

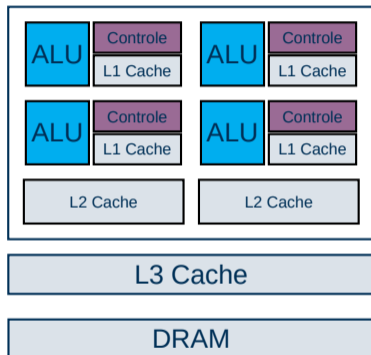
# GPU COMPUTING

## An Introduction

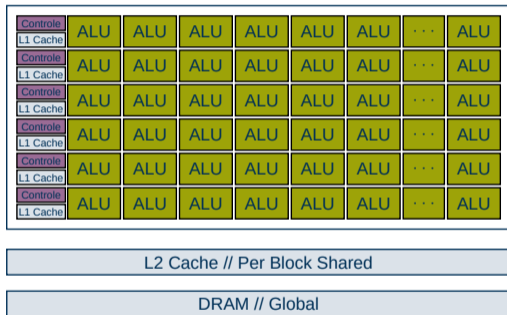
August 10, 2021 | Luis Altenkort, Marius Neumann, Marcel Rodekamp, Christian Schmidt, Xin Wu

- A graphics processing unit (GPU) is a processor featuring a highly parallel structure, making it efficient at processing large blocks of data.
- Large data sets can be worked on by multiple cores. Two methods are commonly used:
  - multi CPU ( $\sim 10 - 100$  cores)
  - GPUs ( $\sim 5000$  cores)
- Where are GPUs used?
  - video/graphics applications
  - linear algebra
  - machine learning
  - general purpose GPU (GPGPU)



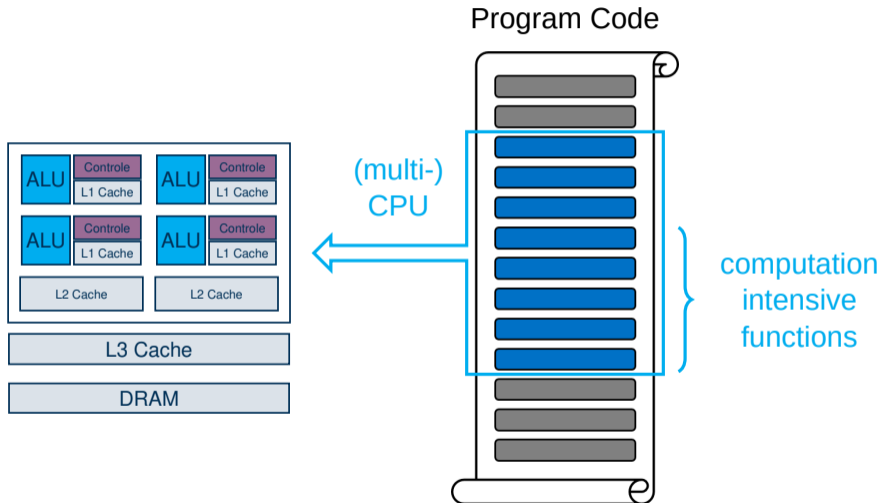


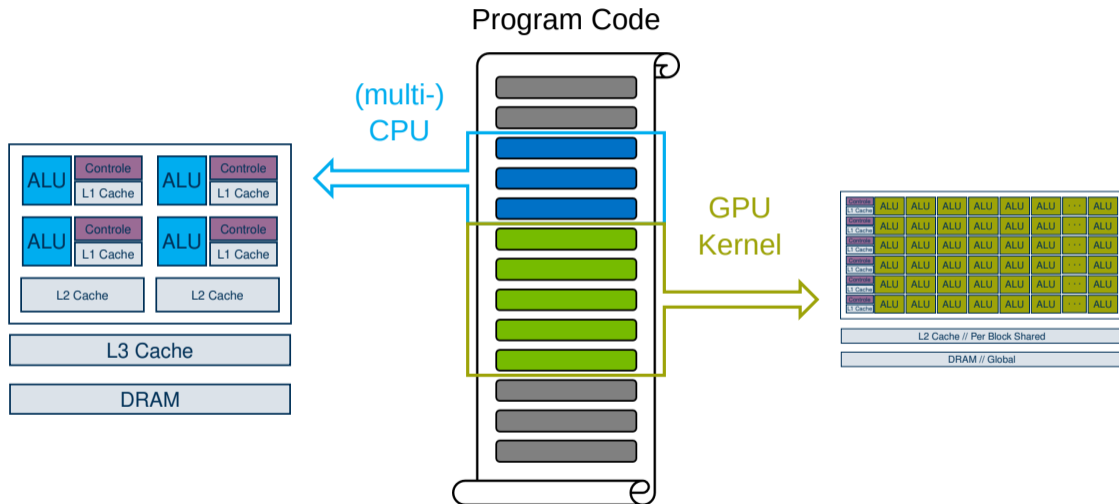
- $\mathcal{O}(10)$  Algorithmic-Logic-Units (ALUs) with broad instruction set
- On Chip:
  - Sophisticated Controle Units
  - Local Per Core L1 Cache
  - Per Core L2 Cache
  - L3 Cache
- Off Chip DRAM
- low latency optimized
- versatile usage



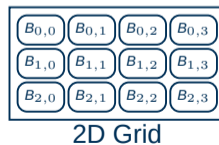
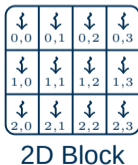
- $\mathcal{O}(10^9)$  ALUs with limited instruction set
- On Chip:
  - Basic Controle Units
  - Local Per Thread L1 Cache
  - Per Block L2 Cache
- Off Chip DRAM
- high throughput optimized
- limited usage

	CPU	GPU
company	Intel	Nvidia
processor	Xeon Silver 4114	Tesla V100
launch (year)	2017	2017
clock freq. (GHz)	2.20 - 3.00	1.4
no. of cores	10	5120
cache (MB)	13.75	6.0
memory size (GB)	404	32
TDP (W)	85	300
SP peak (TFLOP/s)	> 0.7	> 15.7
bandwidth (GB/s)	4	300





- The kernel spawns Threads, Blocks and Grids
  - Thread: Smallest executing unit
    - ALU with limited instruction set, but many available
  - Block: collection of threads
    - share memory
  - Grid: collection of blocks
    - do not share memory

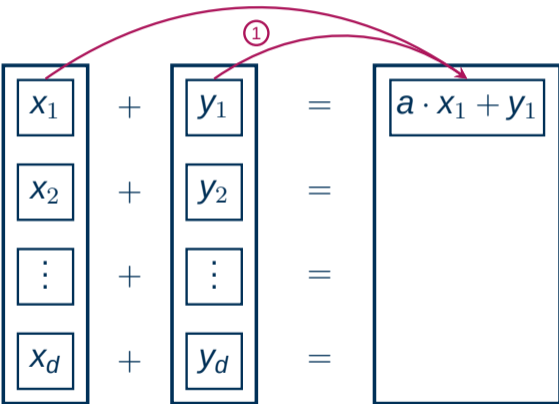




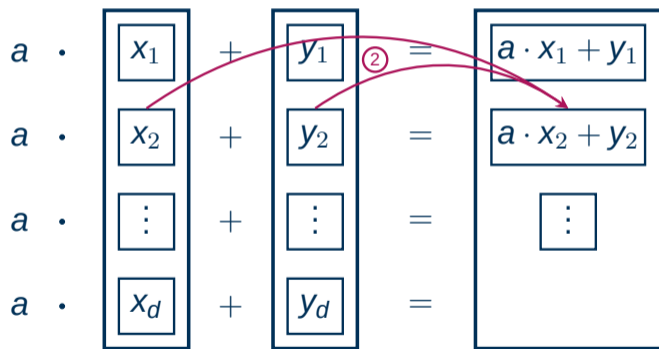
$$a \cdot \begin{array}{|c|} \hline x_1 \\ \hline x_2 \\ \hline \vdots \\ \hline x_d \\ \hline \end{array} + \begin{array}{|c|} \hline y_1 \\ \hline y_2 \\ \hline \vdots \\ \hline y_d \\ \hline \end{array} = \begin{array}{|c|} \hline a \cdot x_1 + y_1 \\ \hline a \cdot x_2 + y_2 \\ \hline \vdots \\ \hline a \cdot x_d + y_d \\ \hline \end{array}$$

- SAXPY: **S**ingle Precision **A** times **X** Plus **Y**
  - Compute  $a \cdot \vec{x} + \vec{y}$

$$\begin{array}{r}
 a \cdot \\
 a \cdot \\
 a \cdot \\
 a \cdot
 \end{array}
 \begin{array}{|c|}
 \hline x_1 \\
 \hline
 \end{array}
 +
 \begin{array}{|c|}
 \hline y_1 \\
 \hline
 \end{array}
 =
 \begin{array}{|c|}
 \hline a \cdot x_1 + y_1 \\
 \hline
 \end{array}$$



- SAXPY: **S**ingle Precision **A** times **X** Plus **Y**
  - Compute  $a \cdot \vec{x} + \vec{y}$
- Sequential
  1. Compute  $a \cdot x_1 + y_1$



– SAXPY: **S**ingle Precision **A** times **X** Plus **Y**

– Compute  $a \cdot \vec{x} + \vec{y}$

– Sequential

1. Compute  $a \cdot x_1 + y_1$

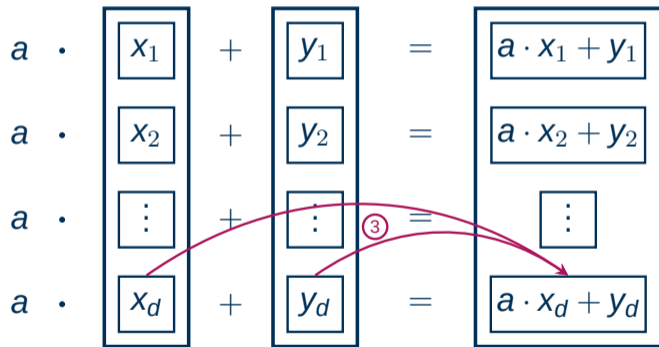
2. Compute  $a \cdot x_2 + y_2$

...

$$\begin{array}{r}
 a \cdot \\
 a \cdot \\
 a \cdot \\
 a \cdot
 \end{array}
 \begin{array}{|c|}
 \hline x_1 \\
 \hline
 \end{array}
 +
 \begin{array}{|c|}
 \hline y_1 \\
 \hline
 \end{array}
 =
 \begin{array}{|c|}
 \hline a \cdot x_1 + y_1 \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 a \cdot \\
 a \cdot \\
 a \cdot \\
 a \cdot
 \end{array}
 \begin{array}{|c|}
 \hline x_2 \\
 \hline
 \end{array}
 +
 \begin{array}{|c|}
 \hline y_2 \\
 \hline
 \end{array}
 =
 \begin{array}{|c|}
 \hline a \cdot x_2 + y_2 \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 a \cdot \\
 a \cdot \\
 a \cdot \\
 a \cdot
 \end{array}
 \begin{array}{|c|}
 \hline \vdots \\
 \hline
 \end{array}
 +
 \begin{array}{|c|}
 \hline \vdots \\
 \hline
 \end{array}
 =
 \begin{array}{|c|}
 \hline \vdots \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 a \cdot \\
 a \cdot \\
 a \cdot \\
 a \cdot
 \end{array}
 \begin{array}{|c|}
 \hline x_d \\
 \hline
 \end{array}
 +
 \begin{array}{|c|}
 \hline y_d \\
 \hline
 \end{array}
 =
 \begin{array}{|c|}
 \hline a \cdot x_d + y_d \\
 \hline
 \end{array}$$


– SAXPY: **S**ingle Precision **A** times **X** Plus **Y**

– Compute  $a \cdot \vec{x} + \vec{y}$

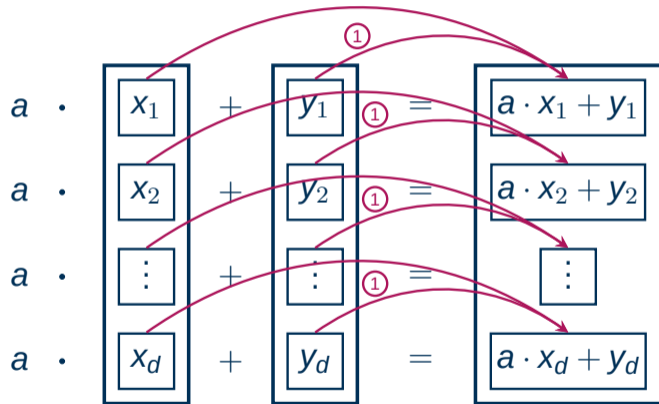
– Sequential

1. Compute  $a \cdot x_1 + y_1$

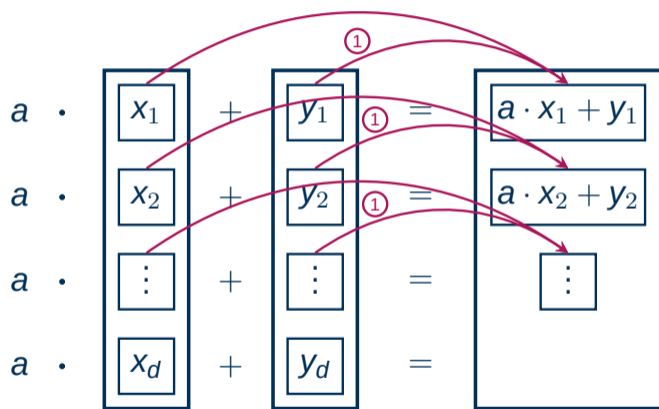
2. Compute  $a \cdot x_2 + y_2$

...

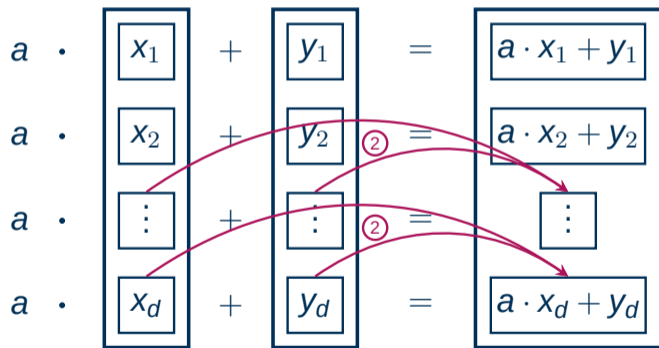
3. Compute  $a \cdot x_d + y_d$



- SAXPY: **S**ingle Precision **A** times **X** Plus **Y**
  - Compute  $a \cdot \vec{x} + \vec{y}$
- Sequential
- Parallel
  1. Compute  $a \cdot x_i + y_i$  for all  $i \in [1, d]$  at once



- SAXPY: **S**ingle Precision **A** times **X** Plus **Y**
- Sequential
- Parallel
  1. Compute  $a \cdot x_1 + y_1$
  - ...
  - Compute  $a \cdot x_4 + y_4$



– SAXPY: **S**ingle Precision **A** times **X** Plus **Y**

– Sequential

– Parallel

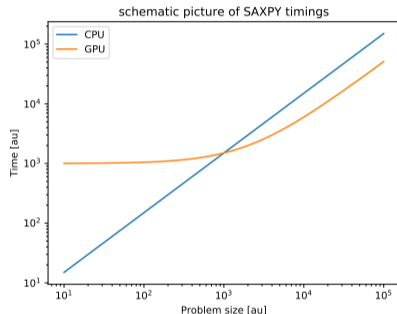
2. Compute  $a \cdot x_5 + y_5$

...

Compute  $a \cdot x_8 + y_8$

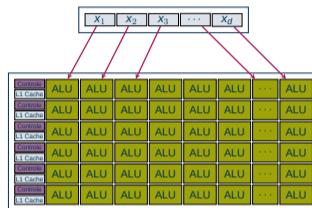
3. ...

- What speedups do we expect for CPU and GPU?
- CPU
  - neglectible overhead
  - linear
- GPU
  - high overhead
  - faster than CPU once overhead is small compared to problem size

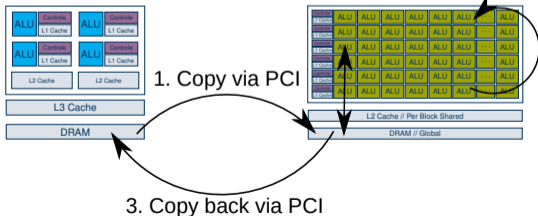




- Why is the GPU sometimes slower?
  - A GPU core is less performant than a CPU core
  - Data copy required
  - Arithmetic intensity
  - SAXPY requires too less flops
- ⇒ Memory bound



## 2. Execute the device program



- Several languages we are going to look at:
  - cuBLAS
  - openACC
  - OpenMP
  - Thrust
  - CUDA C++
  - CUDA Fortran
  - Julia + CUDA.jl
  - CUDA Python
  - OpenCL
  - Sycl
  - OneAPI
  - Magma
  - Kokkos