# Lecture 19: Parallel architectures

Tuesday, March 13, 2018    10:46 AM

## Outline

- Types of parallel architectures
- Taxonomy for parallel architectures
- Emerging parallel architectures

Why <u>parallelism</u>?

  Incr throughput w/ multiple tasks

  Incr performace/ reduce latency

2 choices → double freq → h/w limits
     or double cores → lower powr

$$\rightarrow P = \tfrac{1}{2} C v^2 f$$

- <u>Scalar</u> → not parallel single inst stream & single data stream

<u>Types of Parallel architectures</u>    →

how to program?

how to communicate?

how are instructions executed?

☆ - Accelerators + CPU   (GPU)

   → send a task to accelerator

   → Send data to accel explicitly

   → <u>offload</u> computation
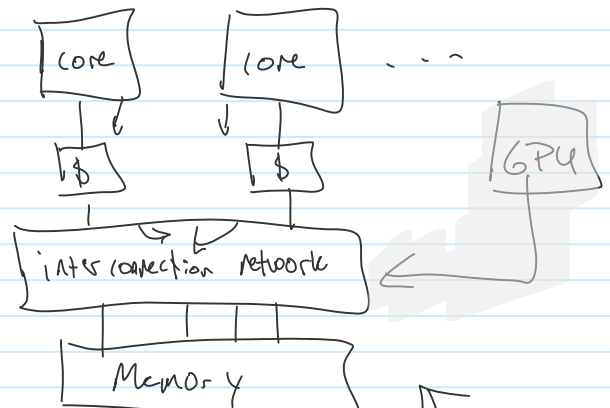
   → explicitly copy results back

☆ <u>Multi core</u>

   → Shared memory
     ↳ because all processors see
      same data at same address

   → reading/ writing memory for comm.

   → Programming?
     threads → has a program counter

→ Programming?

threads → has a program counter
but in the same virtual
address space
usually executing same code

==Single program multiple data== SPMD

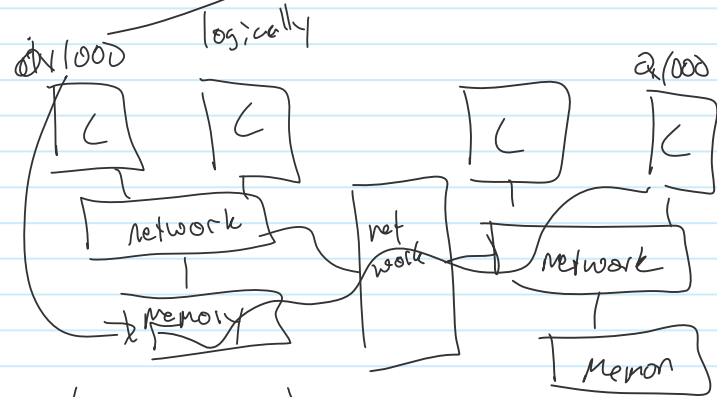SPMD is a type of MIMD system
↳ ==multiple inst. multiple data==

☆ Multi processor / Multi-chip

has multiple processor "sockets"

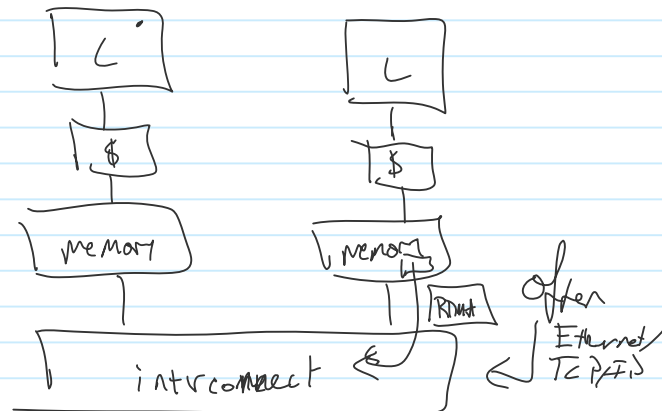Shared memory

- memory access time is non-uniform

- programmers must think about data
locality

☆ Warehouse-scale computers (cloud/datacenter/super computer)

- communicate with ==message passing==
request for data
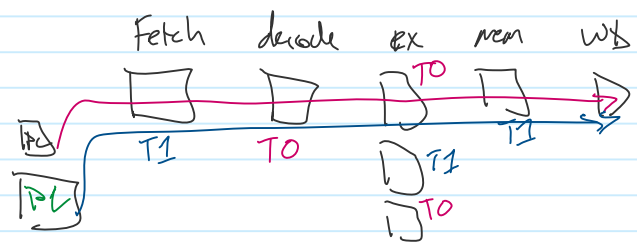respond w/ data

- remote DMA (Direct memory access)

☆ Multi-threaded process (Intel hyperthreading)

4 cores / 8 threads

→ increases throughput
at cost of latency

SPMD
shared memory

☆ - vector architectures

→ take advantage of ==data parallelism==   $\vec{a} + \vec{b} = \vec{c}$

→ take advantage of  ==data parallelism==    $\vec{a} + \vec{b} = \vec{c}$

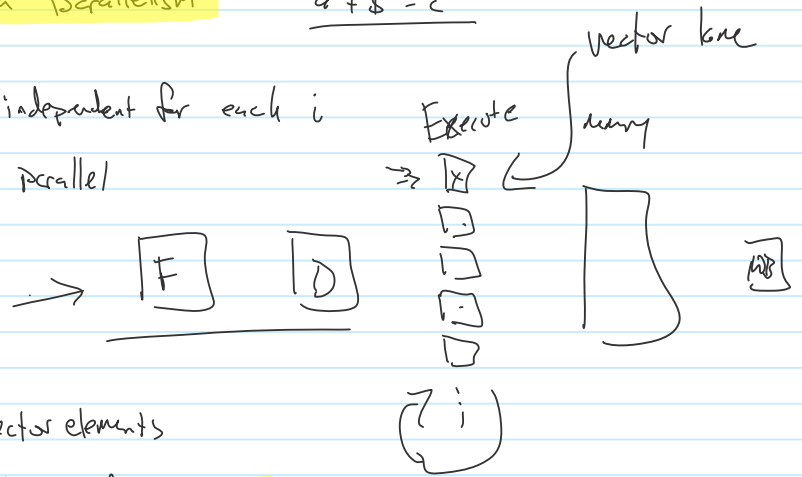for (i = 0 → n)
→    c[i] = a[i] + b[i]    independent for each i

execute all vector elements in parallel

vector lane

Execute } many

→ 

- incr throughput + perf
- easy to reason about

| F | | D |

→

- fetch/decode shared
   decode once execute for all vector elements

(∑ i)

- ==Single instruction stream  multiple data  streams==
   SIMD

Inst. Stream

| | | Single | multiple |
|---|---|---|---|
| data streams | Single | SISD Scalar | MISD ? not really... |
| | multiple | SIMD vectors/ GPUs | MIMD SPMD  threading Shared memory message passing |

## Examples of parallel programs

```
void daxpy(double *X, double *Y, double alpha, const int N)
{
    for (int i = 0; i < N; i++)          scalar
    {
        Y[i] = alpha * X[i] + Y[i];
    }
}

void daxpy_simd(double *X, double *Y, double alpha, const int N)
{
    __m256d *vx = (__m256d*)X;
    __m256d *vy = (__m256d*)Y;
    __m256d va = {alpha, alpha, alpha, alpha};

    for (int i = 0; i < N/4; i++)
    {
        __m256d mult = _mm256_mul_pd(vx[i], va);
        __m256d add = _mm256_add_pd(vy[i], mult);
        vy[i] = add;
    }
}
```

```
}

void daxpy_fmad(double *X, double *Y, double alpha, const int N)
{
    __m256d *vx = (__m256d*)X;
    __m256d *vy = (__m256d*)Y;
    __m256d va = {alpha, alpha, alpha, alpha};

    for (int i = 0; i < N/4; i++)
    {
        vy[i] = _mm256_fmadd_pd(va, vx[i], vy[i]);
    }
}

void daxpy_threads(double *X, double *Y, double alpha, const int N)
{
    int num_threads = std::thread::hardware_concurrency();

    std::thread threads[num_threads];

    int chunk = N / num_threads;

    for (int i=0; i<num_threads; i++) {
        threads[i] = std::thread(daxpy, X+i*chunk, Y+i*chunk, alpha, chunk);   ⟵
    }

    for (int i=0; i<num_threads; i++) {
        threads[i].join();
    }
}

void daxpy_omp(double *X, double *Y, double alpha, const int N)
{
    #pragma omp parallel for
    for (int i = 0; i < N; i++)           Open MP
    {
        Y[i] = alpha * X[i] + Y[i];
    }
}

void daxpy_opencl(double *X, double *Y, double alpha, const int N)
{
    int tid = get_global_id(0);
    for (int i = tid; i < N; i += get_global_size(0)) {          GPU
        Y[i] = alpha * X[i] + Y[i];
    }
}
```
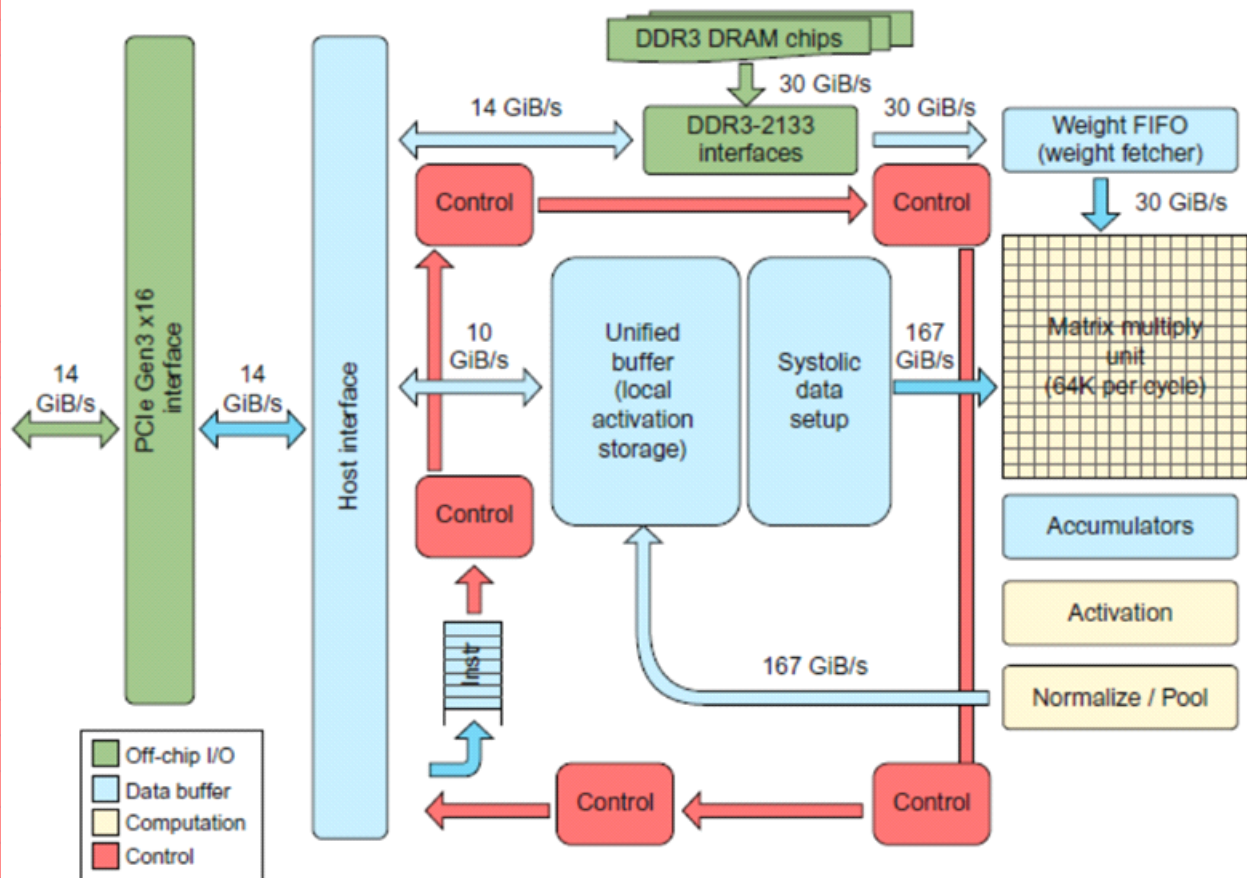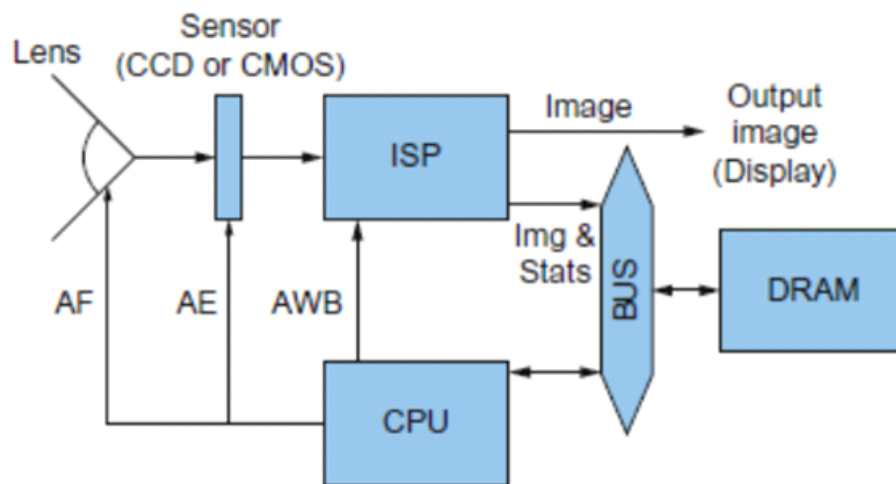
# New parallel architectures

## Google's TPU

## Pixel visual core

5 x 5 stencil