

# INTRODUCTION TO GPU PROGRAMMING WITH THE CUDA PLATFORM

---

Pedro Bruel  
[phrb@ime.usp.br](mailto:phrb@ime.usp.br)  
4th INFIERI, 2017



Instituto de Matemática e Estatística  
Universidade de São Paulo

# ABOUT



Pedro Bruel  
[phrb@ime.usp.br](mailto:phrb@ime.usp.br)



Alfredo Goldman  
[gold@ime.usp.br](mailto:gold@ime.usp.br)



Marco Dimas Gubitoso  
[gubi@ime.usp.br](mailto:gubi@ime.usp.br)

# PART I

1. Introduction
2. Heterogeneous Computing
3. GPUs
4. CUDA Platform
5. CUDA C

## PART II

6. Recapitulation
7. Compiling CUDA Applications
8. Optimization Best Practices
9. Analysing CUDA Applications
10. Optimizing CUDA Applications
11. Conclusion

# SLIDES



This presentation and all source code are available at [GitHub](#):

- [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda)

# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

- Associated with a **host** processor

# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

- Associated with a **host** processor
- Separate **control** and **memory**

# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

- Associated with a **host** processor
- Separate **control** and **memory**
- Differ on **specialization** and **configurability**

# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

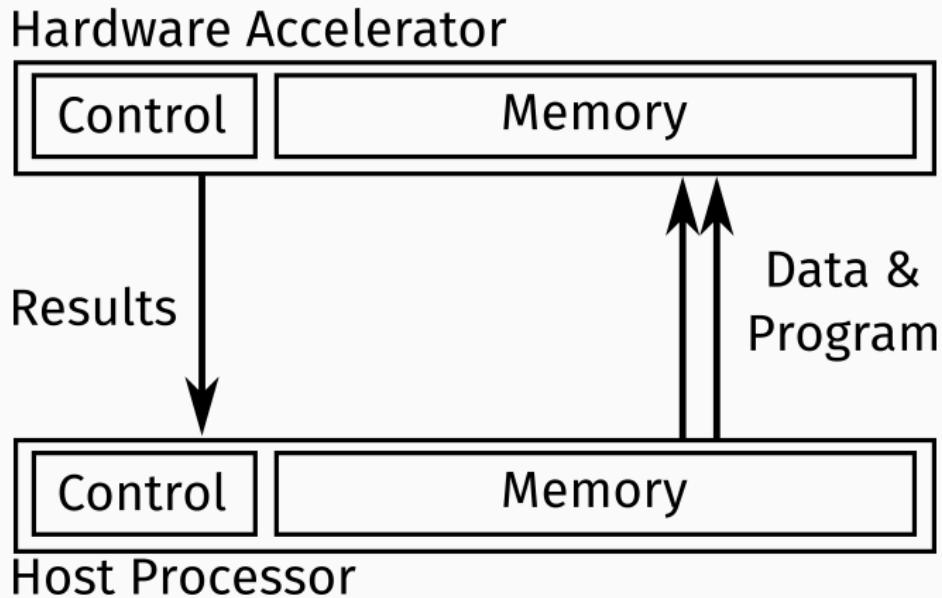
- Associated with a **host** processor
- Separate **control** and **memory**
- Differ on **specialization** and **configurability**
- **GPUs**, DSPs, FPGAs, ASICs

# HARDWARE ACCELERATION

Use cases:

- Machine Learning
- Digital Signal Processing
- Bioinformatics
- Cryptography
- Meteorology
- Simulations
- ...

# HARDWARE ACCELERATION



# HARDWARE ACCELERATION

Percentage of systems with accelerators in the *Top500*:

## ACCELERATORS/CO-PROCESSORS

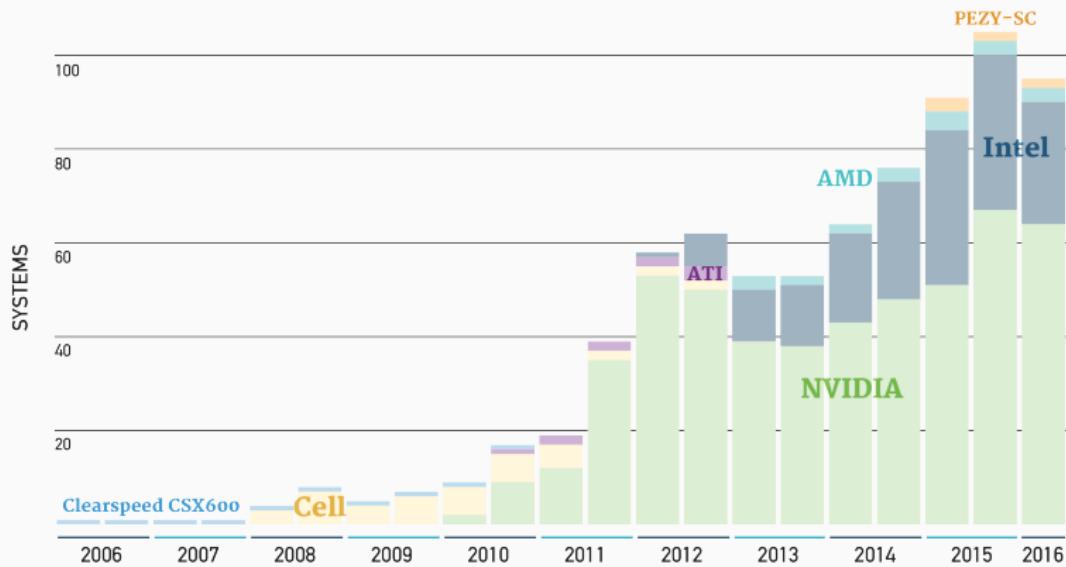


Image: [top500.org/lists/2016/06/download/TOP500\\_201606\\_Poster.pdf](http://top500.org/lists/2016/06/download/TOP500_201606_Poster.pdf) [Accessed in 29/07/16]

# HETEROGENEOUS COMPUTING



# HETEROGENEOUS COMPUTING



Given **heterogeneous** computational resources, **datasets** and a set of **computations**, how to distribute data and computations to **optimize** resource usage?

Image: [olcf.ornl.gov/titan](http://olcf.ornl.gov/titan) [Accessed in 29/07/16]

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

- Low latency: CPUs

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

- Low latency: CPUs
- High Throughput: GPUs

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

- Low latency: CPUs
- High Throughput: GPUs
- Reconfigurable: FPGAs

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

- Low latency: CPUs
- High Throughput: GPUs
- Reconfigurable: FPGAs
- Specialized: ASICs, DSPs

# HETEROGENEOUS COMPUTING

Heterogeneous computational resources:

- Low latency: CPUs
- High Throughput: GPUs
- Reconfigurable: FPGAs
- Specialized: ASICs, DSPs
- Memory?

# HETEROGENEOUS COMPUTING

## Datasets:

- *Big Data*
- *Data Streams*
- ...

# HETEROGENEOUS COMPUTING

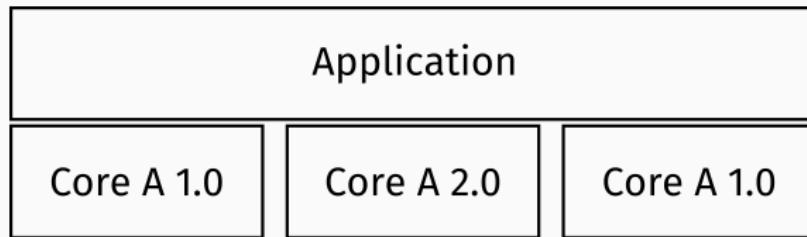
Datasets:

- *Big Data*
- *Data Streams*
- ...

Programming models (**computations**):

- *MapReduce*
- *Task Parallelism*
- ...

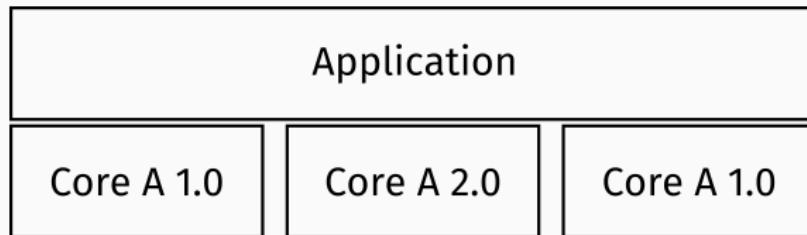
# HETEROGENEOUS COMPUTING: SCALABILITY



**Scalability** means not losing performance running in:

- Multiple identical cores

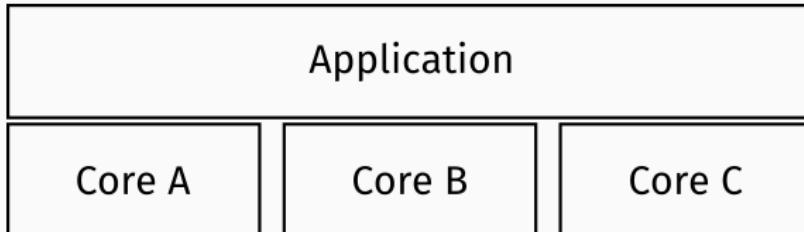
# HETEROGENEOUS COMPUTING: SCALABILITY



Scalability means not losing performance running in:

- Multiple identical cores
- New versions of the same core

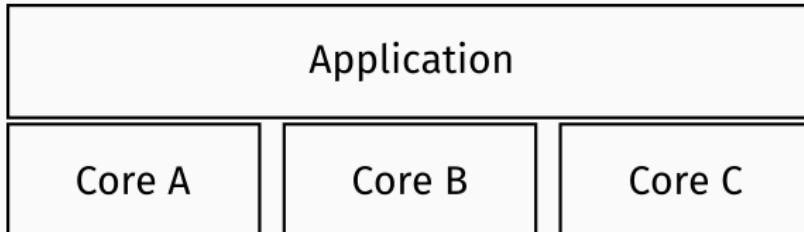
# HETEROGENEOUS COMPUTING: PORTABILITY



Portability means not losing performance running in different:

- Cores or accelerators

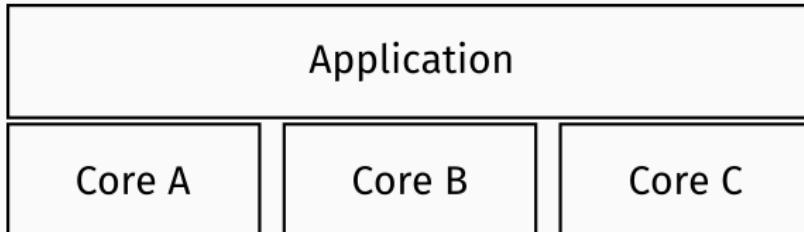
# HETEROGENEOUS COMPUTING: PORTABILITY



**Portability** means not losing performance running in different:

- Cores or **accelerators**
- Memory models

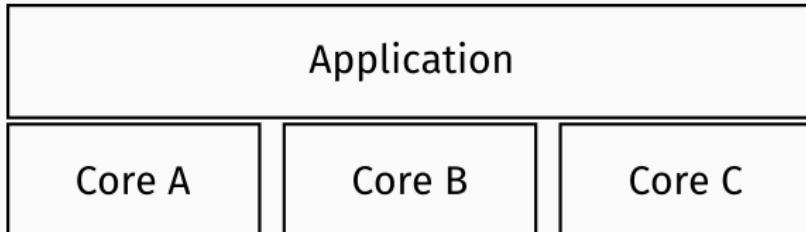
# HETEROGENEOUS COMPUTING: PORTABILITY



**Portability** means not losing performance running in different:

- Cores or **accelerators**
- Memory models
- Parallelism models

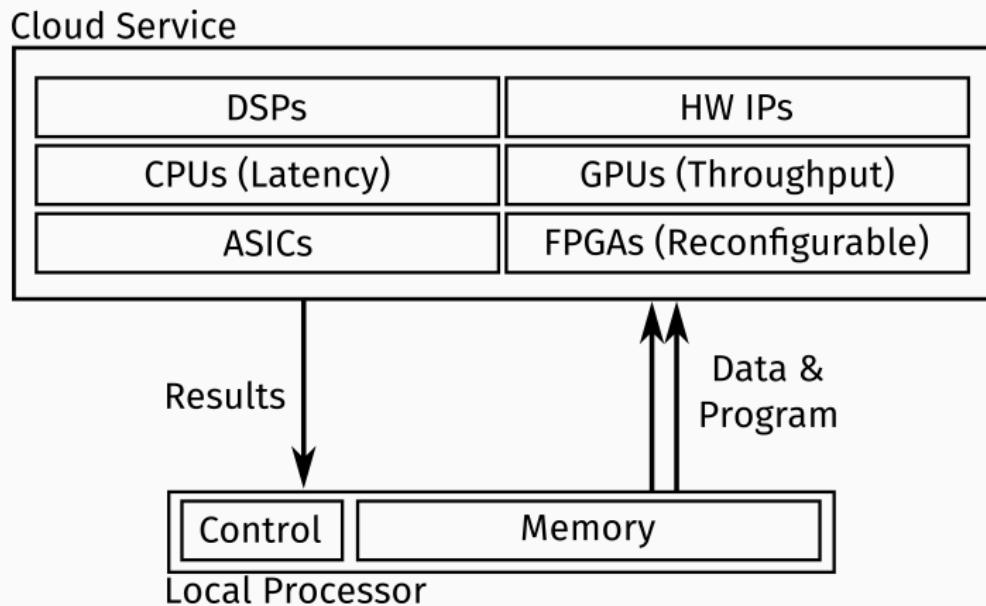
# HETEROGENEOUS COMPUTING: PORTABILITY



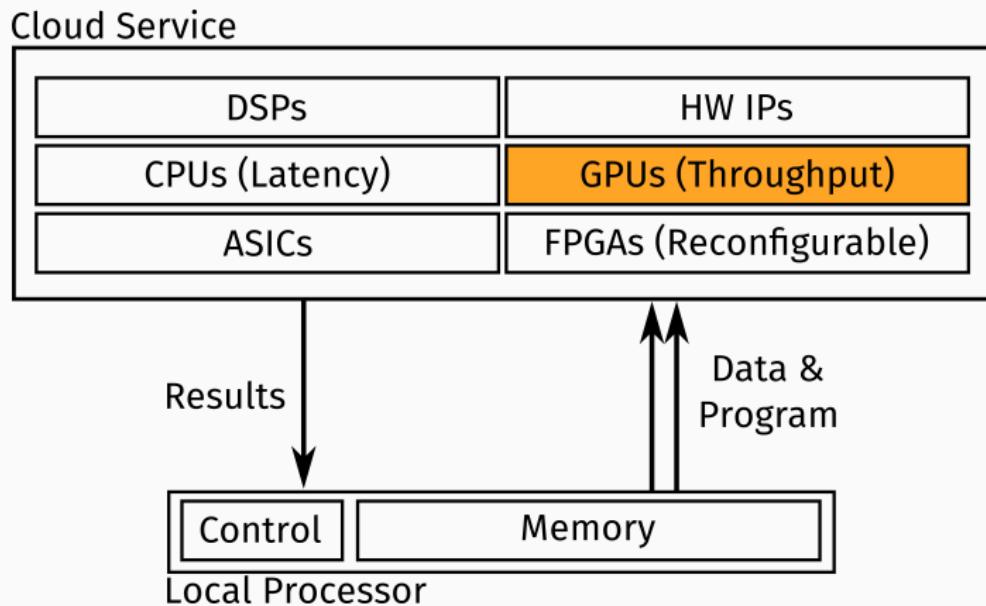
Portability means not losing performance running in different:

- Cores or accelerators
- Memory models
- Parallelism models
- Instruction sets

# HETEROGENEOUS COMPUTING



# HETEROGENEOUS COMPUTING



# GRAPHICS PROCESSING UNITS



# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

- Small caches (*kilobytes*)

# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

- Small caches (*kilobytes*)
- No **branch prediction**

# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

- Small caches (*kilobytes*)
- No **branch prediction**
- Thousands of **higher latency** ALUs

# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

- Small caches (*kilobytes*)
- No **branch prediction**
- Thousands of **higher latency** ALUs
- Execution **pipelines**

# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

- Small caches (*kilobytes*)
- No **branch prediction**
- Thousands of **higher latency** ALUs
- Execution **pipelines**
- 114 688 concurrent **threads** in Pascal architecture

# GRAPHICS PROCESSING UNITS

Now, largely used in general-purpose computing, also called *General Purpose GPUs* (GPGPUs):

# GRAPHICS PROCESSING UNITS

Now, largely used in general-purpose computing, also called *General Purpose GPUs* (GPGPUs):

- *OpenGL*: 1992
- *DirectX*: 1995
- **CUDA**: 2007
- *OpenCL*: 2009
- *Vulkan*: 2016

# HARDWARE MODEL

Flynn's taxonomy:

- *Single Instruction Multiple Data (SIMD)*

# HARDWARE MODEL

Flynn's taxonomy:

- *Single Instruction Multiple Data (SIMD)*
- *Single Instruction Multiple Thread (SIMT)?*

# HARDWARE MODEL

Flynn's taxonomy:

- *Single Instruction Multiple Data (SIMD)*
- *Single Instruction Multiple Thread (SIMT)?*

Scheduling and execution:

- *Streaming Multiprocessor (SM)*

# HARDWARE MODEL

Flynn's taxonomy:

- *Single Instruction Multiple Data (SIMD)*
- *Single Instruction Multiple Thread (SIMT)?*

Scheduling and execution:

- *Streaming Multiprocessor (SM)*
- *Warps*

# HARDWARE MODEL

Flynn's taxonomy:

- *Single Instruction Multiple Data (SIMD)*
- *Single Instruction Multiple Thread (SIMT)?*

Scheduling and execution:

- *Streaming Multiprocessor (SM)*
- *Warps*
- *Grids, blocks and threads*

# HARDWARE MODEL

Higher-Level scheduling:

- *Pipelines*
- *Texture Processing Cluster (TPC)*
- *Graphics Processing Cluster (GPC)*

# HARDWARE MODEL

Higher-Level scheduling:

- *Pipelines*
- *Texture Processing Cluster (TPC)*
- *Graphics Processing Cluster (GPC)*

Memory:

- Shared: L1, L2 caches; *megabytes* (NVIDIA Pascal)
- Global: volatile, GDDR5; *gigabytes*

# HARDWARE MODEL

Higher-Level scheduling:

- *Pipelines*
- *Texture Processing Cluster (TPC)*
- *Graphics Processing Cluster (GPC)*

Memory:

- Shared: L1, L2 caches; *megabytes* (NVIDIA Pascal)
- Global: volatile, GDDR5; *gigabytes*
- SSD: **non-volatile**; **terabytes** (AMD SSG Fiji, 2017)

AMD SSG: [anandtech.com/show/10518/amd-announces-radeon-pro-ssg-fiji-with-m2-ssds-onboard](http://anandtech.com/show/10518/amd-announces-radeon-pro-ssg-fiji-with-m2-ssds-onboard) [Accessed in 30/07/16]

## *COMPUTE CAPABILITY AND SIMT*

The **Compute Capability** specifies *Single Instruction Multiple Thread* (SIMT) architectures:

## *COMPUTE CAPABILITY AND SIMT*

The **Compute Capability** specifies *Single Instruction Multiple Thread* (SIMT) architectures:

- *Warp*: Group of *threads* that runs in **parallel**
- *Streaming Multiprocessor (SM)*: Creates, schedules and runs *warps*

# *COMPUTE CAPABILITY AND SIMT*

The **Compute Capability** specifies *Single Instruction Multiple Thread* (SIMT) architectures:

- *Warp*: Group of *threads* that runs in **parallel**
- *Streaming Multiprocessor* (SM): Creates, schedules and runs *warps*
- Multiple *warps* run **concurrently** in a SM

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#hardware-implementation](http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#hardware-implementation) [Accessed in 29/07/16]

# *COMPUTE CAPABILITY AND SIMT*

*Warps:*

## *COMPUTE CAPABILITY AND SIMT*

*Warps:*

- Unit of execution and scheduling
- 32 *threads*

## *COMPUTE CAPABILITY AND SIMT*

*Warps:*

- Unit of execution and scheduling
- 32 *threads*
- 1 **common instruction** at a time

## *COMPUTE CAPABILITY AND SIMT*

*Warps:*

- Unit of execution and scheduling
- 32 *threads*
- 1 **common instruction** at a time
- *Active threads*: in the current execution branch
- *Inactive threads*: **not** in the current execution branch

## *COMPUTE CAPABILITY AND SIMT*

*Warps:*

- Unit of execution and scheduling
- 32 *threads*
- 1 **common instruction** at a time
- *Active threads*: in the current execution branch
- *Inactive threads*: **not** in the current execution branch
- Inactive threads run **sequentially**

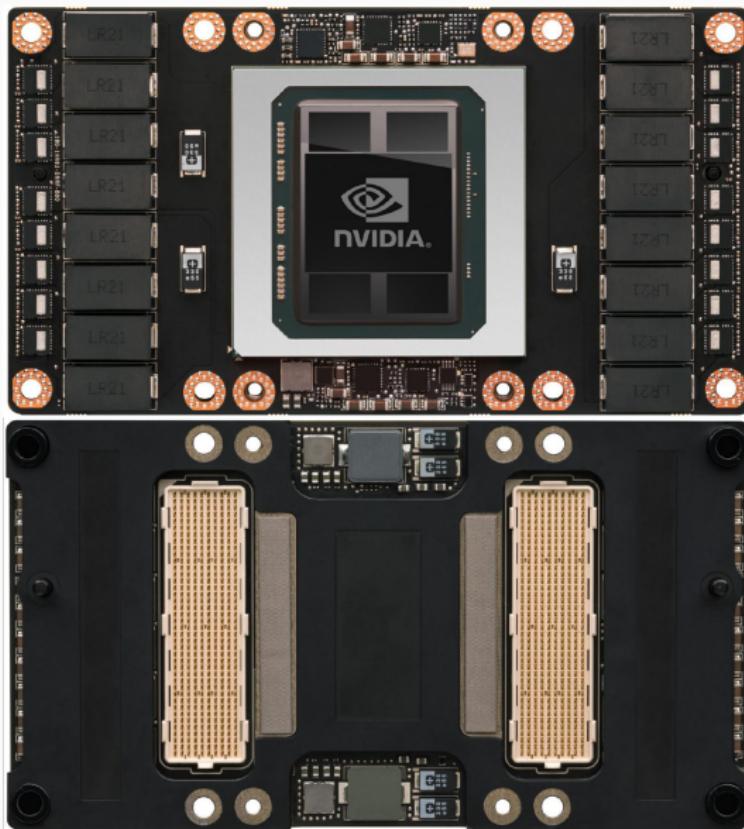
# COMPUTE CAPABILITY AND SIMT

Warps:

- Unit of execution and scheduling
- 32 *threads*
- 1 **common instruction** at a time
- *Active threads*: in the current execution branch
- *Inactive threads*: **not** in the current execution branch
- Inactive threads run **sequentially**
- Efficiency: no execution divergence inside *Warps*

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#hardware-implementation](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#hardware-implementation) [Accessed in 29/07/16]

# PASCAL GP100 ARCHITECTURE



# PASCAL GP100 ARCHITECTURE: COMPUTE CAPABILITY 6.0

GPU	Kepler GK110	Maxwell GM200	Pascal GP100
Compute Capability	3.5	5.2	6.0
Threads / Warp	32	32	32
Max Warps / Multiprocessor	64	64	64
Max Threads / Multiprocessor	2048	2048	2048
Max Thread Blocks / Multiprocessor	16	32	32
Max 32-bit Registers / SM	65536	65536	65536
Max Registers / Block	65536	32768	65536
Max Registers / Thread	255	255	255
Max Thread Block Size	1024	1024	1024
Shared Memory Size / SM	16 KB/32 KB/48 KB	96 KB	64 KB

Image: [images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf](http://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf) [Accessed in 29/07/16]

# PASCAL GP100 ARCHITECTURE: SM



Image: [images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf](http://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf) [Accessed in 29/07/16]

# PASCAL GP100 ARCHITECTURE

Tesla Products	Tesla K40	Tesla M40	Tesla P100
<b>GPU</b>	GK110 (Kepler)	GM200 (Maxwell)	GP100 (Pascal)
<b>SMs</b>	15	24	56
<b>TPCs</b>	15	24	28
<b>FP32 CUDA Cores / SM</b>	192	128	64
<b>FP32 CUDA Cores / GPU</b>	2880	3072	3584
<b>FP64 CUDA Cores / SM</b>	64	4	32
<b>FP64 CUDA Cores / GPU</b>	960	96	1792
<b>Base Clock</b>	745 MHz	948 MHz	1328 MHz
<b>GPU Boost Clock</b>	810/875 MHz	1114 MHz	1480 MHz
<b>Peak FP32 GFLOPs<sup>1</sup></b>	5040	6840	10600
<b>Peak FP64 GFLOPs<sup>1</sup></b>	1680	210	5300
<b>Texture Units</b>	240	192	224
<b>Memory Interface</b>	384-bit GDDR5	384-bit GDDR5	4096-bit HBM2
<b>Memory Size</b>	Up to 12 GB	Up to 24 GB	16 GB
<b>L2 Cache Size</b>	1536 KB	3072 KB	4096 KB
<b>Register File Size / SM</b>	256 KB	256 KB	256 KB
<b>Register File Size / GPU</b>	3840 KB	6144 KB	14336 KB
<b>TDP</b>	235 Watts	250 Watts	300 Watts
<b>Transistors</b>	7.1 billion	8 billion	15.3 billion
<b>GPU Die Size</b>	551 mm <sup>2</sup>	601 mm <sup>2</sup>	610 mm <sup>2</sup>
<b>Manufacturing Process</b>	28-nm	28-nm	16-nm FinFET

<sup>1</sup> The GFLOPS in this chart are based on GPU Boost Clocks.

## PASCAL GP100 ARCHITECTURE



Image: [images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf](http://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf) [Accessed in 29/07/16]

# CUDA PLATFORM

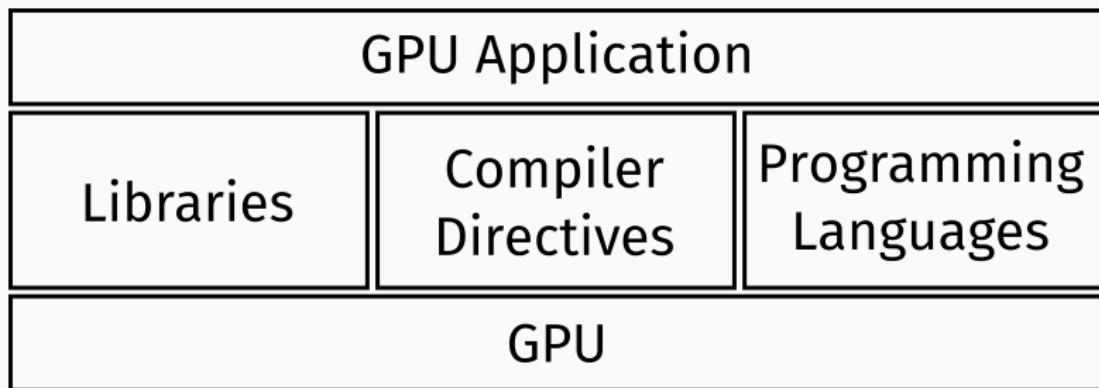


## CUDA PLATFORM

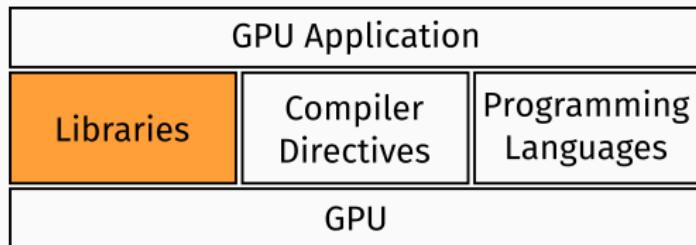


- A platform for **parallel computing**
- *Application Programming Interface* (API)
- *CUDA Toolkit*

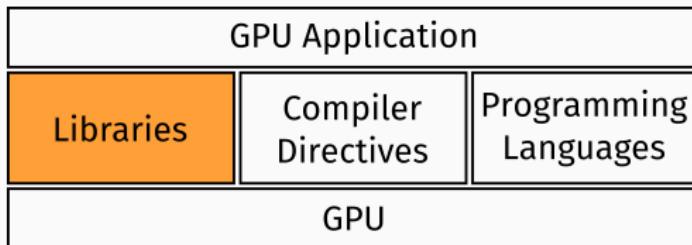
# SOFTWARE ACCELERATION



# LIBRARIES

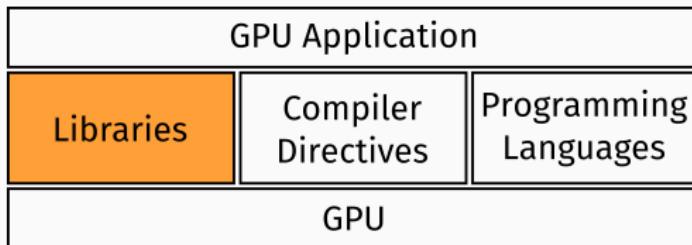


# LIBRARIES



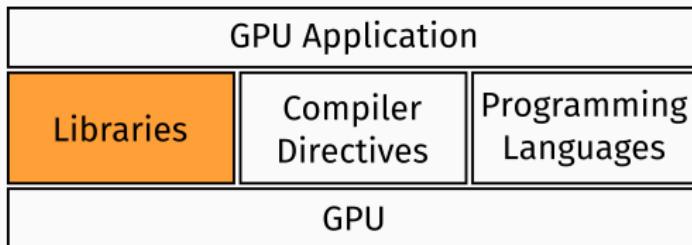
- Easy to use

# LIBRARIES



- Easy to use
- Drop-in acceleration

# LIBRARIES



- Easy to use
- Drop-in acceleration
- Optimized by specialists



# LIBRARIES

Vector sum with the Thrust library:

```
thrust::device_vector<float> device_input1(input_length);
thrust::device_vector<float> device_input2(input_length);
thrust::device_vector<float> device_output(input_length);
```

## LIBRARIES

Vector sum with the Thrust library:

```
thrust::device_vector<float> device_input1(input_lenght);
thrust::device_vector<float> device_input2(input_lenght);
thrust::device_vector<float> device_output(input_lenght);

thrust::copy(host_input1, host_input1 + input_lenght,
            device_input1.begin());

thrust::copy(host_input2, host_input2 + input_lenght,
            device_input2.begin());
```

## LIBRARIES

Vector sum with the Thrust library:

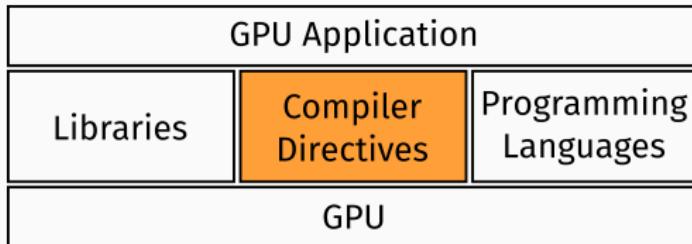
```
thrust::device_vector<float> device_input1(input_lenght);
thrust::device_vector<float> device_input2(input_lenght);
thrust::device_vector<float> device_output(input_lenght);

thrust::copy(host_input1, host_input1 + input_lenght,
            device_input1.begin());

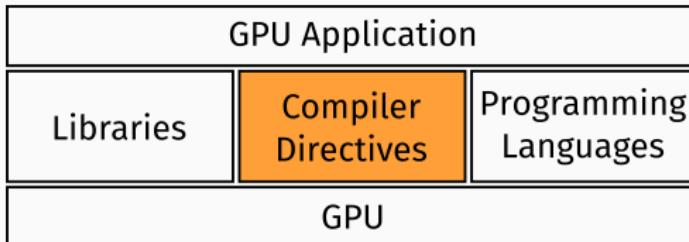
thrust::copy(host_input2, host_input2 + input_lenght,
            device_input2.begin());

thrust::transform(device_input1.begin(),
                 device_input1.end(), device_input2.begin(),
                 device_output.begin(), thrust::plus<float>());
```

# COMPILER DIRECTIVES

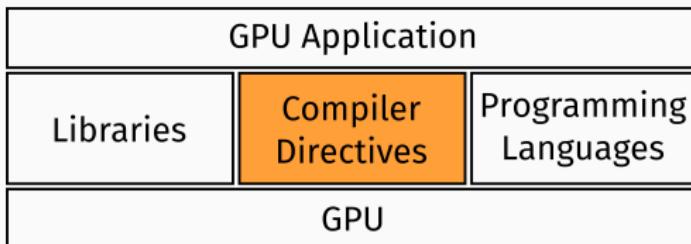


# COMPILER DIRECTIVES



- Easy to use

# COMPILER DIRECTIVES



- Easy to use
- Portable, but performance depends on compiler

# COMPILER DIRECTIVES

Vector sum with OpenACC:

```
#pragma acc data
    copyin(input1[0:input_lenght],input2[0:input_lenght]),
    copyout(output[0:input_lenght])
```

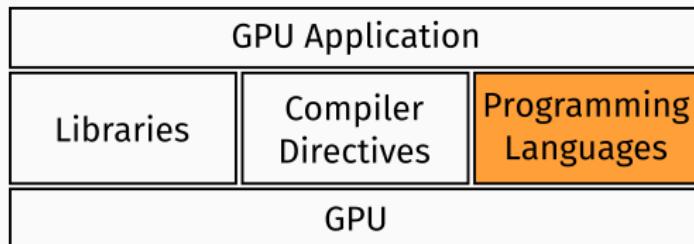
# COMPILER DIRECTIVES

Vector sum with OpenACC:

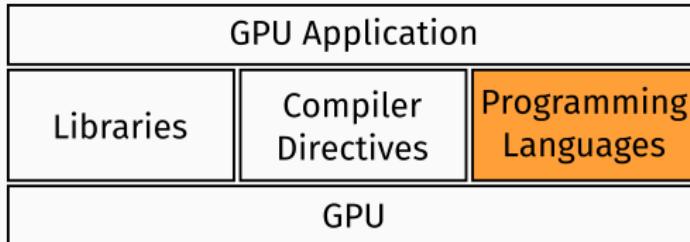
```
#pragma acc data
    copyin(input1[0:input_lenght],input2[0:input_lenght]),
    copyout(output[0:input_lenght])

{
    #pragma acc kernels loop independent
    for(i = 0; i < input_lenght; i++) {
        output[i] = input1[i] + input2[i];
    }
}
```

# PROGRAMMING LANGUAGES

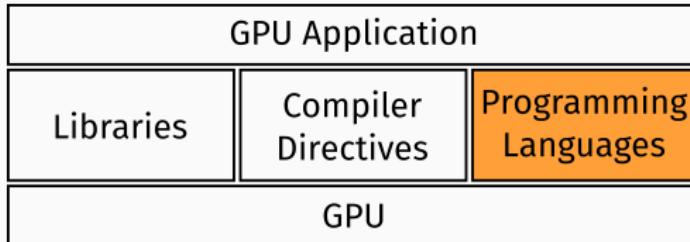


# PROGRAMMING LANGUAGES



- Best performance, but hard to master

# PROGRAMMING LANGUAGES



- Best performance, but hard to master
- Flexible

# PROGRAMMING LANGUAGES

CUDA implementations in:

- C/C++
- Fortran
- Python
- F#

# CUDA C





Programming Model:



Programming Model:

- Kernels



Programming Model:

- Kernels
- Threads hierarchy



Programming Model:

- Kernels
- Threads hierarchy
- Memory hierarchy



Programming Model:

- Kernels
- Threads hierarchy
- Memory hierarchy
- Heterogeneous Programming

# CUDA C: KERNEL

## Kernels:

- **Kernels** are C functions executed by **threads**
- $N$  *threads* run  $N$  *kernels*

# CUDA C: KERNEL

## Kernels:

- **Kernels** are C functions executed by **threads**
- $N$  *threads* run  $N$  *kernels*
- Access thread **IDs** using the **threadIdx** variable

# CUDA C: KERNEL

## Kernels:

- Kernels are C functions executed by threads
- $N$  threads run  $N$  kernels
- Access thread IDs using the `threadIdx` variable
- Defined within the `__global__` scope
- Its type must be `void`

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

## CUDA C: KERNEL

### Kernels:

- Configured and launched using the syntax:
  - `kernel_name<<<...>>>( ... );`

## CUDA C: KERNEL

### Kernels:

- Configured and launched using the syntax:
  - `kernel_name<<<...>>>( ... );`
- Ideally executed in `parallel`

# CUDA C: KERNEL

## Kernels:

- Configured and launched using the syntax:
  - `kernel_name<<<...>>>( ... );`
- Ideally executed in `parallel`
- In practice, parallelism depends on:
  - Number of threads in a warp

# CUDA C: KERNEL

## Kernels:

- Configured and launched using the syntax:
  - `kernel_name<<<...>>>( ... );`
- Ideally executed in `parallel`
- In practice, parallelism depends on:
  - Number of threads in a warp
  - Execution branch coherence

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

# CUDA C: KERNEL

Writing and launching a kernel:

```
#include <cuda_runtime.h>

__global__ void VecAdd(float* A, float* B, float* C) {
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}
```

# CUDA C: KERNEL

Writing and launching a kernel:

```
#include <cuda_runtime.h>

__global__ void VecAdd(float* A, float* B, float* C) {
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}

int main() {
    ...
    VecAdd<<<1, N>>>(A, B, C);
    ...
}
```

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

## CUDA C: THREAD HIERARCHY

### Thread Block:

- Grouping of threads

## CUDA C: THREAD HIERARCHY

### Thread Block:

- Grouping of threads
- Tridimensional:  $Db = (Dx, Dy, Dz)$

## CUDA C: THREAD HIERARCHY

### Thread Block:

- Grouping of threads
- Tridimensional:  $D_b = (D_x, D_y, D_z)$
- Maximum size of 1024 threads

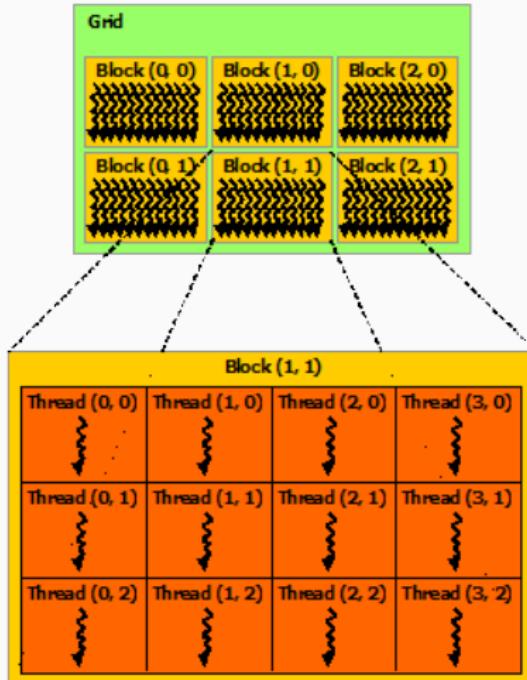
# CUDA C: THREAD HIERARCHY

## Thread Block:

- Grouping of threads
- Tridimensional:  $D_b = (D_x, D_y, D_z)$
- Maximum size of 1024 threads
- A Grid is the tridimensional grouping of *blocks*

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

# CUDA C: THREAD HIERARCHY



Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

# REMEMBERING: PASCAL GP100 SM



Image: [images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf](http://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf) [Accessed in 29/07/16]

## CUDA C: THREAD HIERARCHY

```
__global__ void MatAdd(float A[N][N], float B[N][N], float
C[N][N]) {
    int i = threadIdx.x;
    int j = threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}
```

## CUDA C: THREAD HIERARCHY

```
__global__ void MatAdd(float A[N][N], float B[N][N], float
C[N][N]) {
    int i = threadIdx.x;
    int j = threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main() {
    ...
    int numBlocks = 1;
    dim3 threadsPerBlock(N, N);
    MatAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
    ...
}
```

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

## CUDA C: THREAD HIERARCHY

```
__global__ void MatAdd(float A[N][N], float B[N][N], float
C[N][N]) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < N && j < N) {
        C[i][j] = A[i][j] + B[i][j];
    }
}
```

## CUDA C: THREAD HIERARCHY

```
__global__ void MatAdd(float A[N][N], float B[N][N], float
C[N][N]) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < N && j < N) {
        C[i][j] = A[i][j] + B[i][j];
    }
}

int main() {
    ...
    dim3 threadsPerBlock(16, 16);
    dim3 numBlocks(N / threadsPerBlock.x, N /
                    threadsPerBlock.y);
    MatAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
    ...
}
```

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

## CUDA C: THREAD HIERARCHY

Configuration and Launching syntax `<<< Dg, Db, Ns, S >>>`:

- **dim3** Dg grid dimension and size:

`Dg.x * Dg.y * Dg.z = numBlocks`

## CUDA C: THREAD HIERARCHY

Configuration and Launching syntax `<<< Dg, Db, Ns, S >>>`:

- **dim3** Dg grid dimension and size:

`Dg.x * Dg.y * Dg.z = numBlocks`

- **dim3** Db block dimension and size:

`Db.x * Db.y * Db.z = threadsPerBlock`

## CUDA C: THREAD HIERARCHY

Configuration and Launching syntax `<<< Dg, Db, Ns, S >>>`:

- `dim3` `Dg` grid dimension and size:

$$Dg.x * Dg.y * Dg.z = \text{numBlocks}$$

- `dim3` `Db` block dimension and size:

$$Db.x * Db.y * Db.z = \text{threadsPerBlock}$$

- `size_t` `Ns` = 0: extra bytes in shared memory

- `cudaStream_t` `S` = 0: CUDA stream

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

## CUDA C: MEMORY HIERARCHY

Threads access **multiple memory addresses** during a kernel's execution:

- Local

## CUDA C: MEMORY HIERARCHY

Threads access **multiple memory addresses** during a kernel's execution:

- Local
- Shared with its *block*

## CUDA C: MEMORY HIERARCHY

Threads access **multiple memory addresses** during a kernel's execution:

- Local
- Shared with its *block*
- Global

## CUDA C: MEMORY HIERARCHY

Threads access **multiple memory addresses** during a kernel's execution:

- Local
- Shared with its *block*
- Global: persists between kernels of the **same application**

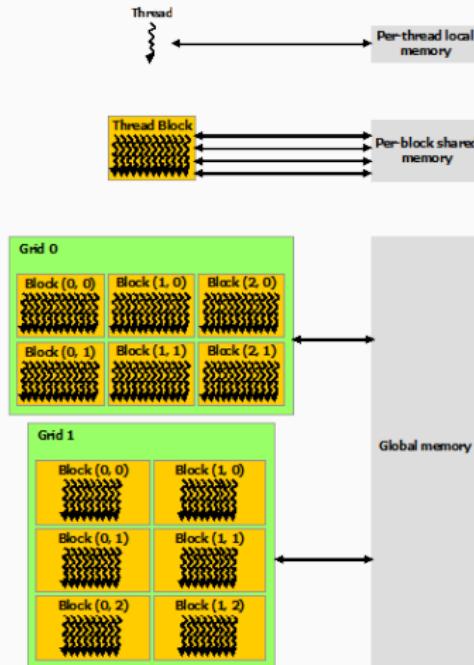
## CUDA C: MEMORY HIERARCHY

Threads access **multiple memory addresses** during a kernel's execution:

- Local
- Shared with its *block*
- Global: persists between kernels of the **same application**
- *Read-only*

Source: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

# CUDA C: MEMORY HIERARCHY



SOURCE: [docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model) [Accessed in 29/07/16]

# CUDA C: HETEROGENEOUS PROGRAMMING

Host's tasks:

- Allocate device memory and resources
- Move data (bottleneck!)

# CUDA C: HETEROGENEOUS PROGRAMMING

Host's tasks:

- Allocate device memory and resources
- Move data (bottleneck!)
- Launch kernel

# CUDA C: HETEROGENEOUS PROGRAMMING

Host's tasks:

- Allocate device memory and resources
- Move data (bottleneck!)
- Launch kernel
- Move results (bottleneck!)
- Release device memory and resources

# CUDA C: HETEROGENEOUS PROGRAMMING

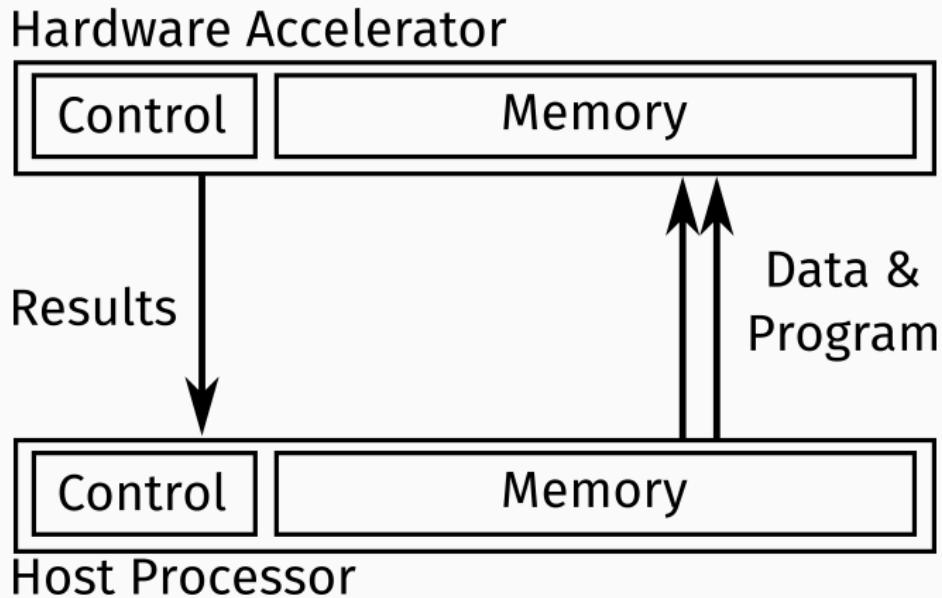
Host's tasks:

- Allocate device memory and resources
- Move data (bottleneck!)
- Launch kernel
- Move results (bottleneck!)
- Release device memory and resources

Device's tasks:

- Run kernel

## CUDA C: HETEROGENEOUS PROGRAMMING



## EXAMPLE: VECTOR SUM

```
--global__ void vectorAdd(const float *A, const float *B,
    float *C, int numElements) {
    int i = blockDim.x * blockIdx.x + threadIdx.x;

    if (i < numElements) {
        C[i] = A[i] + B[i];
    }
}
```

Source: [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda) [Accessed in 29/07/16]

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>
```

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>

float *h_A = (float *) malloc(size);
if (h_A == NULL) { ... };

float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);
err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
if (err != cudaSuccess) { ... };
```

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>

float *h_A = (float *) malloc(size);
if (h_A == NULL) { ... };

float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);
err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
if (err != cudaSuccess) { ... };

int threadsPerBlock = 256;
int blocksPerGrid = (numElements + threadsPerBlock - 1) /
    threadsPerBlock;
```

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>

float *h_A = (float *) malloc(size);
if (h_A == NULL) { ... };

float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);
err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
if (err != cudaSuccess) { ... };

int threadsPerBlock = 256;
int blocksPerGrid = (numElements + threadsPerBlock - 1) /
    threadsPerBlock;

vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C,
    numElements);

err = cudaGetLastError()
err = cudaDeviceSynchronize()
if (err != cudaSuccess) { ... };
```

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>

float *h_A = (float *) malloc(size);
if (h_A == NULL) { ... };

float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);
err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
if (err != cudaSuccess) { ... };

int threadsPerBlock = 256;
int blocksPerGrid = (numElements + threadsPerBlock - 1) /
    threadsPerBlock;

vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C,
    numElements);

err = cudaGetLastError()
err = cudaDeviceSynchronize()
if (err != cudaSuccess) { ... };

err = cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
err = cudaFree(d_A);
if (err != cudaSuccess) { ... };
```

## EXAMPLE: PRACTICAL TIP

Always check for errors!

```
float *d_A = NULL;  
err = cudaMalloc((void **) &d_A, size);  
  
err = cudaGetLastError()  
  
err = cudaDeviceSynchronize()
```

## EXAMPLE: PRACTICAL TIP

Always check for errors!

```
float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);

err = cudaGetLastError()

err = cudaDeviceSynchronize()

if (err != cudaSuccess) {
    fprintf(stderr, "Failed to allocate device vector A
(error code %s)!\n", cudaGetStringError(err));
    exit(EXIT_FAILURE);
}
```

Source: [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda) [Accessed in 29/07/16]

# HANDS ON!

We'll run some CUDA C samples available at  
[github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda):

- `src/cuda-samples/0_Simple/vectorAdd`

# HANDS ON!

We'll run some CUDA C samples available at  
[github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda):

- `src/cuda-samples/0_Simple/vectorAdd`
- `src/cuda-samples/5_Simulations/fluidsGL`
- `src/cuda-samples/5_Simulations/oceanFFT`
- `src/cuda-samples/5_Simulations/smokeParticles`

# HANDS ON!

We'll run some CUDA C samples available at  
[github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda):

- `src/cuda-samples/0_Simple/vectorAdd`
- `src/cuda-samples/5_Simulations/fluidsGL`
- `src/cuda-samples/5_Simulations/oceanFFT`
- `src/cuda-samples/5_Simulations/smokeParticles`
- `src/mandelbrot_numba`

# RESOURCES

This presentation and all source code are available at [GitHub](#):

- [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda)

More resources:

- CUDA C: [docs.nvidia.com/cuda/cuda-c-programming-guide](https://docs.nvidia.com/cuda/cuda-c-programming-guide)
- CUDA Toolkit: [developer.nvidia.com/cuda-toolkit](https://developer.nvidia.com/cuda-toolkit)
- Best Practices Guide:
  - [docs.nvidia.com/cuda/cuda-c-best-practices-guide](https://docs.nvidia.com/cuda/cuda-c-best-practices-guide)
- GPU Teaching Kit: [syllabus.gputeachingkit.com](https://syllabus.gputeachingkit.com)
- iPython: [ipython.org/notebook.html](https://ipython.org/notebook.html)
- Anaconda: [continuum.io/downloads](https://continuum.io/downloads)

# INTRODUCTION TO GPU PROGRAMMING WITH THE CUDA PLATFORM

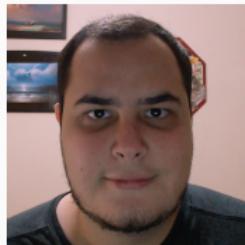
---

Pedro Bruel  
[phrb@ime.usp.br](mailto:phrb@ime.usp.br)  
4th INFIERI, 2017



Instituto de Matemática e Estatística  
Universidade de São Paulo

# ABOUT



Pedro Bruel  
[phrb@ime.usp.br](mailto:phrb@ime.usp.br)



Alfredo Goldman  
[gold@ime.usp.br](mailto:gold@ime.usp.br)



Marco Dimas Gubitoso  
[gubi@ime.usp.br](mailto:gubi@ime.usp.br)

# PART I

1. Introduction
2. Heterogeneous Computing
3. GPUs
4. CUDA Platform
5. CUDA C

## PART II

6. Recapitulation
7. Compiling CUDA Applications
8. Optimization Best Practices
9. Analysing CUDA Applications
10. Optimizing CUDA Applications
11. Conclusion

# SLIDES



This presentation and all source code are available at [GitHub](#):

- [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda)

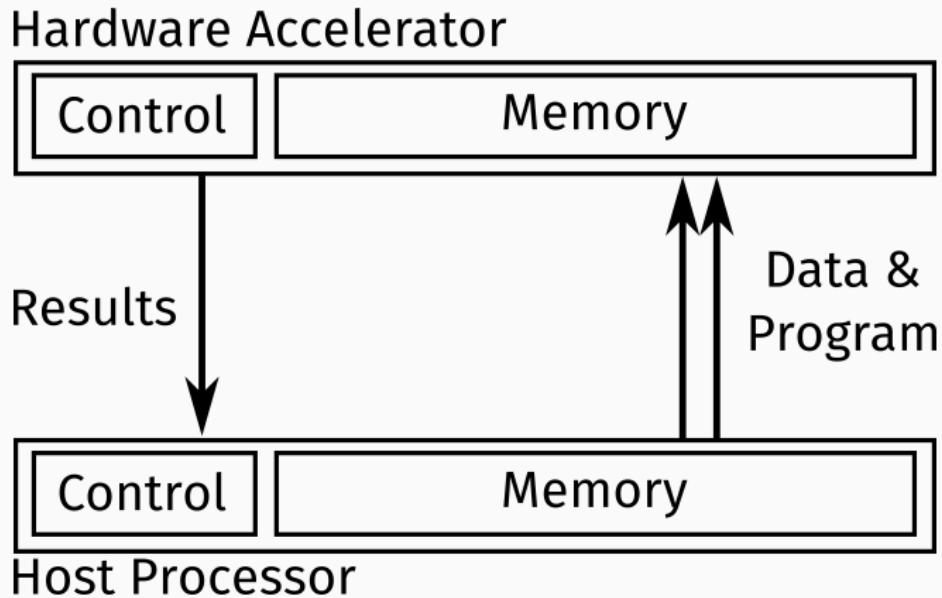
# HARDWARE ACCELERATION



The use of **devices** to accelerate computations on large datasets:

- Associated with a **host** processor
- Separate **control** and **memory**
- Differ on **specialization** and **configurability**
- **GPUs**, DSPs, FPGAs, ASICs

# HARDWARE ACCELERATION



# HARDWARE ACCELERATION

Percentage of systems with accelerators in the *Top500*:

## ACCELERATORS/CO-PROCESSORS

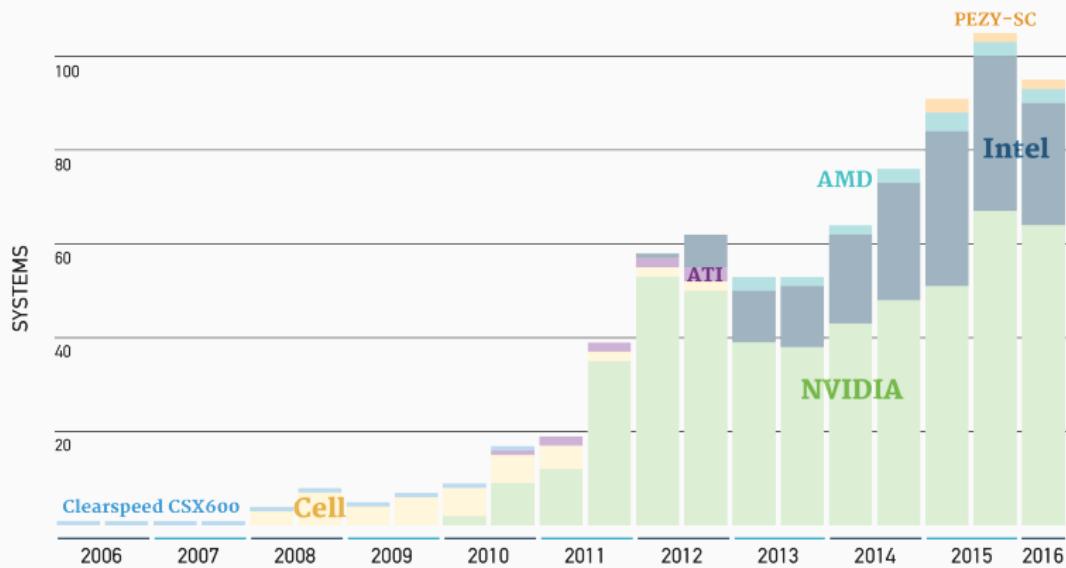


Image: [top500.org/lists/2016/06/download/TOP500\\_201606\\_Poster.pdf](http://top500.org/lists/2016/06/download/TOP500_201606_Poster.pdf) [Accessed in 29/07/16]

# HETEROGENEOUS COMPUTING



Given **heterogeneous** computational resources, **datasets** and a set of **computations**, how to distribute data and computations to **optimize** resource usage?

Image: [olcf.ornl.gov/titan](http://olcf.ornl.gov/titan) [Accessed in 29/07/16]

# GRAPHICS PROCESSING UNITS



# GRAPHICS PROCESSING UNITS



Originally specialized for **graphical processing**, optimized for large datasets and **high throughput**:

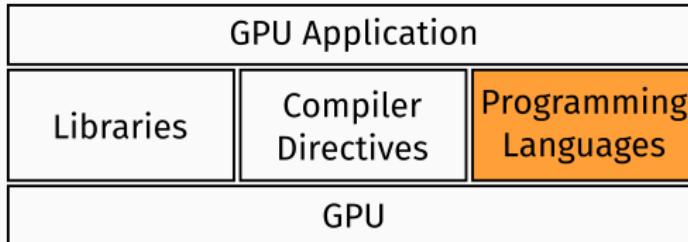
- Small caches (*kilobytes*)
- No **branch prediction**
- Thousands of **higher latency** ALUs
- Execution **pipelines**
- 114 688 concurrent **threads** in Pascal architecture

# CUDA PLATFORM



- A platform for **parallel computing**
- *Application Programming Interface* (API)
- *CUDA Toolkit*

# PROGRAMMING LANGUAGES



- Best performance, but hard to master
- Flexible



## Programming Model:

- Kernels
- Threads hierarchy
- Memory hierarchy
- Heterogeneous Programming

## EXAMPLE: VECTOR SUM

```
#include <cuda_runtime.h>

float *h_A = (float *) malloc(size);
if (h_A == NULL) { ... };

float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);
err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
if (err != cudaSuccess) { ... };

int threadsPerBlock = 256;
int blocksPerGrid = (numElements + threadsPerBlock - 1) /
    threadsPerBlock;

vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C,
    numElements);

err = cudaGetLastError()
err = cudaDeviceSynchronize()
if (err != cudaSuccess) { ... };

err = cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
err = cudaFree(d_A);
if (err != cudaSuccess) { ... };
```

# HANDS ON!

We'll run some CUDA C samples available at  
[github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda):

- `src/cuda-samples/0_Simple/vectorAdd`
- `src/cuda-samples/5_Simulations/fluidsGL`
- `src/cuda-samples/5_Simulations/oceanFFT`
- `src/cuda-samples/5_Simulations/smokeParticles`
- `src/mandelbrot_numba`

## EXAMPLE: PRACTICAL TIP

Always check for errors!

```
float *d_A = NULL;
err = cudaMalloc((void **) &d_A, size);

err = cudaGetLastError()

err = cudaDeviceSynchronize()

if (err != cudaSuccess) {
    fprintf(stderr, "Failed to allocate device vector A
(error code %s)!\n", cudaGetStringError(err));
    exit(EXIT_FAILURE);
}
```

Source: [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda) [Accessed in 29/07/16]

# NVIDIA CUDA COMPILER (nvcc)



# NVIDIA CUDA COMPILER (nvcc)



- Proprietary system

# NVIDIA CUDA COMPILER (nvcc)



- Proprietary system
- Uses the host's C++ compiler

# NVIDIA CUDA COMPILER (nvcc)



- Proprietary system
- Uses the host's C++ compiler
- Compiles CUDA code to *Parallel Thread Execution* (PTX ISA)

# NVIDIA CUDA COMPILER (nvcc)



- Proprietary system
- Uses the host's C++ compiler
- Compiles CUDA code to *Parallel Thread Execution* (PTX ISA)
- Host's code + Device's code → fatbinary

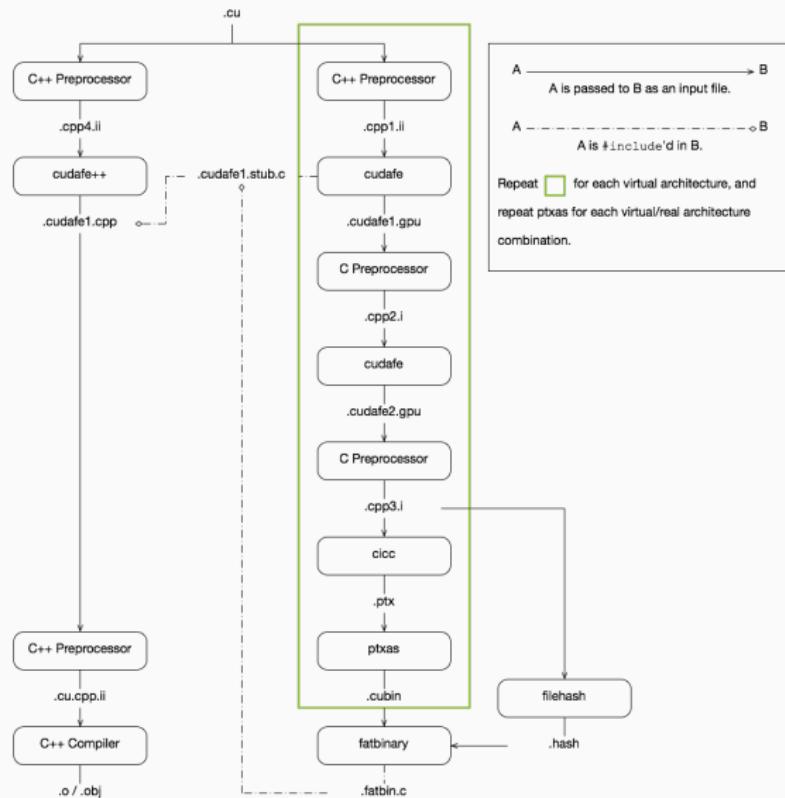
## nvcc: CUDA PROGRAMMING MODEL

- SIMD, SIMT
- *Single Program Multiple Data* (SPMD)

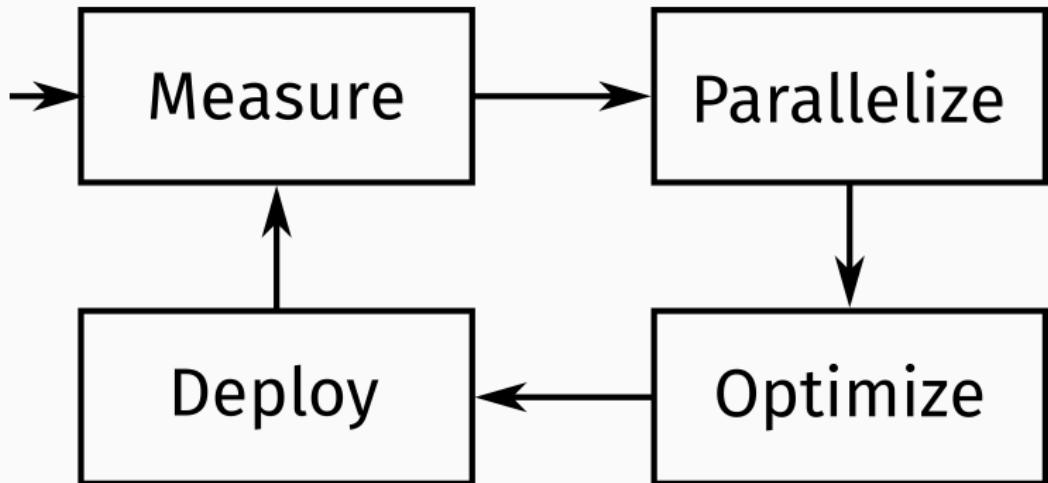
## nvcc: CUDA PROGRAMMING MODEL

- SIMD, SIMT
- *Single Program Multiple Data* (SPMD)
- Device: parallel and independent threads
- Host: does not interfere

# nvcc: COMPILE TRAJECTORY

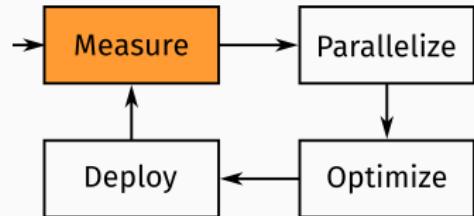


# MEASURE, PARALLELIZE, OPTIMIZE, DEPLOY

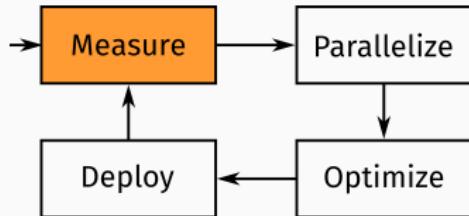


Source: [docs.nvidia.com/cuda/cuda-c-best-practices-guide](http://docs.nvidia.com/cuda/cuda-c-best-practices-guide) [Accessed in 29/07/16]

# MEASURE



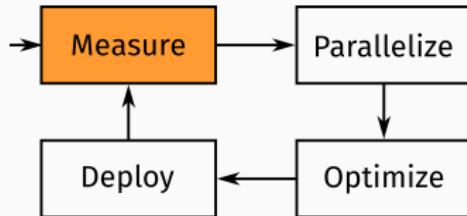
# MEASURE



Given a **sequential program**:

- Which functions do the **heavy lifting** (**bottlenecks**)?

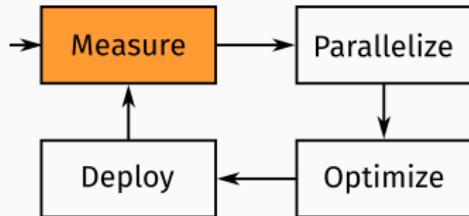
# MEASURE



Given a **sequential program**:

- Which functions do the **heavy lifting** (**bottlenecks**)?
- Can they be **parallelized**?

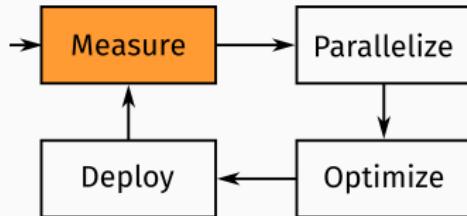
# MEASURE



Given a **sequential program**:

- Which functions do the **heavy lifting** (**bottlenecks**)?
- Can they be **parallelized**?
- How does your program behaves when:
  - The work is divided between more processes? (**strong scaling**)

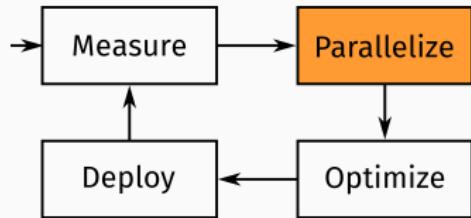
# MEASURE



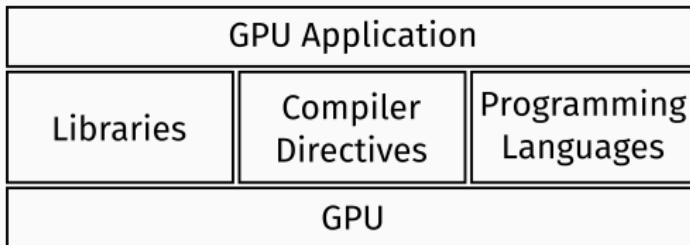
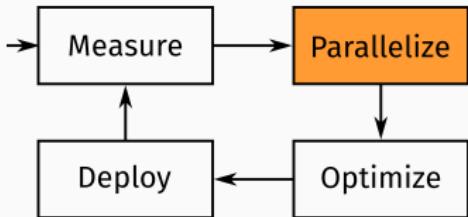
Given a **sequential program**:

- Which functions do the **heavy lifting** (**bottlenecks**)?
- Can they be **parallelized**?
- How does your program behaves when:
  - The work is divided between more processes? (**strong scaling**)
  - There is more work to be done? (**weak scaling**)

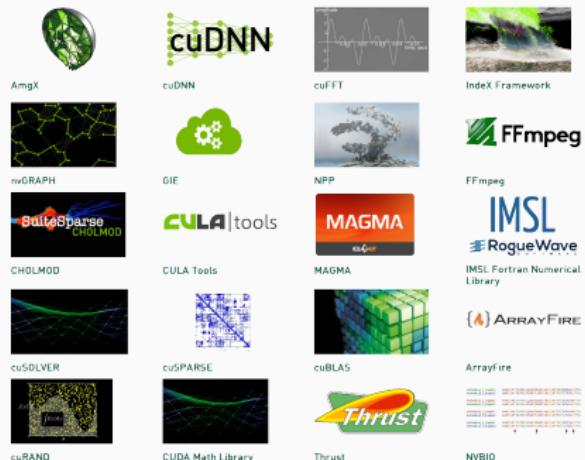
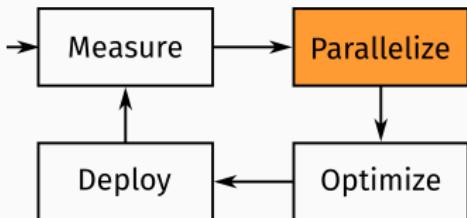
# PARALLELIZE



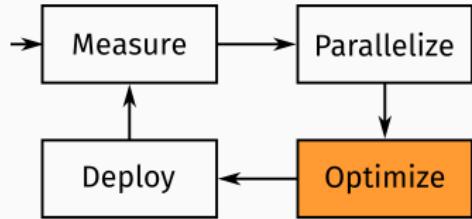
# PARALLELIZE



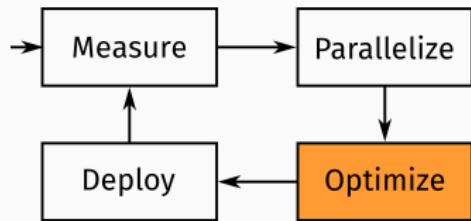
# PARALLELIZE



# OPTIMIZE



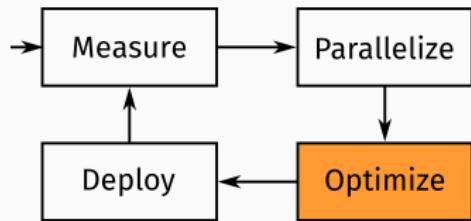
# OPTIMIZE



Optimization:

- Cyclical

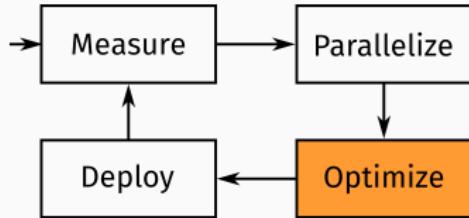
# OPTIMIZE



Optimization:

- Cyclical and incremental

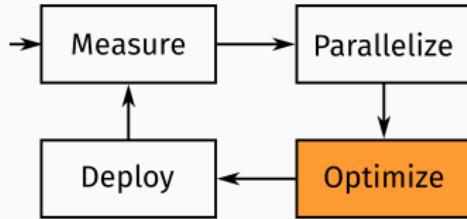
# OPTIMIZE



Optimization:

- Cyclical and incremental
- Applicable to different levels

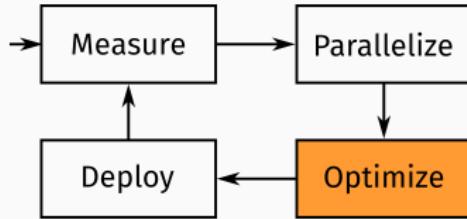
# OPTIMIZE



Optimization:

- Cyclical and incremental
- Applicable to different levels
- Tools are very important

# OPTIMIZE

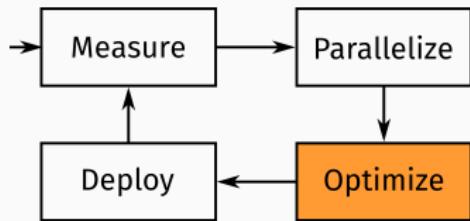


Optimization:

- Cyclical and incremental
- Applicable to different levels
- Tools are very important

Efficiency with algorithms,

# OPTIMIZE

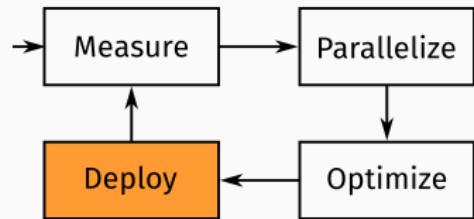


Optimization:

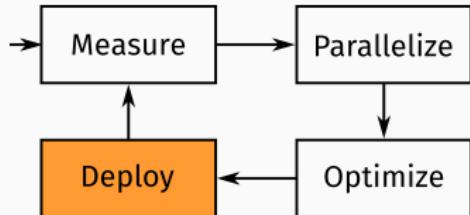
- Cyclical and incremental
- Applicable to different levels
- Tools are very important

Efficiency with algorithms, performance with data structures

# DEPLOY

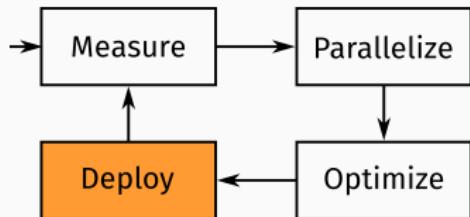


# DEPLOY



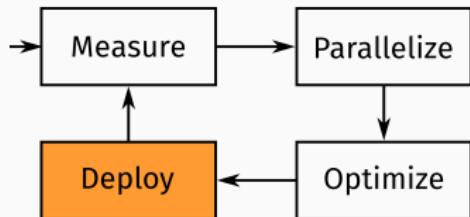
- Prepare application for **measurement**

# DEPLOY



- Prepare application for measurement
- Incremental optimizations and changes

# DEPLOY



- Prepare application for **measurement**
- **Incremental** optimizations and changes
- How much is it still **possible** to optimize?

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential
- Coarse grained

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential
- Coarse grained

Analysis or profiling:

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential
- Coarse grained

Analysis or profiling:

- What are the hotspots (bottlenecks)?

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential
- Coarse grained

Analysis or profiling:

- What are the hotspots (bottlenecks)?
- Can they be parallelized?

# ANALYSING CUDA APPLICATIONS



Parallel and distributed computing is increasingly important, but legacy code is often:

- Sequential
- Coarse grained

Analysis or profiling:

- What are the hotspots (bottlenecks)?
- Can they be parallelized?
- What are the relevant workloads?

## HOST AND DEVICE

Differences between host and device:

- Threading model

## HOST AND DEVICE

Differences between host and device:

- Threading model
- Memory

# HOST AND DEVICE

Differences between host and device:

- Threading model
- Memory

What to run on **GPUs**?

# HOST AND DEVICE

Differences between host and device:

- Threading model
- Memory

What to run on **GPUs**?

- Large **datasets**

# HOST AND DEVICE

Differences between host and device:

- Threading model
- Memory

What to run on **GPUs**?

- Large **datasets**
- **Independent** computations

# HOST AND DEVICE

Differences between host and device:

- Threading model
- Memory

What to run on **GPUs**?

- Large **datasets**
- **Independent** computations

Always consider the impact of **data transferences!**

# *PROFILING*

Monitoring code during **runtime** obtaining information for optimization

# *PROFILING*

Monitoring code during **runtime** obtaining information for optimization

Code **instrumentation**:

- Capture **events**

# PROFILING

Monitoring code during **runtime** obtaining information for optimization

Code **instrumentation**:

- Capture **events**
- Generate **data**

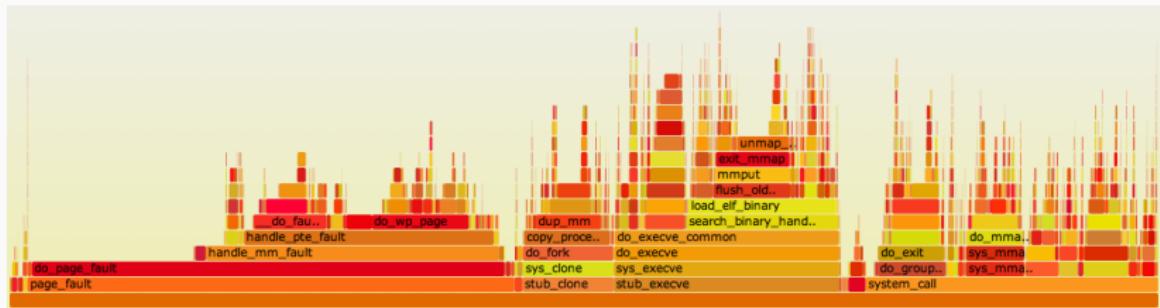
# *PROFILING*

Monitoring code during **runtime** obtaining information for optimization

Code **instrumentation**:

- Capture **events**
- Generate **data**
- Tools

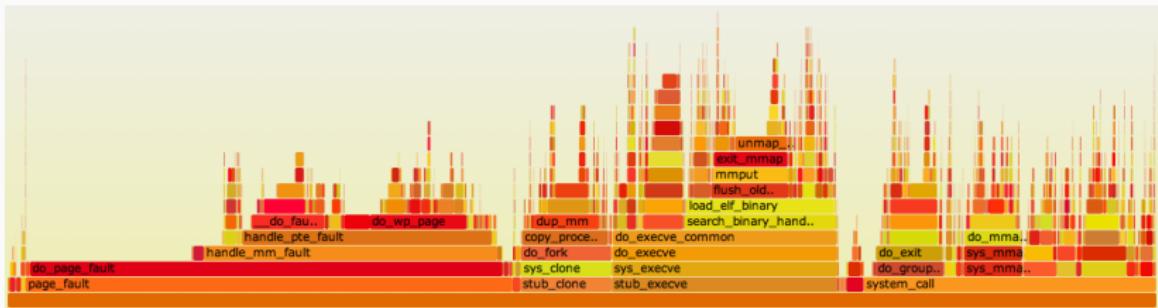
# PROFILING: FLAMEGRAPHS



Help analyse:

- Performance *bottlenecks, hotspots*

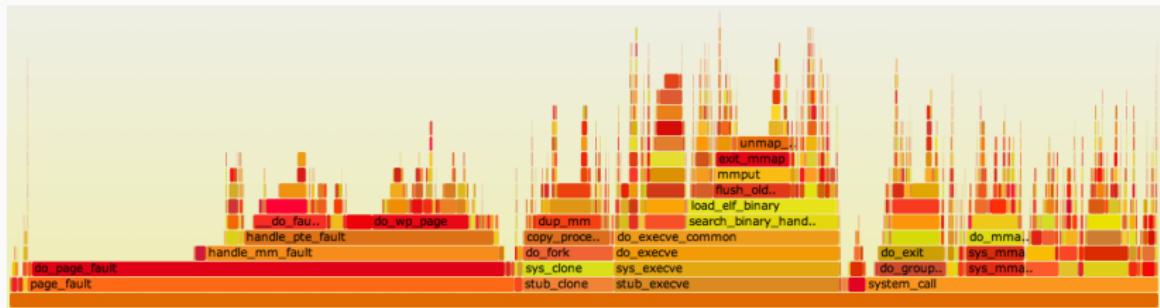
# PROFILING: FLAMEGRAPHS



Help analyse:

- Performance *bottlenecks, hotspots*
- Frequency and duration of **function calls**: *On-CPU*

# PROFILING: FLAMEGRAPHS



Help analyse:

- Performance *bottlenecks, hotspots*
- Frequency and duration of **function calls**: *On-CPU*
- **Idle** time (*waiting*): *Off-CPU*

## PROFILING: nvprof

A **command line** profiler that enables collecting CUDA-related events:

## PROFILING: nvprof

A **command line** profiler that enables collecting CUDA-related events:

- CPU and GPU

## PROFILING: nvprof

A **command line** profiler that enables collecting CUDA-related events:

- CPU and GPU
- Kernel execution

## PROFILING: nvprof

A **command line** profiler that enables collecting CUDA-related events:

- CPU and GPU
- Kernel execution
- Memory usage and transferences

## PROFILING: nvprof

A **command line** profiler that enables collecting CUDA-related events:

- CPU and GPU
- Kernel execution
- Memory usage and transferences
- CUDA API

## PROFILING: nvprof

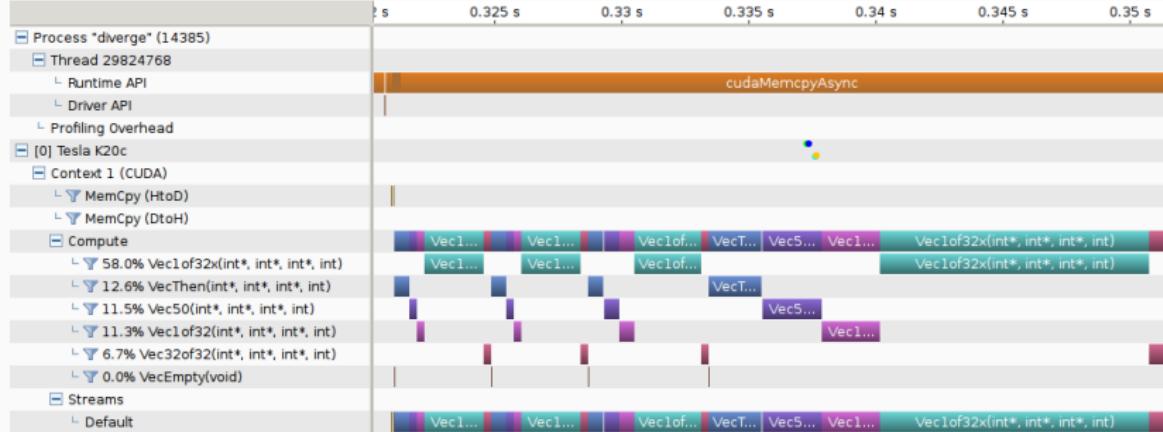
A **command line** profiler that enables collecting CUDA-related events:

- CPU and GPU
- Kernel execution
- Memory usage and transferences
- CUDA API

Generates data for **posterior visualization**

Source: [docs.nvidia.com/cuda/profiler-users-guide](https://docs.nvidia.com/cuda/profiler-users-guide) [Accessed in 29/07/16]

## PROFILING: nvvp

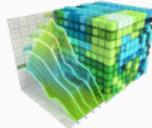


Source: docs.nvidia.com/cuda/profiler-users-guide [Accessed in 29/07/16]

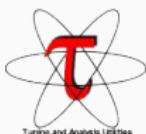
# PROFILING: TOOLS



NVIDIA® Nsight™



NVIDIA Visual Profiler



Tuning and Analysis Utilities

TAU Performance System®



VampirTrace



The PAPI CUDA Component



The NVIDIA CUDA Profiling Tools Interface

Source: [developer.nvidia.com/performance-analysis-tools](http://developer.nvidia.com/performance-analysis-tools) [Accessed in 29/07/16]

## *DEBUGGING*

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

# *DEBUGGING*

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- **Step-by-step** execution

# *DEBUGGING*

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints

# *DEBUGGING*

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints
- Different threads

# *DEBUGGING*

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints
- Different threads
- CUDA **gdb** and **memcheck**

# DEBUGGING

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints
- Different threads
- CUDA **gdb** and **memcheck**

Help fixing:

- Memory **leaks**

# DEBUGGING

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints
- Different threads
- CUDA **gdb** and **memcheck**

Help fixing:

- Memory **leaks**
- Control **flow** problems

# DEBUGGING

Monitoring code during **runtime** obtaining information for **error correction (bugs)**

Code **instrumentation**:

- Step-by-step execution
- Breakpoints
- Different threads
- CUDA **gdb** and **memcheck**

Help fixing:

- Memory **leaks**
- Control **flow** problems
- ...

# OPTIMIZING CUDA APPLICATIONS



Different abstraction and priority levels:

# OPTIMIZING CUDA APPLICATIONS



Different **abstraction** and **priority** levels:

- Memory

# OPTIMIZING CUDA APPLICATIONS



Different abstraction and priority levels:

- Memory
- Control flow

# OPTIMIZING CUDA APPLICATIONS



Different **abstraction** and **priority** levels:

- Memory
- Control flow
- Execution configuration

# OPTIMIZING CUDA APPLICATIONS



Different abstraction and priority levels:

- Memory
- Control flow
- Execution configuration
- Instruction selection

# MEMORY OPTIMIZATIONS

The most important for high performance

## MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations
  - Create and destroy intermediate data structures

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations
  - Create and destroy intermediate data structures
  - Grouped (coalescent) memory accesses

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations
  - Create and destroy intermediate data structures
  - Grouped (coalescent) memory accesses
- Host:

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations
  - Create and destroy intermediate data structures
  - Grouped (coalescent) memory accesses
- Host:
  - Delegate computations

# MEMORY OPTIMIZATIONS

The most important for high performance

Objectives:

- Maximize bandwidth
- Minimize transferences between device and host

Best practices:

- Device:
  - Prioritize computations
  - Create and destroy intermediate data structures
  - Grouped (coalescent) memory accesses
- Host:
  - Delegate computations
  - Minimize data transferences

# CONTROL FLOW OPTIMIZATIONS

High priority for high performance

# CONTROL FLOW OPTIMIZATIONS

High priority for high performance

Objectives:

- Avoid Control Flow divergence inside warps

# CONTROL FLOW OPTIMIZATIONS

High priority for high performance

Objectives:

- Avoid Control Flow divergence inside warps
- `if, switch, do, for, while`

# CONTROL FLOW OPTIMIZATIONS

High priority for high performance

Objectives:

- Avoid Control Flow divergence inside warps
- `if, switch, do, for, while`
- Use `threadIdx` efficiently

# EXECUTION CONFIGURATION OPTIMIZATIONS

How to use the **device**'s potential?

# EXECUTION CONFIGURATION OPTIMIZATIONS

How to use the **device**'s potential?

- Maximize **occupancy**

# EXECUTION CONFIGURATION OPTIMIZATIONS

How to use the **device**'s potential?

- Maximize **occupancy**
- Efficiently allocate **blocks** and **threads**

# EXECUTION CONFIGURATION OPTIMIZATIONS

How to use the `device`'s potential?

- Maximize `occupancy`
- Efficiently allocate `blocks` and `threads`
- Independent `kernel` executions

# INSTRUCTION SELECTION OPTIMIZATIONS

Low priority for high performance:

# INSTRUCTION SELECTION OPTIMIZATIONS

Low priority for high performance:

- Low level

# INSTRUCTION SELECTION OPTIMIZATIONS

Low priority for high performance:

- Low level
- Applied to hotspots

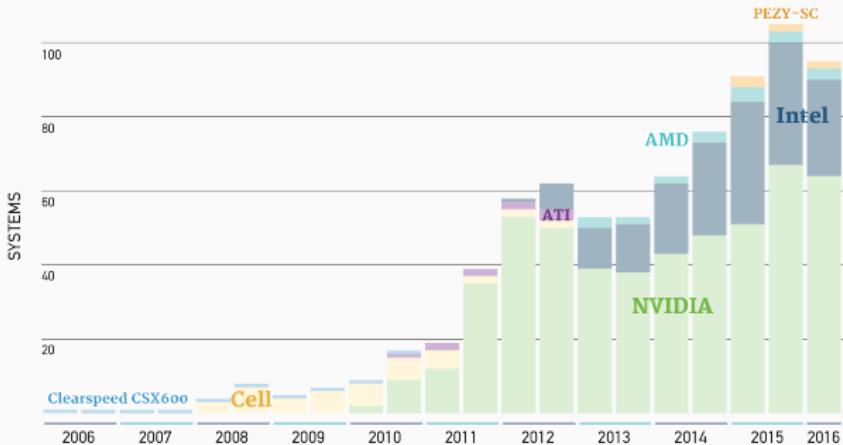
# INSTRUCTION SELECTION OPTIMIZATIONS

Low priority for high performance:

- Low level
- Applied to hotspots
- Applied by experts
- Only after exhausting other optimizations

# CONCLUSION

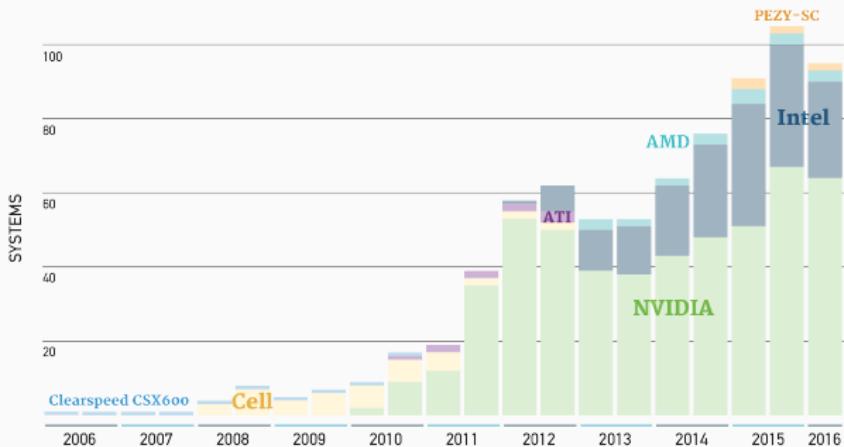
## ACCELERATORS/CO-PROCESSORS



Achieve **high performance** using **accelerators**

# CONCLUSION

## ACCELERATORS/CO-PROCESSORS



Achieve **high performance** using **accelerators (GPUs)**

# CONCLUSION



# CONCLUSION



Implementation, analysis and optimization costs

# CONCLUSION



Implementation, analysis and optimization costs

CUDA Platform tools

# RESOURCES

This presentation and all source code are available at [GitHub](#):

- [github.com/phrb/intro-cuda](https://github.com/phrb/intro-cuda)

More resources:

- CUDA C: [docs.nvidia.com/cuda/cuda-c-programming-guide](https://docs.nvidia.com/cuda/cuda-c-programming-guide)
- CUDA Toolkit: [developer.nvidia.com/cuda-toolkit](https://developer.nvidia.com/cuda-toolkit)
- Best Practices Guide:
  - [docs.nvidia.com/cuda/cuda-c-best-practices-guide](https://docs.nvidia.com/cuda/cuda-c-best-practices-guide)
- GPU Teaching Kit: [syllabus.gputeachingkit.com](https://syllabus.gputeachingkit.com)
- iPython: [ipython.org/notebook.html](https://ipython.org/notebook.html)
- Anaconda: [continuum.io/downloads](https://continuum.io/downloads)

# INTRODUCTION TO GPU PROGRAMMING WITH THE CUDA PLATFORM

---

Pedro Bruel  
[phrb@ime.usp.br](mailto:phrb@ime.usp.br)  
4th INFIERI, 2017



Instituto de Matemática e Estatística  
Universidade de São Paulo