

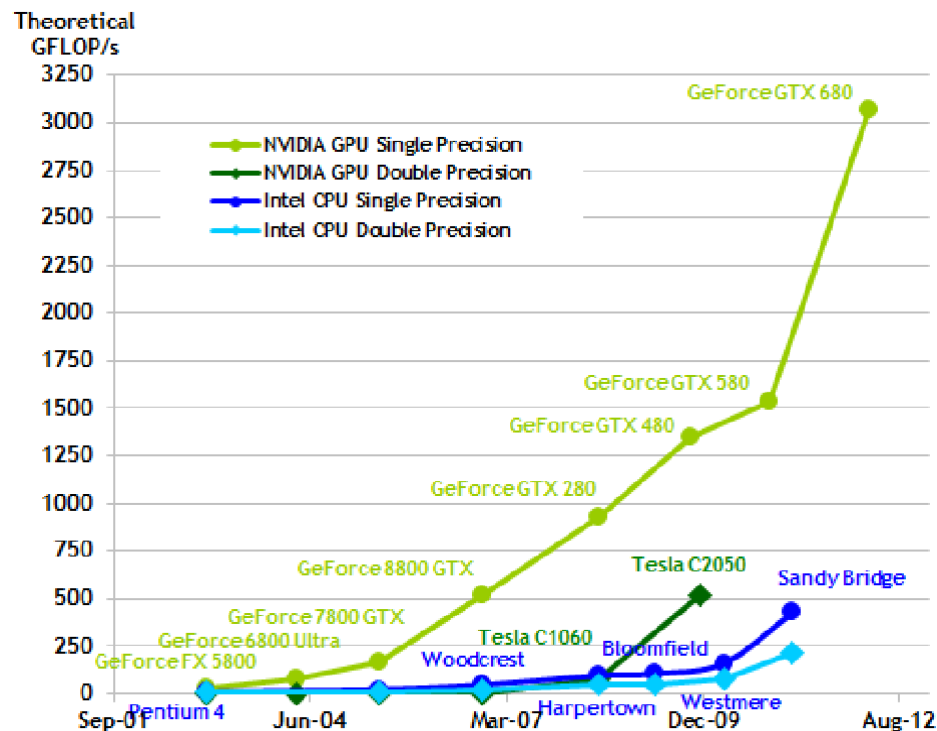
Лекция 1

Содержание лекции

- линейка графических процессоров NVIDIA;
- аппаратные особенности графических процессоров;
- архитектура CUDA – основные свойства и принципы;
- программная модель: хост, устройства, ядра, иерархия нитей (*threads*);
- установка CUDA Toolkit

Рост производительности GPU (август 2012)

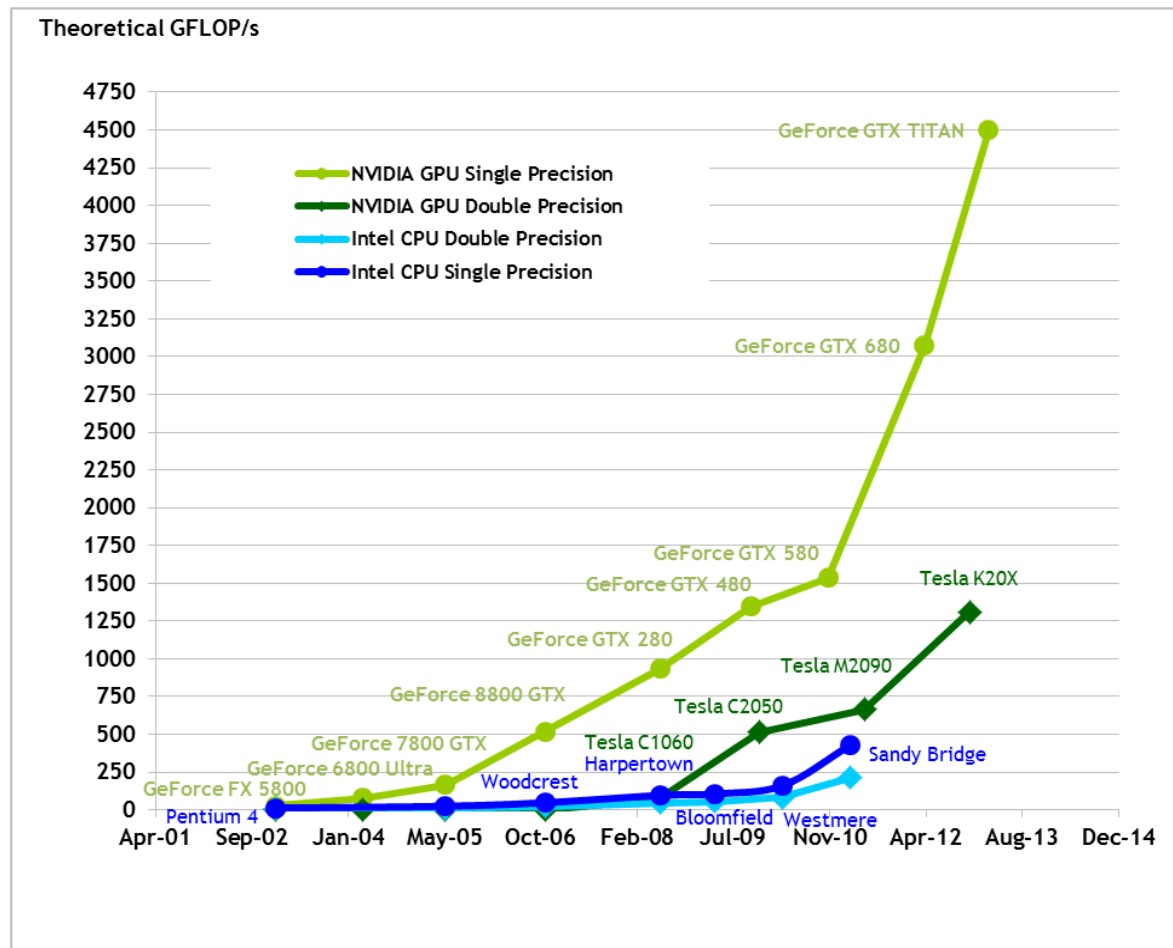
<http://www.nvidia.ru/object/cuda-parallel-computing-ru.html>



Intel(R) Core(TM) i5-2500K CPU @ 3.30GH - **100Gflop/s**
GeForce GTX 560 Ti - **1260 Gflop/s**

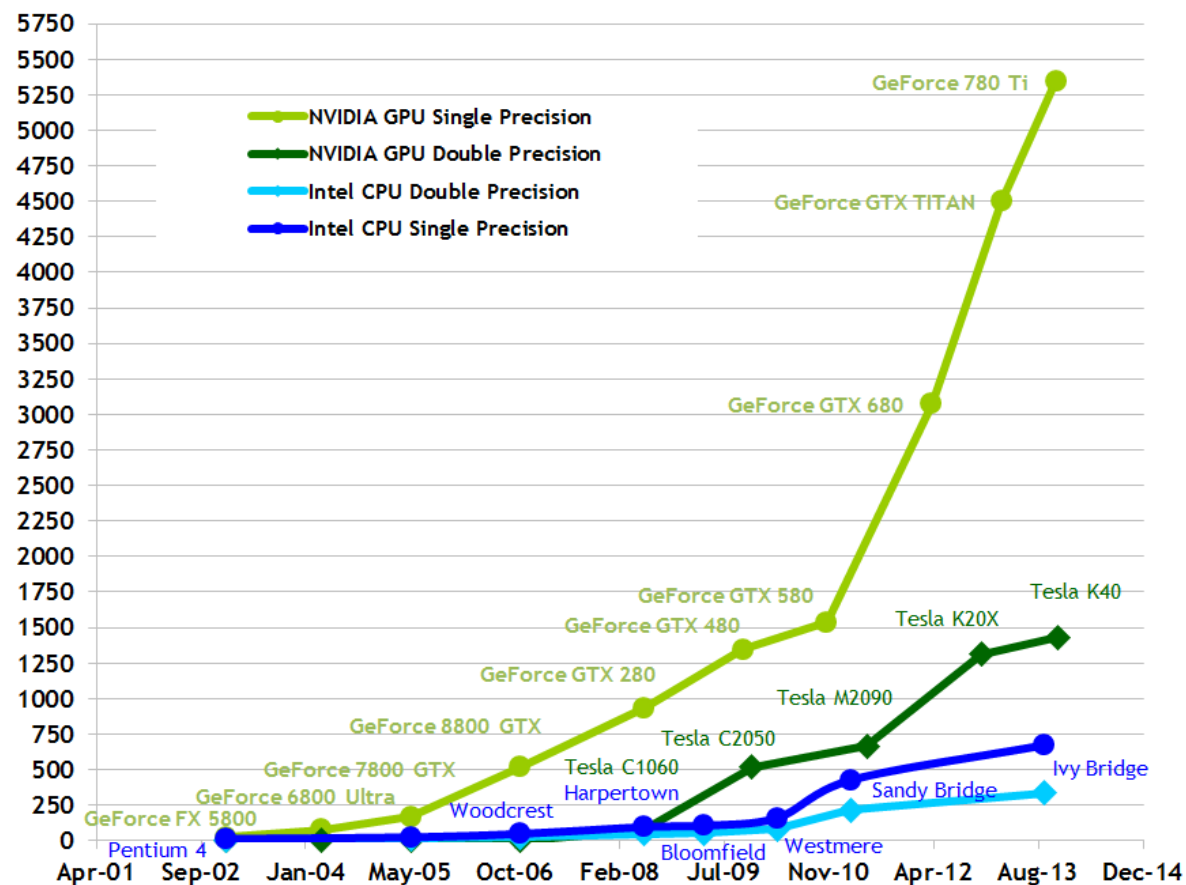
Производительность CPU и GPU.

Рост производительности GPU (август 2013)

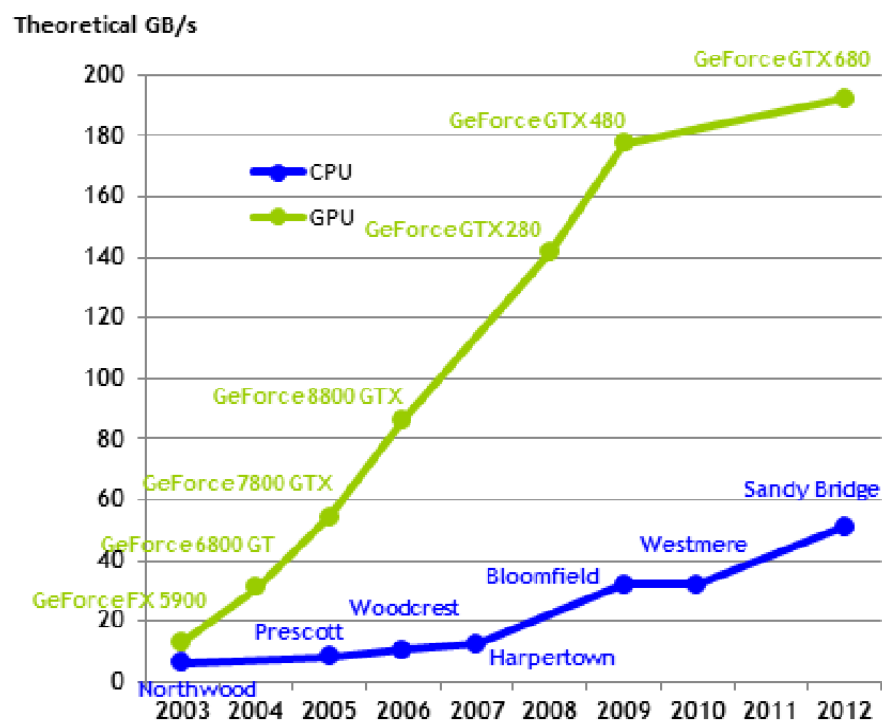


Рост производительности GPU (декабрь 2014)

Theoretical GFLOP/s



Рост пропускной способности GPU



Пропускная способность CPU и GPU.

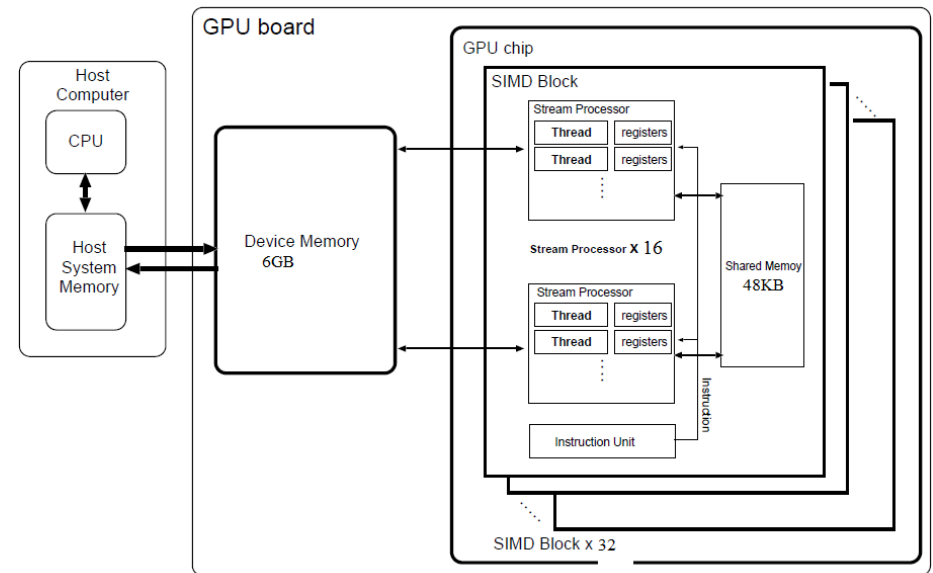
PCI Express 1.0	4Gb
PCI Express 2.0	16Gb
PCI Express 3.0	25.6Gb

Пропускная способность внешнего интерфейса.

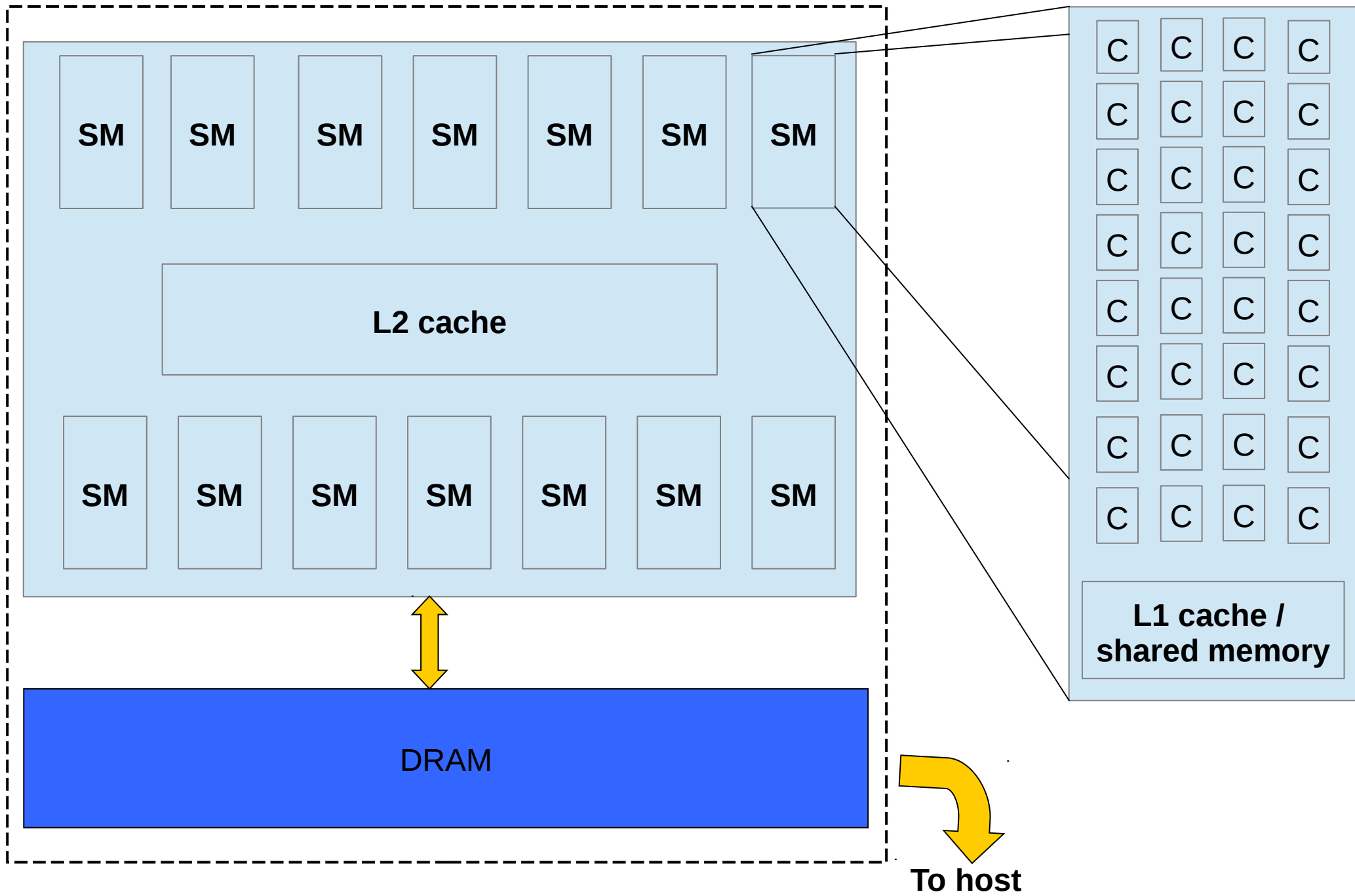
Физическое представление GPU

В общих чертах графический процессор состоит из набора SIMD (*Single Instruction – Multiple Data*) блоков или multiprocessors, каждый из которых включает в себя определенное количество stream processors (обычно 8-16), несколько ALU (Arithmetic Logical Unit) и единственный instruction unit. Микросхема SIMD блока также содержит небольшую быструю память, выполняющую роль кэша – shared memory, и небольшую память для чтения – кэши константной и текстурной памяти. Каждый из stream processors имеет набор регистров, куда загружаются данные, обрабатываемые одной инструкцией одновременно.

На плате GPU располагается микросхема памяти большого объема. GPU соединяется с узлом вычислительной системы посредством интерфейса PCI Express.



GPU (Fermi architecture)



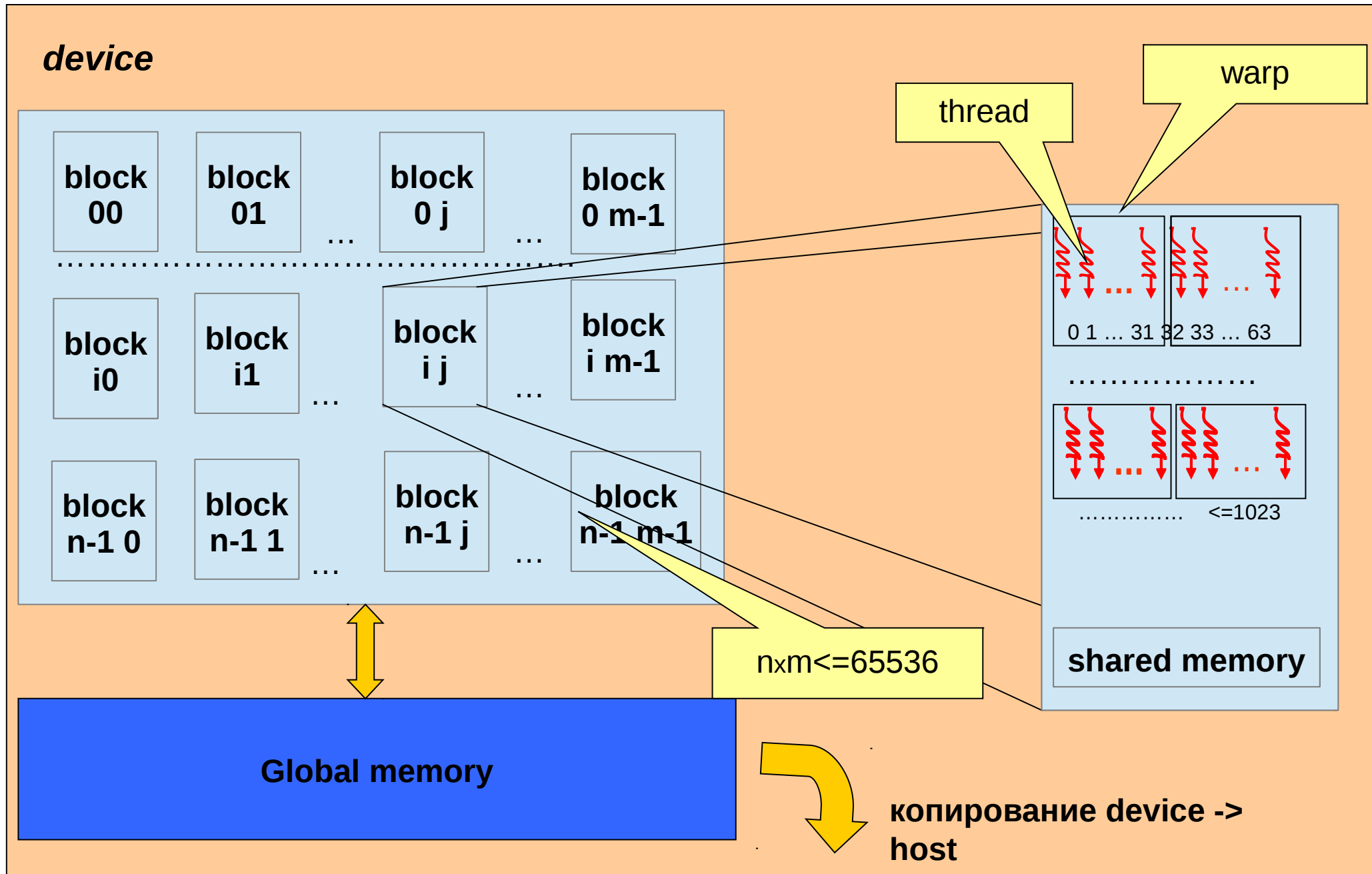
Логическое представление GPU

Активное использование графических процессоров (GPU) для прикладных расчетов научно-технического назначения во многом связано с предоставлением компанией NVIDIA технологии **CUDA** (*Cuda Unified Device Architecture*) (позже конкурентом CUDA выступила технология *OpenCL*, представляющая универсальный подход для программирования специализированных вычислительных устройств, и поддерживаемая компаниями *AMD*, *Intel* и *Nvidia*). Технология CUDA предоставляет понятную для прикладного программиста абстракцию графического процессора (GPU) и простой интерфейс прикладного программирования (*API – Application Programming Interface*).

По терминологии *CUDA* вычислительный узел с *CPU* и *main memory* называется **host**, *GPU* называется **device**. Программа, выполняемая на *host*'е содержит код – ядро (**kernel**), который загружается на *device* в виде многочисленных копий. Все копии загруженного кода – нити (**threads**), объединяются в блоки (**blocks**) по 512-1024 нити в каждом. Все блоки объединяются в сеть (**grid**) с максимальным количеством блоков 65536. Все нити имеют совместный доступ на запись/чтение к памяти большого объема - *global memory*, на чтение к кэшируемым **constant memory** и **texture memory**. Нити одного блока имеют доступ к быстрой памяти небольшого объема – **shared memory**.

CUDA (Compute Unified Device Architecture)

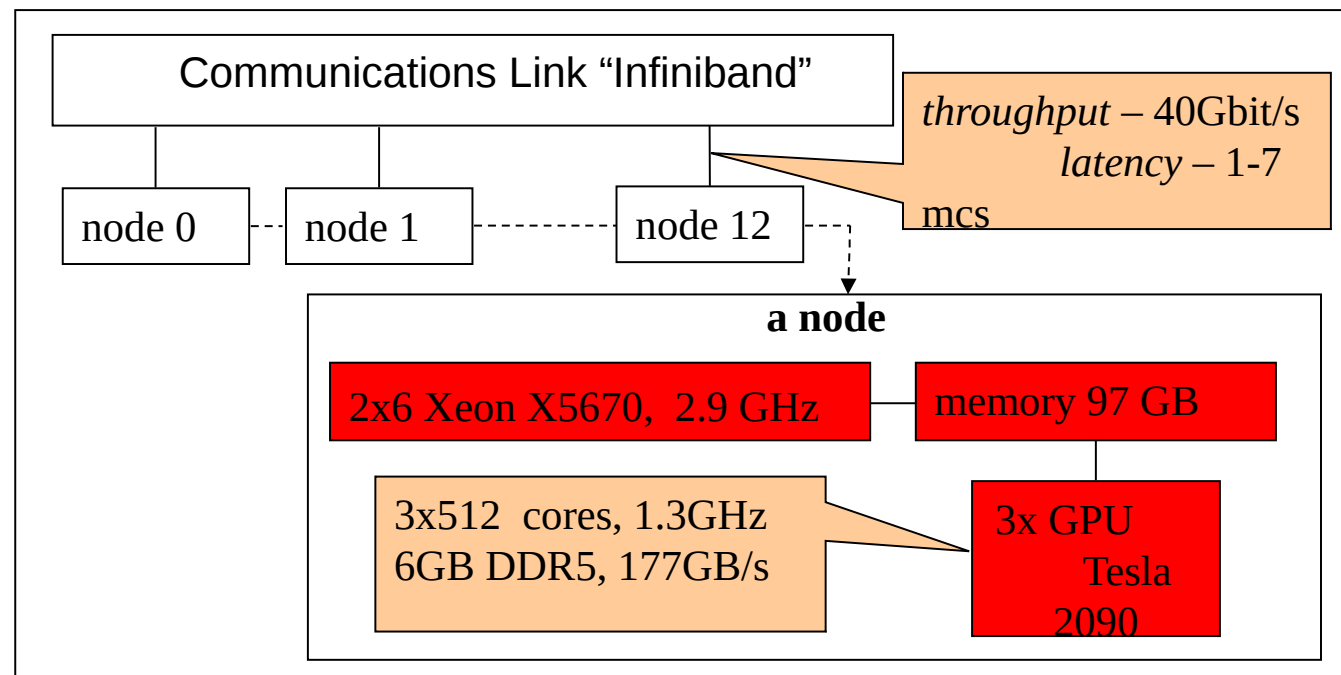
- cuda предоставляет абстракцию GPU для программистов



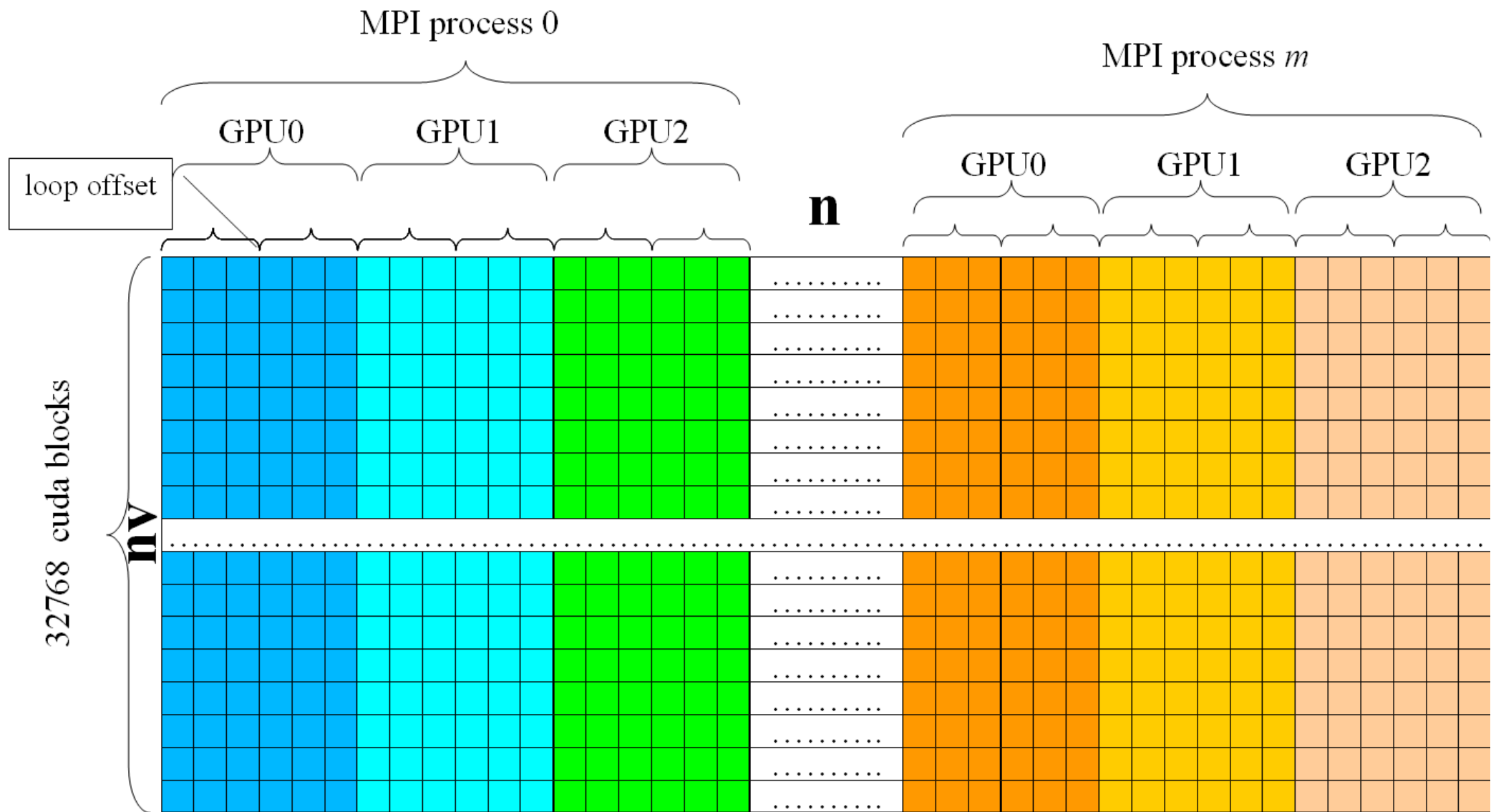
Гибридный кластер НГУ (<http://nusc.ru>)

Кластер состоит из 12 узлов *HP SL390s G7*, каждый из которых содержит два 6-ядерных CPU *Xeon X5670* и три графические карты *NVIDIA Tesla M2090*.

Каждый GPU имеет **512 ядер** (cores) с частотой 1.3GHz и память размером **6GB** с пропускной способностью (bandwidth) 177 GB/s.



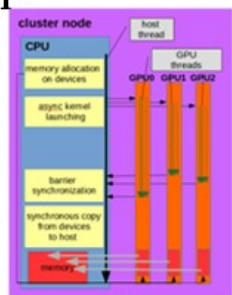
decomposition



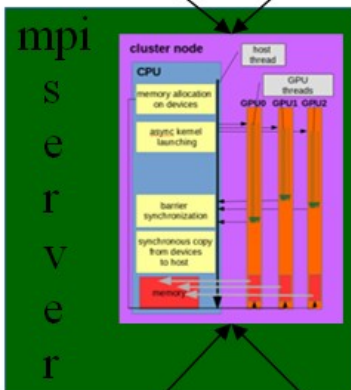
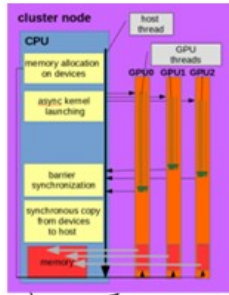
Typical value of *loop offset* – from 32 to 512, correspond to number of cuda threads in a cuda block.

select=any:ncpus=1:
mpiprocs=1:gpus=1..3

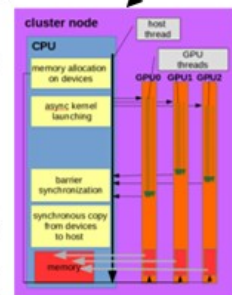
mpi
client



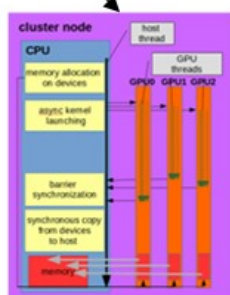
mpi
client



mpi
client



mpi
client



cluster node

CPU

memory allocation
on devices

async kernels
launching

barrier
synchronization

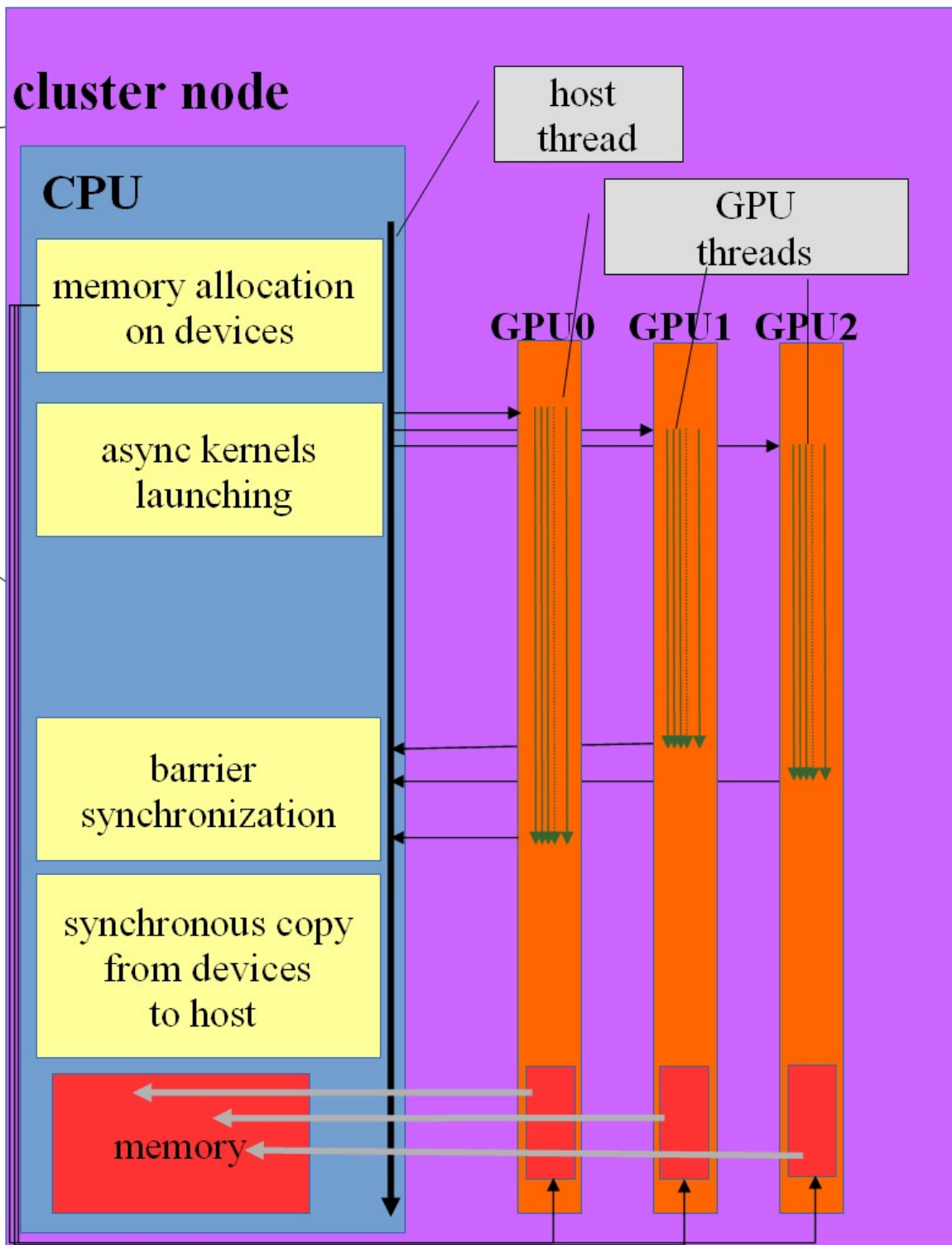
synchronous copy
from devices
to host

memory

host
thread

GPU
threads

GPU0 GPU1 GPU2



(one gpu) x (six loops)

NG=1 gCalculateParams Elapsed time: 1.55024

gClearStoss gClearStoss Elapsed time: 0.213664

gTransposeDf Elapsed time: 5.2569

gStossCalc Elapsed time: **246532**

gTransposeDfInverse Elapsed time: 7.51043

gClearStoss gClearStoss Elapsed time: 0.217504

gTransposeDf Elapsed time: 5.25312

gStossCalc Elapsed time: **246536**

gTransposeDfInverse Elapsed time: 7.53414

gClearStoss gClearStoss Elapsed time: 0.217088

gTransposeDf Elapsed time: 5.25322

gStossCalc Elapsed time: **246535**

gTransposeDfInverse Elapsed time: 7.51981

gClearStoss gClearStoss Elapsed time: 0.216192

gTransposeDf Elapsed time: 5.25834

gStossCalc Elapsed time: **246533**

gTransposeDfInverse Elapsed time: 7.46163

gClearStoss gClearStoss Elapsed time: 0.216544

gTransposeDf Elapsed time: 5.25523

gStossCalc Elapsed time: **246538**

gTransposeDfInverse Elapsed time: 7.49238

gClearStoss gClearStoss Elapsed time: 0.215904

gTransposeDf Elapsed time: 5.24982

gStossCalc Elapsed time: **246536**

gTransposeDfInverse Elapsed time: 7.52288

real 24m47.123s

user 24m45.149s

sys 0m0.944s

25,165,824
grid nodes
in phase space

every loop handles
4,194,304 grid nodes
in phase space

(two gpu) x (three loops)

NG=2 gCalculateParams Elapsed time: 1.53939

NG=2 gCalculateParams Elapsed time: 1.49997

gClearStoss gClearStoss Elapsed time: 0.211264

gClearStoss gClearStoss Elapsed time: 0.19552

gTransposeDf Elapsed time: 5.2217

gTransposeDf Elapsed time: 5.23258

gStossCalc Elapsed time: **248925**

gStossCalc Elapsed time: **248926**

gTransposeDfInverse Elapsed time: 7.48387

gTransposeDfInverse Elapsed time: 7.48378

gClearStoss gClearStoss Elapsed time: 0.200736

gClearStoss gClearStoss Elapsed time: 0.209472

gTransposeDf Elapsed time: 5.2145

gTransposeDf Elapsed time: 5.22477

gStossCalc Elapsed time: **248722**

gStossCalc Elapsed time: **248723**

gTransposeDfInverse Elapsed time: 7.51037

gTransposeDfInverse Elapsed time: 7.51152

gClearStoss gClearStoss Elapsed time: 0.199392

gClearStoss gClearStoss Elapsed time: 0.208352

gTransposeDf Elapsed time: 5.21789

gTransposeDf Elapsed time: 5.22787

gStossCalc Elapsed time: **248926**

gStossCalc Elapsed time: **248926**

gTransposeDfInverse Elapsed time: 7.51379

gTransposeDfInverse Elapsed time: 7.51504

real 12m35.874s

user 12m33.123s

25,165,824
grid nodes
in phase

space

every loop handles
4,194,304 grid nodes
in phase space

two-times speedup
vs. 1 GPU

(three gpu) x (two loops)

NG=3	gCalculateParams	Elapsed time: 1.47718
NG=3	gCalculateParams	Elapsed time: 1.53101
NG=3	gCalculateParams	Elapsed time: 1.5304
gClearStoss	gClearStoss	Elapsed time: 0.140192
gClearStoss	gClearStoss	Elapsed time: 0.155264
gClearStoss	gClearStoss	Elapsed time: 0.124384
gTransposeDf		Elapsed time: 4.94176
gTransposeDf		Elapsed time: 5.1377
gTransposeDf		Elapsed time: 5.14787
gStossCalc		Elapsed time: 164576
gStossCalc		Elapsed time: 164589
gStossCalc		Elapsed time: 164590
gTransposeDfInverse		Elapsed time: 7.12509
gTransposeDfInverse		Elapsed time: 7.43171
gTransposeDfInverse		Elapsed time: 7.44042
gClearStoss	gClearStoss	Elapsed time: 0.128416
gClearStoss	gClearStoss	Elapsed time: 0.13792
gClearStoss	gClearStoss	Elapsed time: 0.147904
gTransposeDf		Elapsed time: 4.93661
gTransposeDf		Elapsed time: 5.13046
gTransposeDf		Elapsed time: 5.13981
gStossCalc		Elapsed time: 164575
gStossCalc		Elapsed time: 164590
gStossCalc		Elapsed time: 164591
gTransposeDfInverse		Elapsed time: 7.11542
gTransposeDfInverse		Elapsed time: 7.44173
gTransposeDfInverse		Elapsed time: 7.4505
real		5m38.815s
user		5m36.017s

25,165,824
grid nodes
in phase

space

every loop handles
4,194,304 grid nodes
in phase space

superlinear speedup:
aprox 5-times (not 3 !) speedup
vs. 1 GPU

the former estimation of speedup

Parallel computations with one spatial cell ensure a **12-fold speedup as compared with sequential computations on a processor with a frequency of 2.66 GHz**. The results of relaxation computations with different numbers of spatial cells (threads in a block) performed on the GPGPU Tesla M2090 are summarized in the table:

Number of threads	Execution time (ms)	Speedup
1	676	1
2	523	2.59
4	462	5.85
8	416	13
16	349	30.99
32	523	41.36
64	1045	41.4

saturation

SPEEDUP: $12 \times 41 \times 5 = 2460$ (3 GPU on one node)

Compute capabilities (вычислительные возможности)

Compute capabilities (вычислительные возможности) представляют спецификацию GPU. Особенности “вычислительных возможностей” включают допустимость операций с плавающей точкой, допустимость атомарных операций, возможность синхронизации нитей, кэшируемость глобальной памяти и т.д. Описание различных версий Compute capabilities можно найти, например, в CUDA C Programming Guid – руководстве по CUDA C компании NVIDIA.

Архитектура GPU	Compute capabilities	Версия CUDA
Tesla	1.*	CUDA 2.*-3.*
Fermi	2.*	CUDA 4.*-5.*
Kepler	3.*	CUDA 5.*
Maxwell	5.*	CUDA 6.*-7.*
Pascal	6.*	CUDA 8
Volta	7.*	

Характеристики GPU на серверах 'home' и 'dew'

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GTX 560 Ti"

CUDA Driver Version / Runtime Version 5.0 / 5.0

CUDA Capability Major/Minor version number: 2.1

Total amount of global memory: 2048 MBytes (2147024896 bytes)

(8) Multiprocessors x (48) CUDA Cores/MP: 384 CUDA Cores

GPU Clock rate: 1645 MHz (1.64 GHz)

Memory Clock rate: 2004 Mhz

Memory Bus Width: 256-bit

L2 Cache Size: 524288 bytes

Max Texture Dimension Size (x,y,z) 1D=(65536), 2D=(65536,65535), 3D=(2048,2048,2048)

Max Layered Texture Size (dim) x layers 1D=(16384) x 2048, 2D=(16384,16384) x 2048

Total amount of constant memory: 65536 bytes

Total amount of shared memory per block: 49152 bytes

Total number of registers available per block: 32768

Warp size: 32

Maximum number of threads per multiprocessor: 1536

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x 1024 x 64

Maximum sizes of each dimension of a grid: 65535 x 65535 x 65535

Maximum memory pitch: 2147483647 bytes

Texture alignment: 512 bytes

Concurrent copy and kernel execution: Yes with 1 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

Device supports Unified Addressing (UVA): Yes

Device PCI Bus ID / PCI location ID: 1 / 0

Compute Mode:

< Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GTX 650"

CUDA Driver Version / Runtime Version 5.5 / 5.5

CUDA Capability Major/Minor version number: 3.0

Total amount of global memory: 2048 MBytes (2147155968 bytes)

(2) Multiprocessors x (192) CUDA Cores/MP: 384 CUDA Cores

GPU Clock rate: 1110 MHz (1.11 GHz)

Memory Clock rate: 2500 Mhz

Memory Bus Width: 128-bit

L2 Cache Size: 262144 bytes

Max Texture Dimension Size (x,y,z) 1D=(65536), 2D=(65536,65536), 3D=(4096,4096,4096)

Max Layered Texture Size (dim) x layers 1D=(16384) x 2048, 2D=(16384,16384) x 2048

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: 2048

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x 1024 x 64

Maximum sizes of each dimension of a grid: 65535 x 65535 x 65535

Maximum memory pitch: 2147483647 bytes

Texture alignment: 512 bytes

Concurrent copy and kernel execution: Yes with 1 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

Device supports Unified Addressing (UVA): Yes

Device PCI Bus ID / PCI location ID: 1 / 0

Compute Mode:

< Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

Характеристики GPU на гибридном кластере 'nusc.ru'

Found 3 CUDA Capable device(s)

Device 0: "Tesla M2090"

CUDA Driver Version / Runtime Version 4.2 / 4.1

CUDA Capability Major/Minor version number: 2.0

Total amount of global memory: 5375 MBytes (5636554752 bytes)

(16) Multiprocessors x (32) CUDA Cores/MP: 512 CUDA Cores

GPU Clock Speed: 1.30 GHz

Memory Clock rate: 1848.00 Mhz

Memory Bus Width: 384-bit

L2 Cache Size: 786432 bytes

Max Texture Dimension Size (x,y,z) 1D=(65536),

2D=(65536,65535), 3D=(2048,2048,2048)

Max Layered Texture Size (dim) x layers 1D=(16384) x 2048,

2D=(16384,16384) x 2048

Total amount of constant memory: 65536 bytes

Total amount of shared memory per block: 49152 bytes

Total number of registers available per block: 32768

Warp size: 32

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x 1024 x 64

Maximum sizes of each dimension of a grid: 65535 x 65535 x 65535

Maximum memory pitch: 2147483647 bytes

Texture alignment: 512 bytes

Concurrent copy and execution: Yes with 2 copy engine(s)

Run time limit on kernels: No

Integrated GPU sharing Host Memory: No

Support host page-locked memory mapping: Yes

Concurrent kernel execution: Yes

Alignment requirement for Surfaces: Yes

Device has ECC support enabled: Yes

Device is using TCC driver mode: No

Device supports Unified Addressing (UVA): Yes

Device PCI Bus ID / PCI location ID: 6 / 0

Compute Mode:

< Exclusive Process (many threads in one process is able to use

::cudaSetDevice() with this device) >

Device 1: "Tesla M2090"

CUDA Driver Version / Runtime Version 4.2 / 4.1

CUDA Capability Major/Minor version number: 2.0

Total amount of global memory: 5375

MBytes (5636554752 bytes)

(16) Multiprocessors x (32) CUDA Cores/MP: 512

CUDA Cores

GPU Clock Speed: 1.30 GHz

Memory Clock rate: 1848.00 Mhz

Memory Bus Width: 384-bit

L2 Cache Size: 786432 bytes

Max Texture Dimension Size (x,y,z)

1D=(65536), 2D=(65536,65535), 3D=(2048,2048,2048)

Max Layered Texture Size (dim) x layers

1D=(16384) x 2048, 2D=(16384,16384) x 2048

Total amount of constant memory: 65536

bytes

Total amount of shared memory per block: 49152

bytes

Total number of registers available per block: 32768

Warp size: 32

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x

1024 x 64

Maximum sizes of each dimension of a grid: 65535 x

65535 x 65535

.....

Device 2: "Tesla M2090"

CUDA Driver Version / Runtime Version 4.2 / 4.1

CUDA Capability Major/Minor version number: 2.0

Total amount of global memory: 5375

MBytes (5636554752 bytes)

(16) Multiprocessors x (32) CUDA Cores/MP: 512

CUDA Cores

GPU Clock Speed: 1.30 GHz

Memory Clock rate: 1848.00 Mhz

Memory Bus Width: 384-bit

L2 Cache Size: 786432 bytes

Max Texture Dimension Size (x,y,z)

1D=(65536), 2D=(65536,65535), 3D=(2048,2048,2048)

Max Layered Texture Size (dim) x layers

1D=(16384) x 2048, 2D=(16384,16384) x 2048

Total amount of constant memory: 65536

bytes

Total amount of shared memory per block: 49152

bytes

Total number of registers available per block: 32768

Warp size: 32

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x

1024 x 64

Maximum sizes of each dimension of a grid: 65535 x

65535 x 65535

.....

Установка CUDA Toolkit (docs.nvidia.com/cuda/)

LINUX: docs.nvidia.com/cuda/cuda-installation-guide-linux/index.html

rpm-накет

Click on the green buttons that describe your target platform. Only supported platforms will be shown.

Operating System	Windows	Linux	Mac OSX			
Architecture ⁱ	x86_64	ppc64le				
Distribution	Fedora	OpenSUSE	RHEL	CentOS	SLES	SteamOS
	Ubuntu					
Version	13.2					
Installer Type ⁱ	runfile (local)	rpm (local)	rpm (network)			

Download Target Installer for Linux OpenSUSE 13.2 x86_64

cuda-repo-opensuse132-7-5-local-7-5-18.x86_64.rpm (md5sum: e6708f4b38cc9ed546d3a04f7c731698)

Download (1.1 GB)

Installation Instructions:

1. ``sudo rpm -i cuda-repo-opensuse132-7-5-local-7-5-18.x86_64.rpm``
2. ``sudo zypper refresh``
3. ``sudo zypper install cuda``

For further information, see the [Installation Guide for Linux](#) and the [CUDA Quick Start Guide](#).

После установки драйвера: `sudo nvidia-xconfig`

run-файл

Click on the green buttons that describe your target platform. Only supported platforms will be shown.

Operating System	Windows	Linux	Mac OSX			
Architecture ⁱ	x86_64	ppc64le				
Distribution	Fedora	OpenSUSE	RHEL	CentOS	SLES	SteamOS
	Ubuntu					
Version	13.2					
Installer Type ⁱ	runfile (local)	rpm (local)	rpm (network)			

Download Target Installer for Linux OpenSUSE 13.2 x86_64

cuda_7.5.18_linux.run (md5sum: 4b3bcecf0dfc35928a0898793cf3e4c6)

Download (1.1 GB)

Installation Instructions:

1. Run ``sudo sh cuda_7.5.18_linux.run``
2. Follow the command-line prompts

The GPU Deployment Kit is available as a separate download [here](#).

For further information, see the [Installation Guide for Linux](#) and the [CUDA Quick Start Guide](#).

1. Выгрузить X-server (`sudo init 3`)
2. Блокировать загрузку драйвера **nouveau**

Установка CUDA Toolkit (docs.nvidia.com/cuda/)

Проверка установки:

- **Перегрузка**
- **Запись в ~/.bashrc:**
\$ export PATH=/usr/local/cuda-7.5/bin:\$PATH
\$ export LD_LIBRARY_PATH=/usr/local/cuda-7.5/lib64:\$LD_LIBRARY_PATH
- **Копирование файлов в папку пользователя:** \$ cuda-install-samples-7.5.sh <dir>
- **Компиляция, компоновка и запуск примера (*make*, затем *./<file_executable>*)**

Программный интерфейс CUDA

CUDA C - расширение языка *C* , и набор функций и структур *CUDA API* предоставляют простой инструмент для программирования на GPU.

Некоторые конструкции программного интерфейса CUDA.

Функция-ядро (*kernel*)

Код, выполняемый на устройстве (ядро), определяется в виде функции типа *void* со спецификатором `__global__`:

```
__global__ void gFunc(<params>){...}
```

Конфигурация нитей

При вызове ядра программист определяет количество нитей в блоке и количество блоков в `grid`. При этом допустима линейная, двумерная или трехмерная индексация нитей:

```
gFunc<<<dim3(bl_xdim, bl_ydim, bl_zdim),  
          dim3(th_xdim, th_ydim, th_zdim)>>>(<params>);
```

Программный интерфейс CUDA (самый простой пример)

```
#include <cuda.h>
#include <stdio.h>

__global__ void gTest(float* a){
    a[threadIdx.x+blockDim.x*blockIdx.x]=(float)(threadIdx.x+blockDim.x*blockIdx.x);
}

int main(){
    float *da, *ha;
    int num_of_blocks=10, threads_per_block=32;
    int N=num_of_blocks*threads_per_block;

    ha=(float*)calloc(N, sizeof(float));
    cudaMalloc((void**)&da, N*sizeof(float));

    gTest<<<dim3(num_of_blocks), dim3(threads_per_block)>>>(da);
    //cudaThreadSynchronize();
    CudaDeviceSynchronize();
    cudaMemcpy(ha,da,N*sizeof(float), cudaMemcpyDeviceToHost);

    for(int i=0;i<N;i++)
        printf("%g\n", ha[i]);

    free(ha);
    cudaFree(da);
    return 0;
}
```

CUDA C - расширение языка C , и набор функций и структур CUDA API предоставляют простой инструмент для программирования на GPU.

Комментарии: использование глобальной памяти

```
cudaError_t cudaMalloc (void ** devPtr, size_t size)
```

```
cudaError_t cudaFree (void * devPtr)
```

```
cudaError_t cudaMemcpy (void * dst, const void * src, size_t count, enum cudaMemcpyKind kind)
```

```
enum cudaError
{
    cudaSuccess = 0,
    cudaErrorMissingConfiguration,
    cudaErrorMemoryAllocation,
    cudaErrorInitializationError,
    cudaErrorLaunchFailure,
    .....
};

typedef enum cudaError cudaError_t;
```

```
enum cudaMemcpyKind
{
    cudaMemcpyHostToHost = 0,
    cudaMemcpyHostToDevice,
    cudaMemcpyDeviceToHost,
    cudaMemcpyDeviceToDevice
};
```

Глобальная память **выделяется только на хосте**, к глобальной памяти возможен **доступ только на устройстве**.

Документация CUDA: <http://docs.nvidia.com/cuda/index.html>

Комментарии: встроенные типы и переменные

- uint3 **threadIdx** – индекс нити в блоке
- dim3 **blockDim** – размер блока
- uint3 **blockIdx** – индекс блока в гриде
- dim3 **gridDim** - размер грида
- int **warpSize** - количество нитей в варпе (warp)

Комментарии: синхронизация всех нитей

Запуск ядра на устройстве (вызов функции с модификатором `__global__`) происходит в асинхронном режиме.

Для синхронизации нитей служат следующие вызовы:

```
//cudaThreadSynchronize(); //устаревшая функция (depricated)  
cudaDeviceSynchronize();
```

Обработка ошибок (код с обработчиком ошибок)

```
#include <cuda.h>
#include <stdio.h>

#define CUDA_CHECK_RETURN(value) {
    cudaError_t _m_cudaStat = value;
    if (_m_cudaStat != cudaSuccess) {
        fprintf(stderr, "Error %s at line %d in file %s\n",
            cudaGetErrorString(_m_cudaStat), __LINE__, __FILE__);
        exit(1);
    } } //макрос для обработки ошибок

__global__ void gTest(float* a){
    a[threadIdx.x+blockDim.x*blockIdx.x]=(float)(threadIdx.x+blockDim.x*blockIdx.x);
}

int main(){
    float *da, *ha;
    int num_of_blocks=10, threads_per_block=32;
    int N=num_of_blocks*threads_per_block;

    ha=(float*)calloc(N, sizeof(float));
    CUDA_CHECK_RETURN(cudaMalloc((void**)&da, N*sizeof(float)));
```

Обработка ошибок (код с обработчиком ошибок, продолжение)

```
gTest<<<dim3(num_of_blocks), dim3(threads_per_block)>>>(da);  
// cudaThreadSynchronize();  
CUDA_CHECK_RETURN(cudaDeviceSynchronize());  
CUDA_CHECK_RETURN(cudaGetLastError());  
CUDA_CHECK_RETURN(cudaMemcpy(ha,da,N*sizeof(float), cudaMemcpyDeviceToHost));  
  
for(int i=0;i<N;i++)  
    printf("%g\n", ha[i]);  
  
free(ha);  
cudaFree(da);  
return 0;  
}
```

Обработка ошибок (комментарии)

char * cudaGetErrorString (cudaError_t code); - возвращает строку в таблице кодов ошибок.

cudaError_t cudaGetLastError (); - возвращает код ошибки при асинхронном вызове.

Контроль производительности (код с контрольными точками)

```
#include <cuda.h>
#include <stdio.h>

#define CUDA_CHECK_RETURN(value) {
    cudaError_t _m_cudaStat = value;
    if (_m_cudaStat != cudaSuccess) {
        fprintf(stderr, "Error %s at line %d in file %s\n",
            cudaGetErrorString(_m_cudaStat), __LINE__, __FILE__);
        exit(1);
    }
}

__global__ void gTest(float* a){
    a[threadIdx.x+blockDim.x*blockIdx.x]=(float)(threadIdx.x+blockDim.x*blockIdx.x);
}

int main(){
    float *da, *ha;
    int num_of_blocks=10, threads_per_block=32;
    int N=num_of_blocks*threads_per_block;

    ha=(float*)calloc(N, sizeof(float));
    CUDA_CHECK_RETURN(cudaMalloc((void**)&da, N*sizeof(float)));
```

Контроль производительности (код с контрольными точками, продолжение)

```
cudaEvent_t start, stop; // встроенный тип данных – структура, для фиксации контрольных точек
float elapsedTime;
```

```
cudaEventCreate(&start); // инициализация
cudaEventCreate(&stop); // событий
```

```
cudaEventRecord(start, 0); // привязка события
```

```
gTest<<<dim3(num_of_blocks), dim3(threads_per_block)>>>(da);
// cudaThreadSynchronize();
//CUDA_CHECK_RETURN(cudaDeviceSynchronize());
//CUDA_CHECK_RETURN(cudaGetLastError());
```

```
cudaEventRecord(stop, 0); // привязка события
cudaEventSynchronize(stop); // синхронизация по событию
cudaEventElapsedTime(&elapsedTime, start, stop); // вычисление затраченного времени
```

```
fprintf(stderr, "gTest took %g\n", elapsedTime);
```

```
cudaEventDestroy(start); // освобождение
cudaEventDestroy(stop); // памяти
```

Контроль производительности (код с контрольными точками, окончание)

```
CUDA_CHECK_RETURN(cudaMemcpy(ha,da,N*sizeof(float), cudaMemcpyDeviceToHost));

for(int i=0;i<N;i++)
    printf("%g\n", ha[i]);

free(ha);
cudaFree(da);
return 0;
}
```

```
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> ./3>tmp
gTest took 0.030656
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> time(./3>tmp)
gTest took 0.030848

real    0m0.086s
user    0m0.006s
sys     0m0.058s
```


Компиляторы

Сравнение компиляторов nvcc и pgi:

```
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> pgcpp 3.cu -o 3p
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> nvcc 3.cu -o 3n
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> time(./3p>tmp)
gTest took 24.0531
```

```
real 0m0.025s
user 0m0.025s
sys 0m0.001s
```

```
ewgenij@linux-715l:~/EDUCATION/workshop/Practicals/Practical1> time(./3n>tmp)
gTest took 0.004992
```

```
real 0m0.113s
user 0m0.011s
sys 0m0.058s
```

Пример строки компиляции nvcc (компиляция и компоновка dll):

```
librelaxation.so: relaxation.cpp kerns_st.cu
nvcc -arch=compute_20 -Xcompiler -fPIC -shared relaxation.cpp kerns_st.cu -I../include
-I/home/ewgenij/common/inc -L/home/ewgenij/common/lib -lcutil_x86_64 -o ../lib/librelaxation.so
```

Упражнения

- Сравнить время выполнения последовательного и параллельного кода инициализации и сложения двух векторов большой размерности на выбранном CPU и GPU, установленными на серверах home, dew и кластере clu.nusc.ru.
- То же самое для задачи транспонирования матрицы.