

# Lecture 20: Parallel performance

Thursday, March 15, 2018 10:05 AM

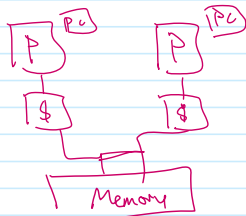
## Outline

- Finish up parallel architectures
- Amdahl's Law
- Some reasons parallelism is hard
- Roofline model

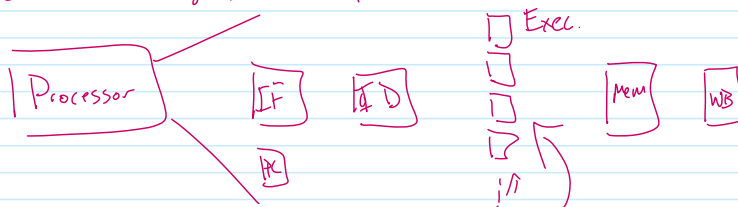
Fract is in wellman 126

Extn office hours Mon. 4:30-6

MIMD  $\rightarrow$  multiple instruction multiple data



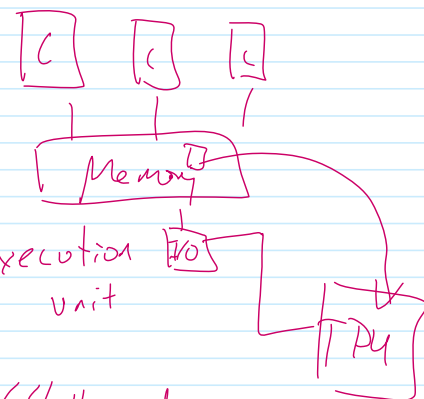
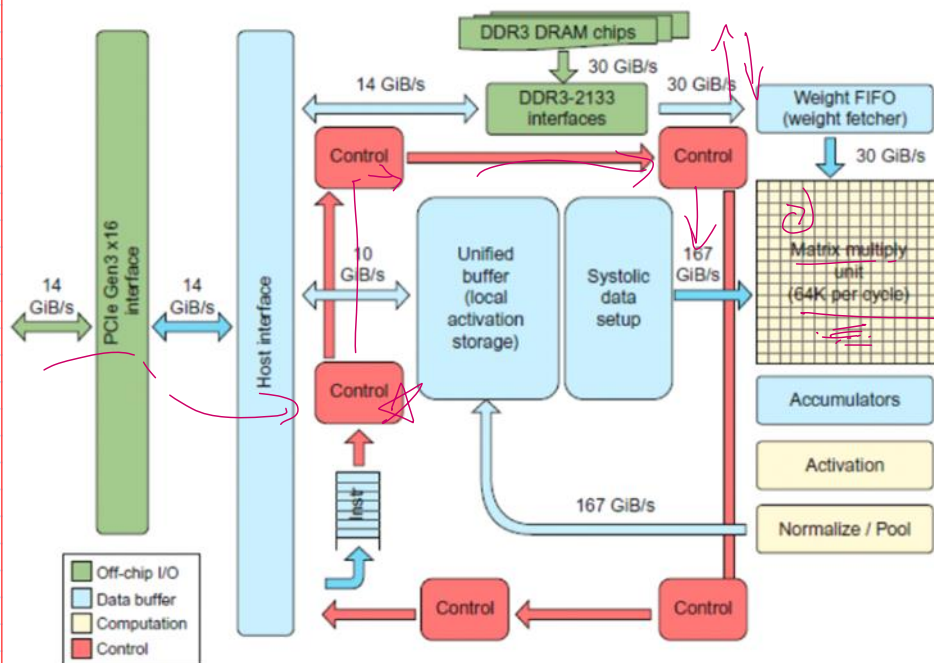
SIMD  $\rightarrow$  single instruction multiple data (vector)



## New parallel architectures

Google's TPU

Tensor Processing Unit



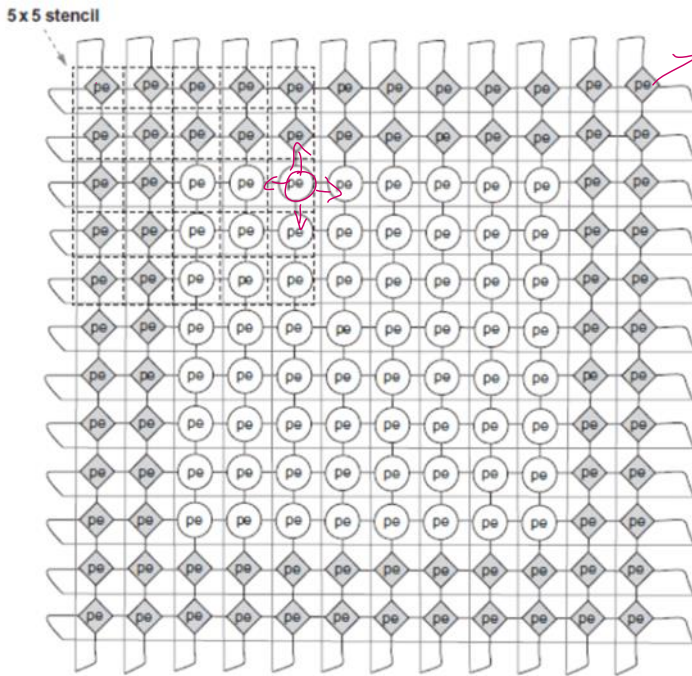
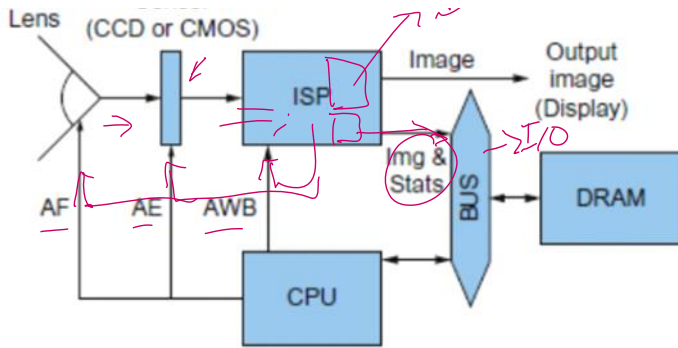
execution unit

64 thousand

MISD

## Pixel visual core





## Amdahl's Law (again)

### Performance

$$\text{Speedup} = \frac{T_{\text{old}}}{T_{\text{new}}} = \frac{T_{1\text{proc.}}}{T_{p\text{proc.}}} = \frac{1}{1 - f_e + \frac{f_e}{s_e}}$$

Diagram illustrating the components of Amdahl's Law:

- Told**: Total time for sequential execution.
- Tnew**: Total time for parallel execution.
- Fraction enhanced**: The portion of the task that can be parallelized.
- Speedup of enhancement**: The speedup achieved by the parallelized portion.

### Parallelism

$$\text{Speedup} = \frac{1}{1 - f + \frac{f}{p}}$$

Where  $p \rightarrow \# \text{ of processors in our system}$   
 $f \rightarrow \text{fraction we can parallelize}$

### 50% parallelizable

Speedup w/ 2 cons?

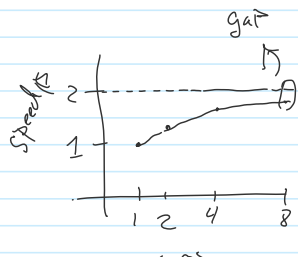
$$S = \frac{1}{1 - 0.5 + \frac{0.5}{2}} = \frac{1}{0.75} = \frac{4}{3} = 1.33x$$

4 cores

$$S = \frac{1}{1 - 0.5 + \frac{0.5}{4}} = \frac{1}{0.625} = \frac{8}{5} = 1.6x$$

8 cores

$$S = \frac{1}{1 - 0.5 + \frac{0.5}{8}} = \frac{1}{0.5625} = \frac{16}{9} = 1.77x$$



$$S = \frac{1}{1.5 + \frac{1}{8}} = \frac{1}{1.625} = \frac{8}{16.25} = \frac{16}{32.5} = \frac{32}{65} = 0.4923x$$

$$8 \text{ cores } S = \frac{1}{1.5 + \frac{1}{8}} = \frac{1}{1.625} = \frac{8}{16.25} = \frac{16}{32.5} = \frac{32}{65} = 0.4923x$$

$$\infty \text{ cores } S = \frac{1}{1.5 + \frac{1}{\infty}} = \frac{1}{1.5} = 0.6667x$$

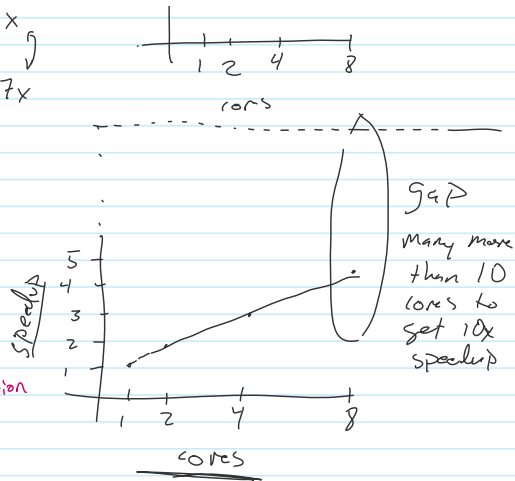
90% parallelizable

$$2 \text{ cores } S = \frac{1}{1.9 + \frac{1}{2}} = \frac{1}{2.4} = 0.4167x$$

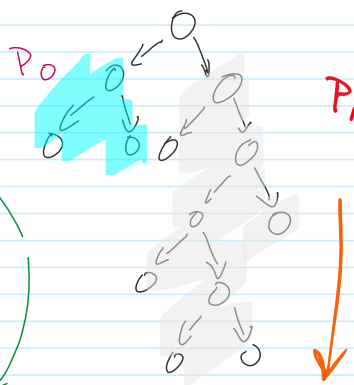
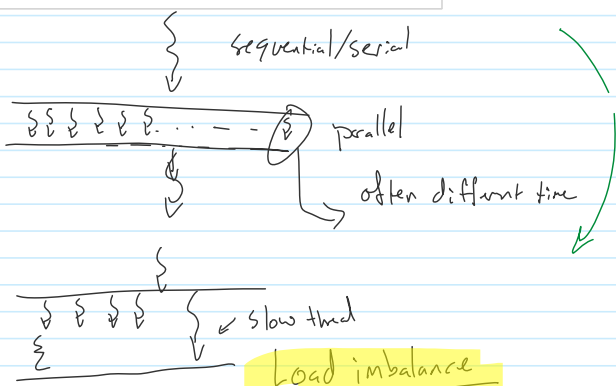
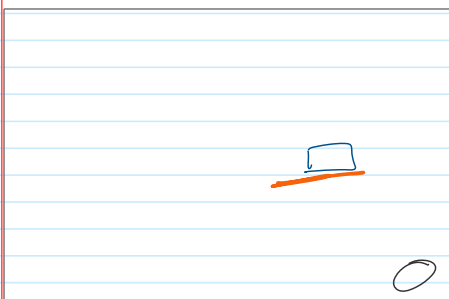
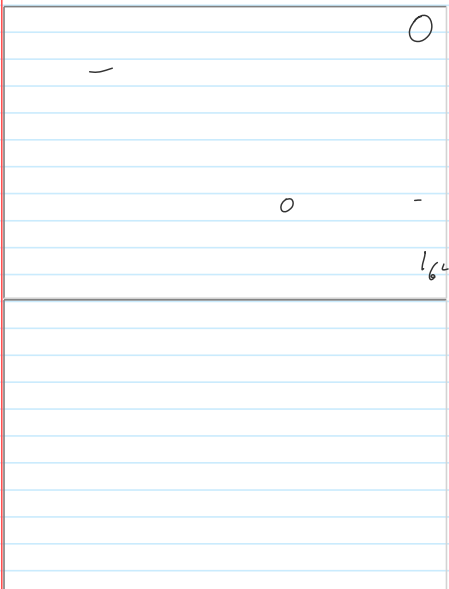
$$4 \text{ cores } S = \frac{1}{1.9 + \frac{1}{4}} = \frac{1}{2.15} = 0.4651x$$

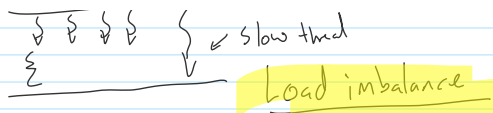
$$8 \text{ cores } S = \frac{1}{1.9 + \frac{1}{8}} = \frac{1}{2.125} = 0.4706x$$

$$\infty \text{ cores peak possible } \rightarrow S = \frac{1}{1.9} = 0.5263x$$



Look @ IPython notebook on github





Solving imbalance

→ make sure all threads have equal partitions

Much worse w/ dynamic parallelism where you "create" work as you go

→ Work queue where all processors put new work & fetch work from here when idle

→ Work stealing (cilk)

## Scaling

→ perform "same" for equal work per processor  
as # processors and amount of work

1 user to 1,000 users → each user has same prod

→ Scaling throughput

→ Weak scaling

## Strong scaling

↳ for a constant problem size incr processors decr latency

Speed → operations per second

↳ floating point

→ FLOPs 6 FLOPs Billions  
7 FLOPs trillions

4 cores 3 SIMD SIMD 256-bit (8 floats) 4 GHz

$4 \cdot 3 \cdot 8 \cdot 4 \cdot 10^9 = 384 \text{ GFLOPs}$  ←

Memory bandwidth for getting data in/out → 34.1 GB/s

## Arithmetic intensity

# flop / byte of data

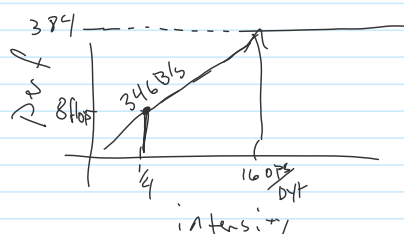
3 ops 3 ld/st

$\frac{3 \text{ ops}}{12 \text{ bytes}}$

4 bytes/float

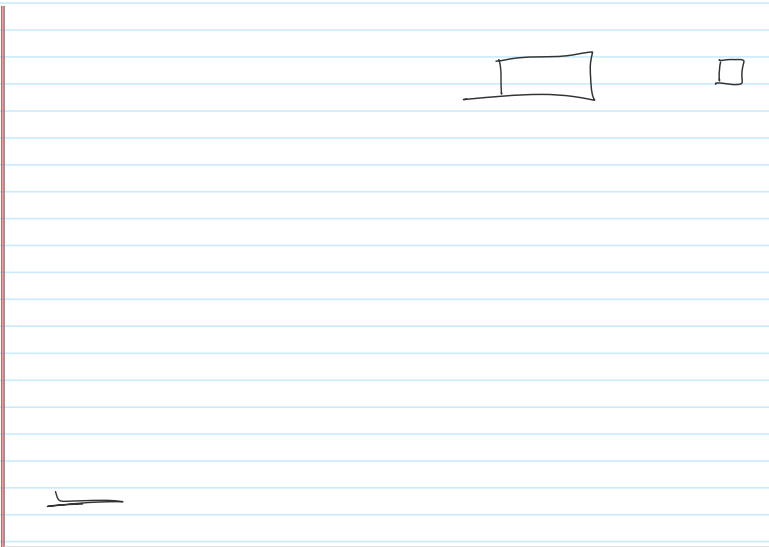
$\frac{1 \text{ op}}{4 \text{ bytes}}$

→



roofline model

Top 500 supercomputers (<https://www.top500.org/lists/2017/11/>)



## Roofline model of performance

