скопировать данные с хоста на девайс
  - Запустить ядро(а) на исполнение
  - Скопировать память обратно на хост
  - Освободить выделенную память

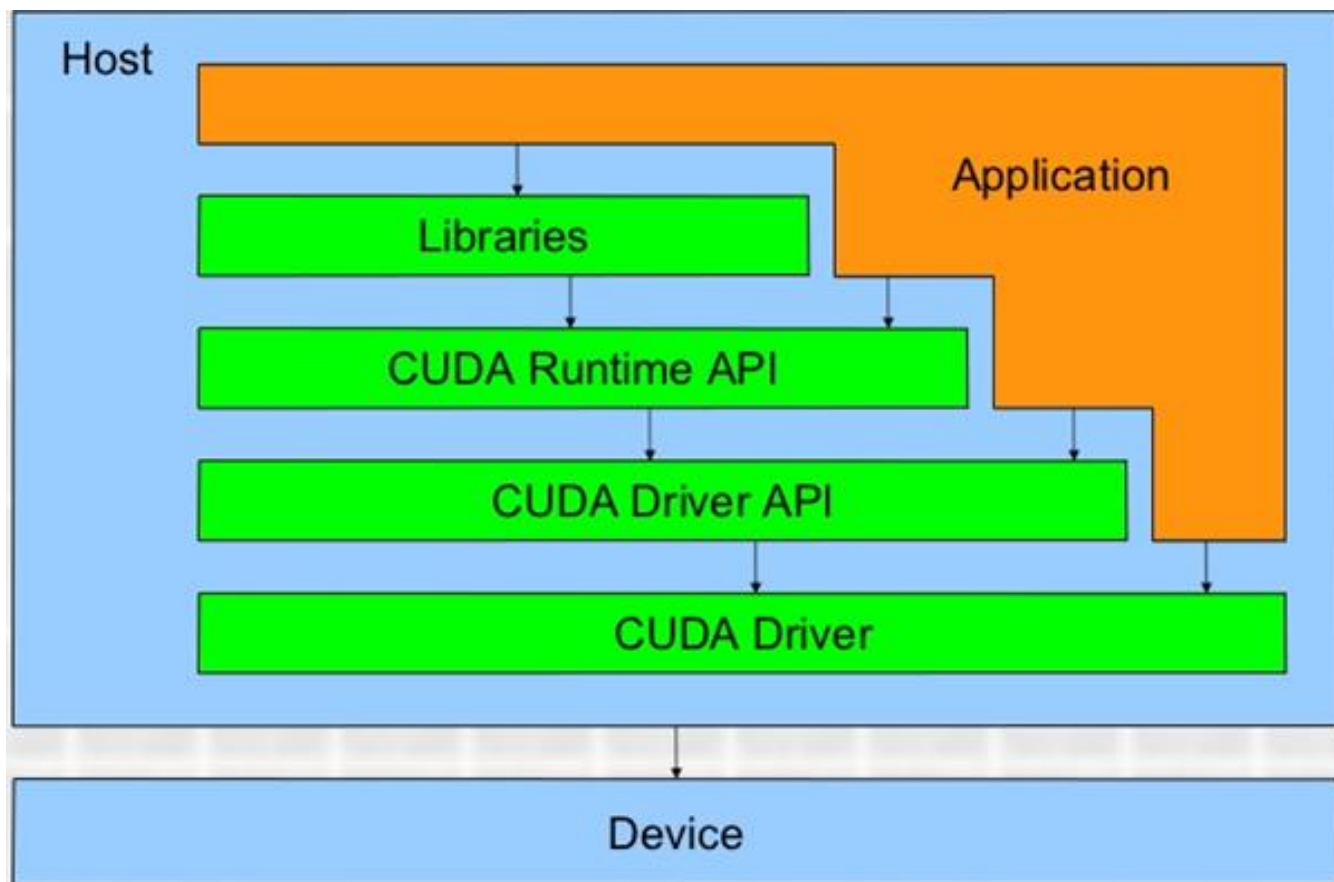


# Программная модель

- ▶ Код на хосте управляет памятью и хоста и девайса, а также запускает ядра для выполнения на девайсе
- ▶ Ядра исполняются на девайсе параллельно множеством нитей (threads)
- ▶ Каждая нить выполняет один и тот же код функции-ядра
- ▶ Как определить, какую часть данных должна обрабатывать каждая нить?



# Программный стек CUDA



# Сетки, блоки и нити

- ▶ Нити организованы в иерархию
  - Нити (threads)
  - Блоки нитей (blocks of threads)
  - Сетка блоков (grid of blocks)
- ▶ При запуске ядра обязательно указывается два параметра
  - Количество нитей в блоке
  - Количество блоков



# Сетки, блоки и нити

- ▶ Нити в блоке и блоки в сетке могут быть организованы в 1-, 2- и 3-х мерные массивы
- ▶ Каждая нить хранит в специальных переменных следующие значения
  - Свой номер в блоке
  - Номер своего блока
  - Количество блоков в сетке
- ▶ Таким образом нить может вычислить свой уникальный порядковый номер





# Сетки, блоки и нити

- ▶ Каждый из блоков сетки выполняется на одном SM
- ▶ В каждый момент времени на SM может выполняться лишь один блок
- ▶ Распределение блоков по SM происходит динамически, управляется логикой GPU
- ▶ Коммуникация и синхронизация нитей возможна лишь в пределах одного блока



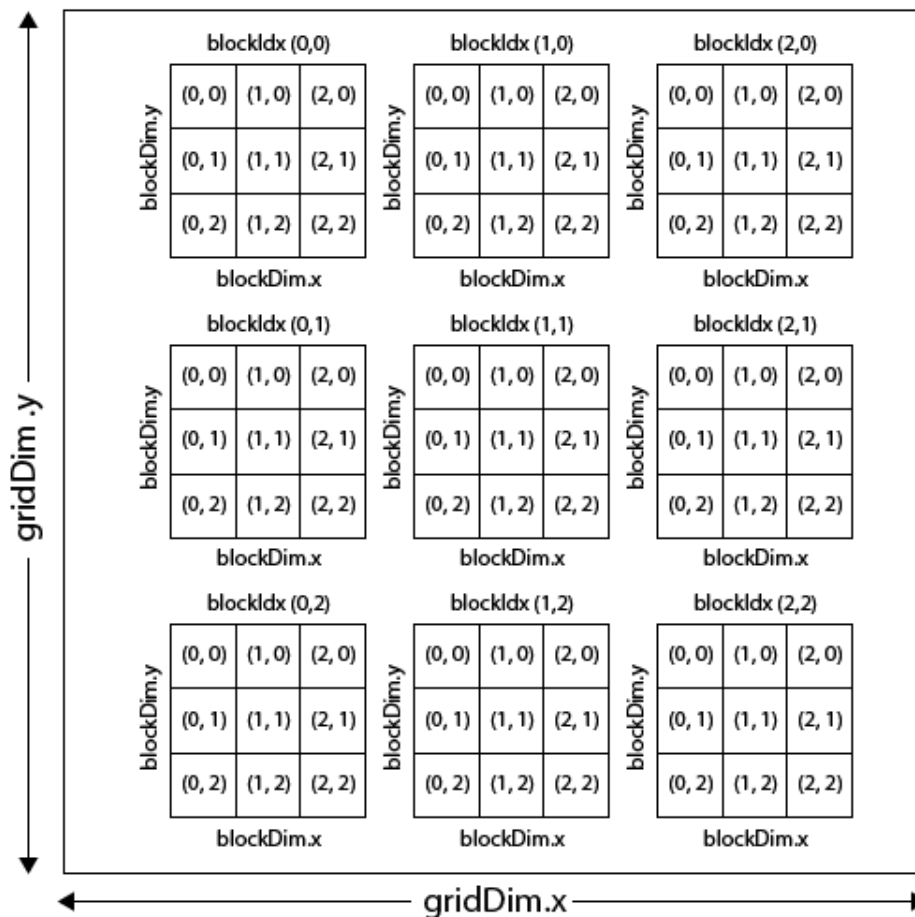
# Сетки, блоки и нити

- ▶ Нити блока делятся на логические группы по 32 нити – варпы (warps)
- ▶ Только нити одного варпа исполняются физически одновременно
- ▶ Распределение нитей по варпам происходит автоматически без участия программиста
- ▶ Если  $n_{threads\_in\_block} \% 32 \neq 0$  один из варпов будет содержать неактивные нити
- ▶  $MAX(n_{threads\_in\_block}) = 1024$



# Сетки, блоки и нити

## CUDA Grid



# Hello, world feat. CUDA

- ▶ Классический пример hello, world для CUDA – вычисление SAXPY
- ▶ SAXPY – Single Precision A \* X Plus Y
- ▶ Поэлементное сложение векторов, один из которых умножен на скаляр

$$\vec{z} = \alpha \vec{x} + \vec{y}$$



# Расширения языка

- ▶ Атрибуты функций
- ▶ Атрибуты переменных
- ▶ Встроенные переменные
- ▶ Дополнительные типы данных
- ▶ Оператор запуска ядра



# Атрибуты функций

- ▶ Атрибут `__global__`
  - Обозначает функцию-ядро - вызывается с хоста исполняется на GPU
  - Ядро должно возвращать значение типа `void`
- ▶ Атрибут `__device__`
  - Обозначает функцию которая вызывается из ядра или другой `device` функции
  - Исполняется на GPU
- ▶ Атрибут `__host__`
  - Обычная функция C/C++



# Атрибуты переменных

- ▶ Атрибуты переменных определяют в какой памяти GPU будет размещаться переменная
- ▶ Скорость и способ доступа к переменной сильно варьируется в зависимости от типа размещения
- ▶ Подробности в лекции о памяти в CUDA



# Встроенные типы

- ▶ 1/2/3/4-мерные векторные типы на основе *char*, *short*, *int*, *long*, *float* и *double*
- ▶ Компоненты векторных типов имеют имена *x*, *y*, *z* и *w*

```
int2 a = make_int2(1, 7); //Создает вектор (1, 7)
```

```
float3 u = make_float3(1, 2, 3.4f); //Создает вектор (1.0f,  
2.0f, 3.4f )
```





# Встроенные типы

- ▶ Для этих типов не поддерживаются автоматические векторные операции
- ▶ Тип *dim3*, используемый для задания размерностей блоков потоков и сеток блоков
- ▶ Этот тип основан на *uint3*, элементы по умолчанию инициализируются 1

```
dim3 blocks ( 16, 16 ); // то же что blocks ( 16, 16, 1 )  
dim3 grid ( 256 ); // то же что grid (256, 1, 1)
```



# Встроенные переменные

- ▶ *gridDim* – количество блоков в сетке по каждому из измерений (тип *dim3*)
- ▶ *blockDim* – количество нитей в блоке по каждому из измерений (тип *dim3*)
- ▶ *blockIdx* – индекс текущего блока в сетке по каждому из измерений (тип *uint3*)
- ▶ *threadIdx* – индекс текущей нити в блоке по каждому из измерений (тип *uint3*)
- ▶ *warpSize* – размер варпа (тип *int*)



# Оператор вызова ядра

```
kernel_name<<<  
    dim3 Nb,  
    dim3 Nt,  
    size_t Ns = 0,  
    cudaStream_t S = 0>>>(args)
```



# Оператор вызова ядра

- ▶ `kernel_name` – это имя или адрес ядра
- ▶ Обязательный параметр `Nb` типа `uint` задает число блоков в сетке блоков
- ▶ Обязательный параметр `Ns` задает число нитей в блоке
- ▶ Параметр `Ms` задает дополнительный объем памяти
- ▶ Параметр `S` ставит вызов ядра в заданную очередь команд `CUDA Stream`



# Обработка ошибок

- ▶ Каждая функция CUDA (кроме запуска ядра) возвращает значение типа *cudaError\_t*
- ▶ При успешном выполнении функции возвращается значение *cudaSuccess*, иначе возвращается код ошибки

```
char* cudaGetErrorString(cudaError_t code); //описание ошибки  
cudaError_t cudaGetLastError(); //получить последнюю ошибку
```



# Асинхронность в CUDA

- ▶ Некоторые функции CUDA являются асинхронными
  - Запуск ядра
  - Асинхронные версии функций копирования и инициализации памяти
  - Функции копирования памяти device ↔ device внутри устройства и между устройствами
- ▶ Для синхронизации устройства используется функция *cudaDeviceSynchronize()*



# CUDA Streams

- ▶ Очередь команд для запуска ядер и обмена данными между GPU и хостом
- ▶ По умолчанию все команды работают в очереди №0
- ▶ Команды с одним значением stream выполняются последовательно, с разными – независимо друг от друга
- ▶ Очереди команд используются для оптимизации работы с памятью



# CUDA events

- ▶ Событие – объект типа `cudaEvent_t` для измерения времени выполнения операций CUDA
- ▶ Функции CUDA API позволяют
  - Создавать и уничтожать события
  - Выделять начало и конец профилируемого кода
  - Проверять/ожидать наступления события
  - Замерять время выполнения







Казанский федеральный  
УНИВЕРСИТЕТ

ВЫСШАЯ ШКОЛА  
информационных технологий  
и информационных систем

# Вопросы

ekhramch@kpfu.ru