# OpenCL Workshop GPU architectures

Ljubljana, May 2019
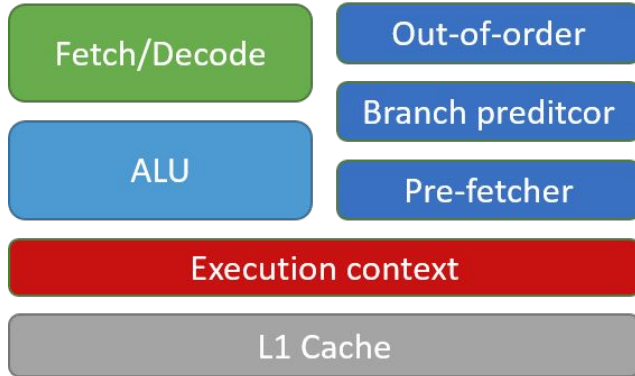
# Outline

- How we got from x86 to GPUs
    - multi-core vs many-core
    - GPU programming model
- OpenCL language
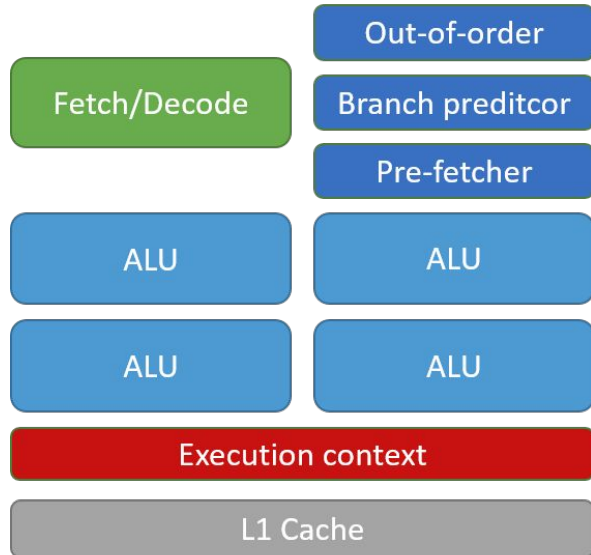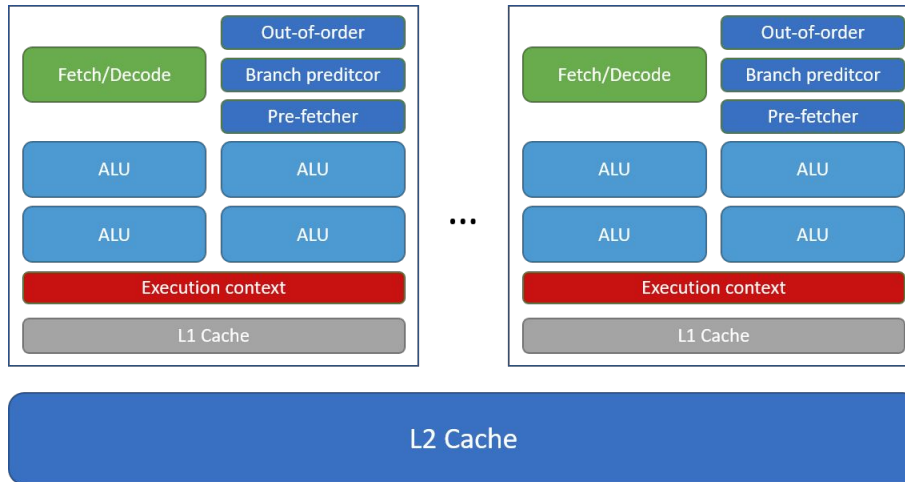    - vector addition
    - image processing

# Single-core CPU

| | |
|---|---|
| Fetch/Decode | Out-of-order |
| | Branch preditcor |
| ALU | Pre-fetcher |
| Execution context | |
| L1 Cache | |

```
int i = 0;
while ( i < N ){
    C[i] = A[i] + B[i];
    i++;
}
```

# SIMD Extensions (1999)

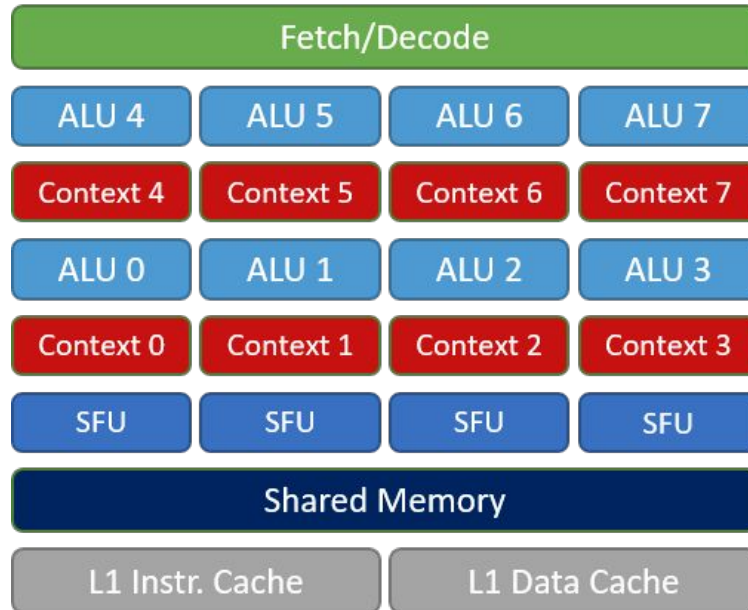| | |
|---|---|
| | Out-of-order |
| Fetch/Decode | Branch preditcor |
| | Pre-fetcher |
| ALU | ALU |
| ALU | ALU |
| Execution context | |
| L1 Cache | |

```
int i = 0;
while ( i < N ){
    C[i] = add4(A[i], B[i]);
    i = i + 4;
}
```
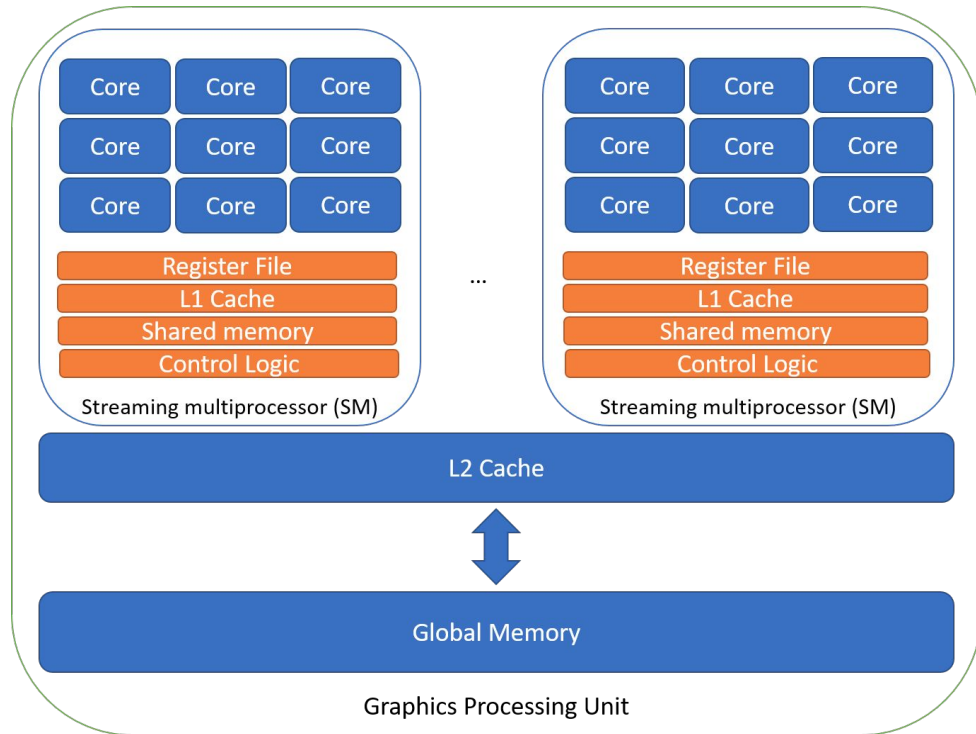
# Multi-core (~2001)



```
i = tID * (N / num_of_cores);
end = (tID +1) * (N /
num_of_cores);
while( i < end ){
        C[i] = A[i] + B[i];
        i++;
}
```

# GPU Streaming multiprocessor (SM)
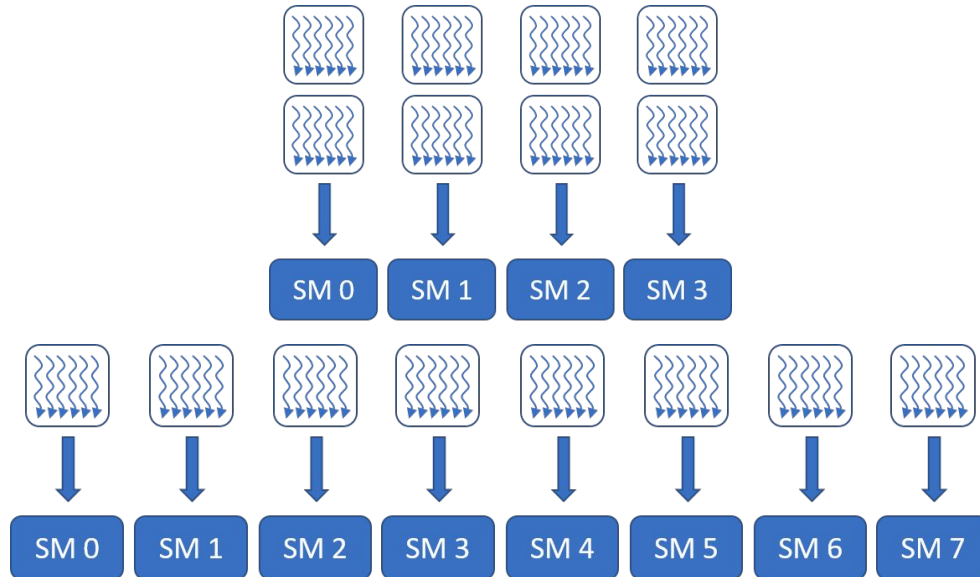
# GPU

# GPU

```
i = tID * (N / num_of_cores); // just a larger number of cores
end = (tID + 1) * (N / num_of_cores);
while( i < end ){
    C[i] = A[i] + B[i];
    i++;
}
```
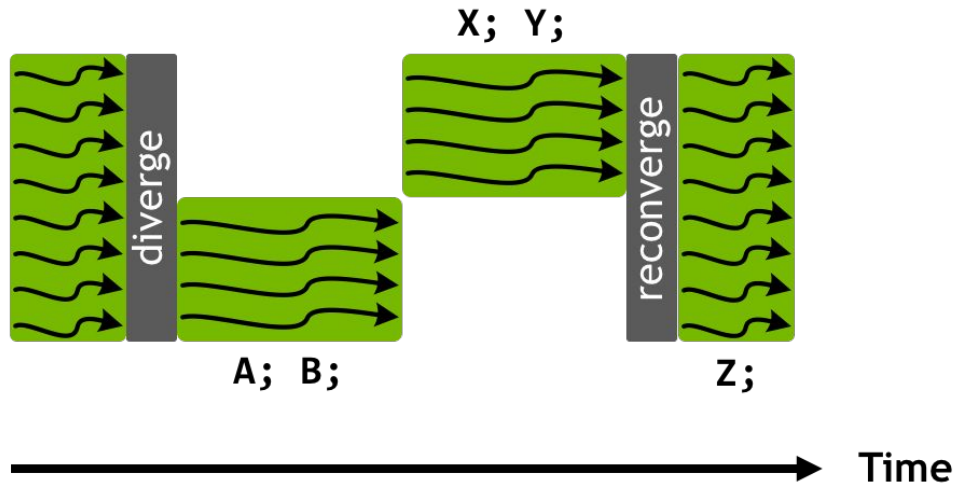
# Thread blocks & scheduling

# Scheduling threads - warps

- SM is issued a single thread-block
  - all threads will be executed in the assigned SM, from start to finish
- Threads are issued to the execution cores in **warps** (typically 32 threads)
  - each instruction issue step 1 warp is selected, usually trying to avoid stalls
- Large number of thread blocks -> parallelism on the SM level
  - all SMs are occupied
- Large thread blocks -> more efficient warp scheduling
  - reduces stalls in the issue steps

# Thread divergence

```
if (threadID < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

# NVIDIA GPU microarchitecture evolution

| Name | Year | # FP32 cores/SM | #FP64 cores/SM | #SM |
|------|------|-----------------|----------------|-----|
| Tesla | 2007 | 8 | 0 | up to 30 |
| Fermi | 2010 | 32 | | up to 32 |
| Kepler | 2012 | 192 | 64 | up to 8 |
| Maxwell | 2014 | 128 | 4 | up to 24 |
| Pascal | 2016 | 64 or 128 | 32 or 4 | up to 30 |
| Volta | 2017 | 64 | 32 | 80 |
| Turing | 2018 | 64 | 2 | 72 |

# GPU architecture trends

Tensor cores

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32     FP16     FP16     FP16 or FP32

The return of scheduling

Improved caches

Mixed precision