

GPU Programming

曾駿馳 Chun-Chih Tseng : Maintainer in Liger-Kernel

劉立行 Li-Hsing Liu : Maintainer in Liger-Kernel

Outline

- Why GPU Programming
- Why Triton
- Hands-on Labs
 - Vector Add – Triton Basics
 - Grayscale – Indexing
 - Fused Softmax – Kernel Fusion
 - Online Softmax – Better Algorithm

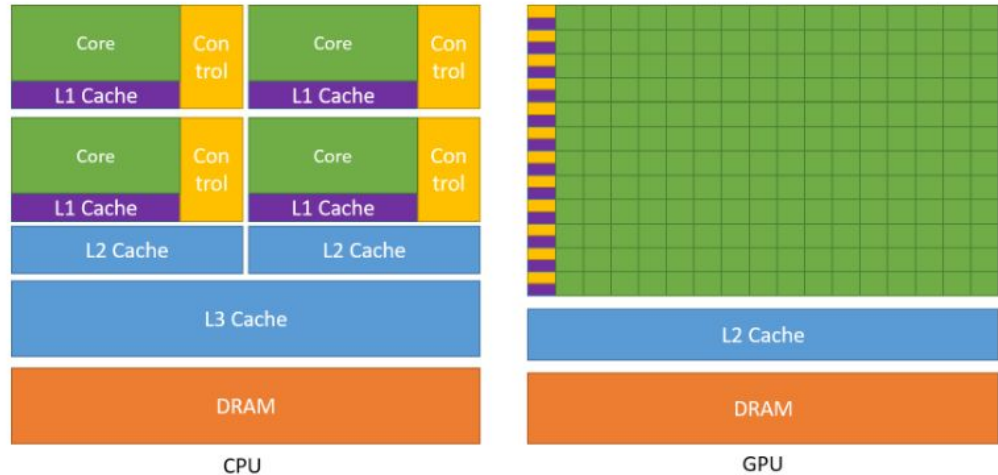
Why GPU Programming – CPU vs GPU

- CPU

- Sequential
- Latency oriented

- GPU

- Parallel
- Throughput oriented
- FLOPs



<https://docs.nvidia.com/cuda/cuda-c-programming-guide/#the-benefits-of-using-gpus>

Why GPU Programming – GPU Architecture

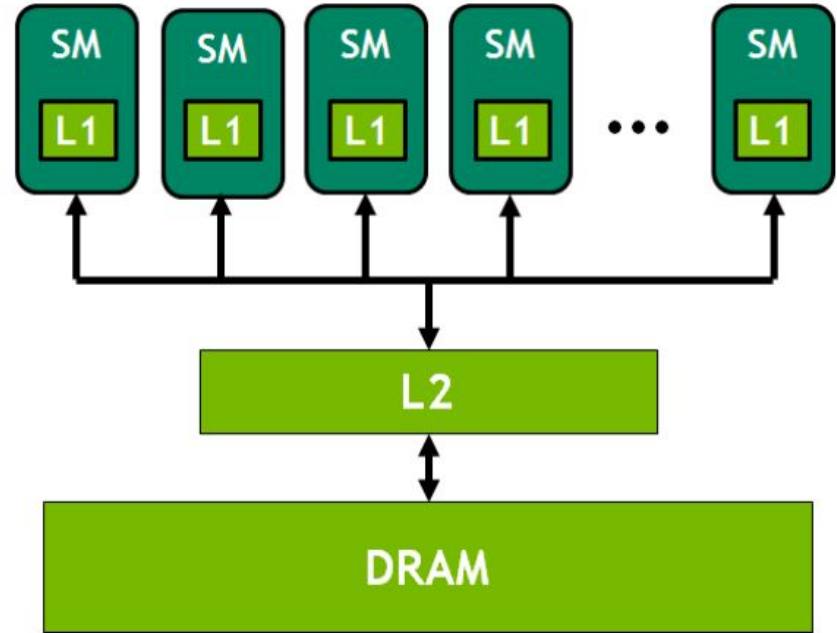


Figure 7. GA100 Streaming Multiprocessor (SM)

<https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9926-tensor-core-performance-the-ultimate-guide.pdf>

Why Triton – Triton vs Cuda



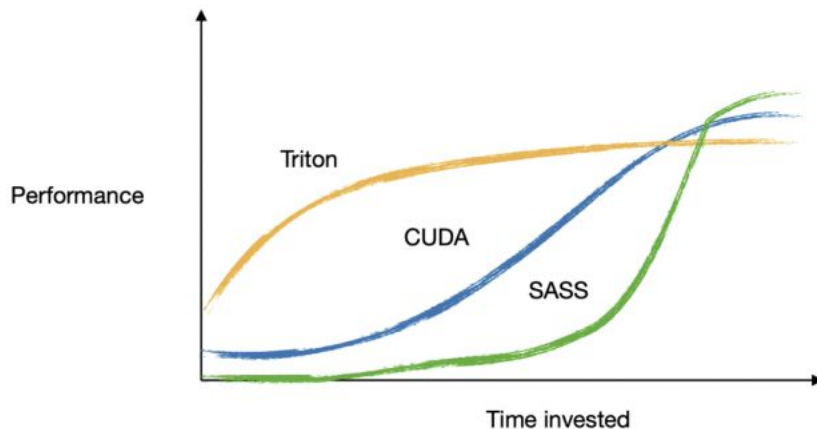
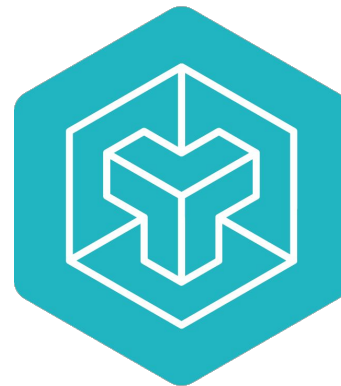
- open-source
- Nvidia, AMD & Intel
- Python
- 90% as fast
- easier



- closed-source
- Nvidia GPUs only
- C/C++
- gold standard for speed
- more difficult

Why Triton

- Automated optimization
- Think in Numpy
- Fast iteration



Triton Basics

2D Indexing

h = 4, w = 7

	0	1	2	3	4	5	6
0	0	1	2	3	4	5	6
1	7	8	9	10	11	12	13
2	14	15	16	17	18	19	20
3	21	22	23	24	25	26	27

offset_0

2	3
---	---

 (1d)

offset_1

4	5
---	---

 (1d)

offset_0[:,None]

2
3

 (2d)

offset_1[None,:]

4	5
---	---

 (2d)

w * offset_0[:,None]

14
21

 (2d)

w * offset_0[:,None] + offset_1[None,:]

18	19
25	26

 (2d)

GPU Performance

GPUs have processing elements (SMs), on-chip memories (e.g. L2 cache), and off-chip DRAM

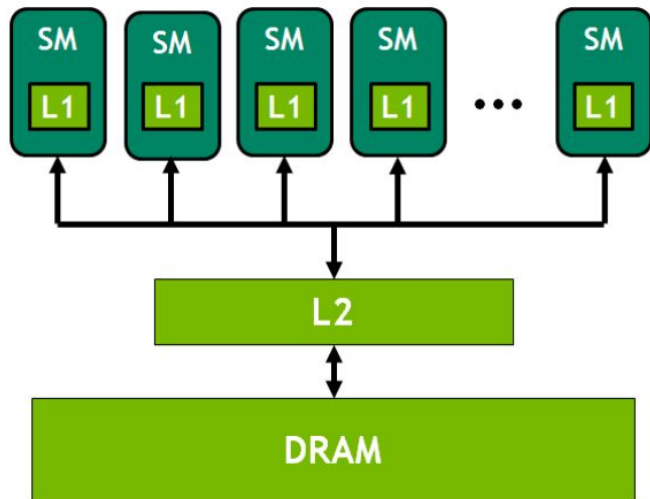
Tesla V100: 125 TFLOPS, 900 GB/s DRAM

What limits the performance of a computation?

$$time_{\text{math operations}} > time_{\text{data movement}}$$

$$\frac{FLOPS}{\text{math throughput}} > \frac{\text{bytes}}{\text{memory bandwidth}}$$

$$\frac{FLOPS}{\text{bytes}} > \frac{\text{math throughput}}{\text{memory bandwidth}}$$



GPU Performance – Arithmetic intensity

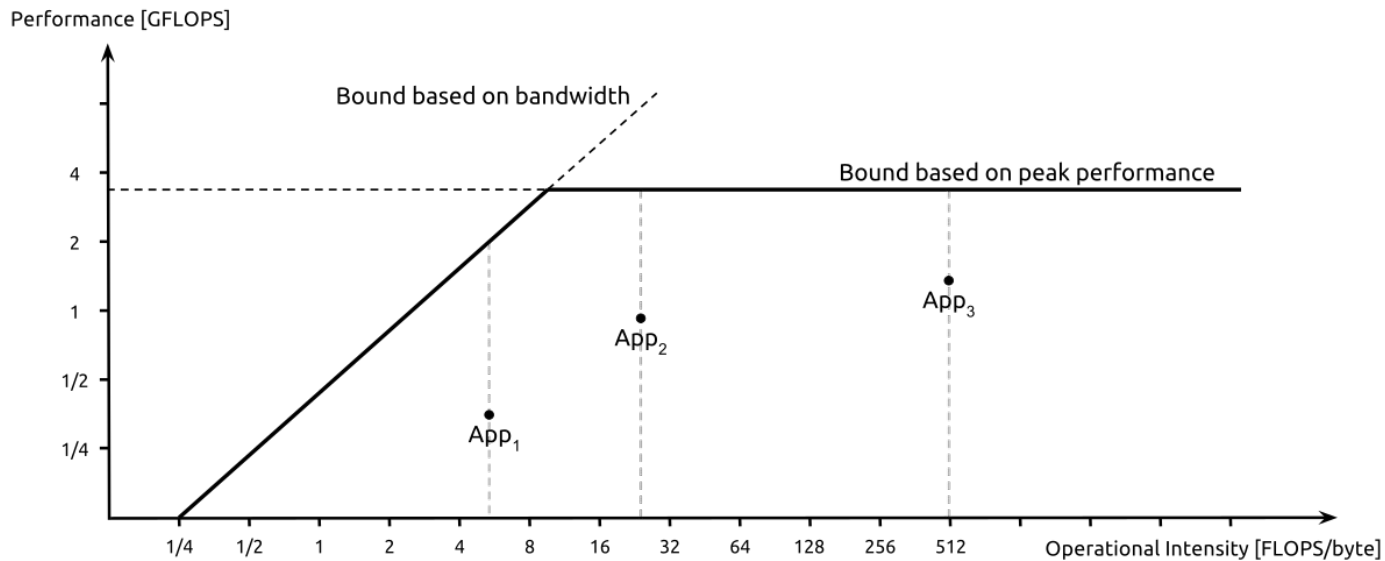
$$\text{Math limited if: } \frac{\text{FLOPS}}{\text{bytes}} > \frac{\text{math throughput}}{\text{memory bandwidth}}$$

Left metric is algorithmic mix of math and memory ops called **arithmetic intensity**

Right metric is the processor's **ops/byte ratio** - e.g. V100 can execute $125/0.9=139$ FLOPS/B

Comparing arithmetic intensity to ops/byte ratio indicates what algorithm is limited by!

GPU Performance – Roofline Model



Online Softmax

Algorithm 2 Safe softmax

```
1:  $m_0 \leftarrow -\infty$ 
2: for  $k \leftarrow 1, V$  do
3:    $m_k \leftarrow \max(m_{k-1}, x_k)$ 
4: end for
5:  $d_0 \leftarrow 0$ 
6: for  $j \leftarrow 1, V$  do
7:    $d_j \leftarrow d_{j-1} + e^{x_j - m_V}$ 
8: end for
9: for  $i \leftarrow 1, V$  do
10:   $y_i \leftarrow \frac{e^{x_i - m_V}}{d_V}$ 
11: end for
```

Algorithm 3 Safe softmax with online normalizer calculation

```
1:  $m_0 \leftarrow -\infty$ 
2:  $d_0 \leftarrow 0$ 
3: for  $j \leftarrow 1, V$  do
4:    $m_j \leftarrow \max(m_{j-1}, x_j)$ 
5:    $d_j \leftarrow d_{j-1} \times e^{m_{j-1} - m_j} + e^{x_j - m_j}$ 
6: end for
7: for  $i \leftarrow 1, V$  do
8:    $y_i \leftarrow \frac{e^{x_i - m_V}}{d_V}$ 
9: end for
```

GPU MODE: <https://www.gpumode.com/>