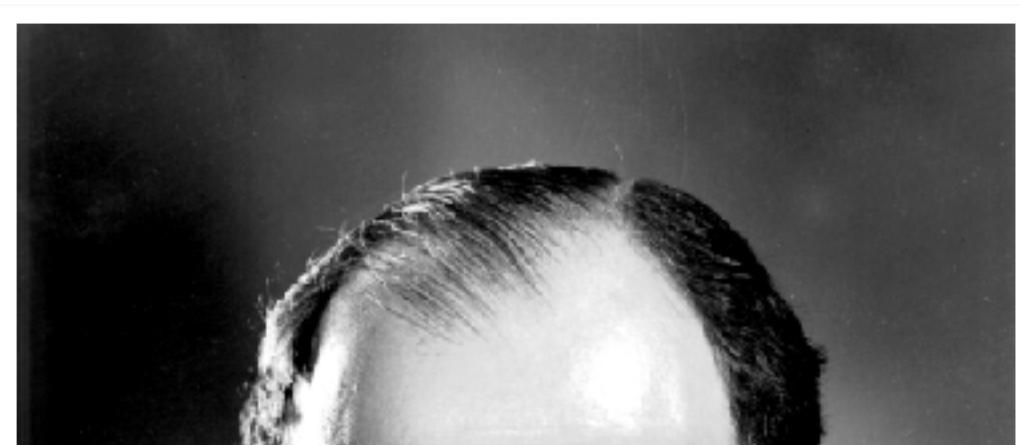


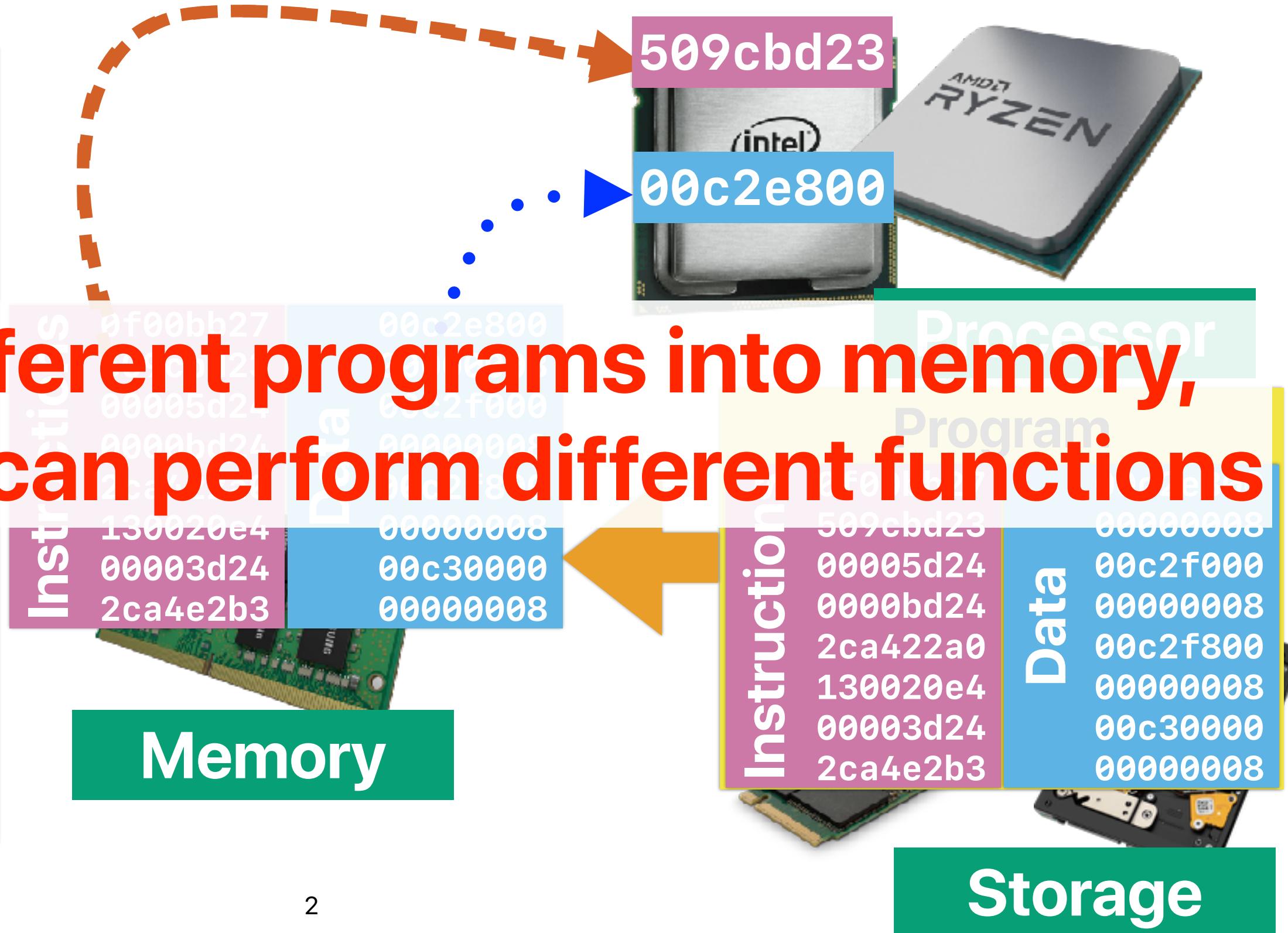
# **Performance (Preview Version)**

Hung-Wei Tseng

# Recap: von Neumann Architecture



By loading different programs into memory,  
your computer can perform different functions



# Recap: Demo

```
if(option)
    std::sort(data, data + arraySize);  $O(n \log_2 n)$ 
```

```
for (unsigned c = 0; c < arraySize*1000; ++c) {
    if (data[c%arraySize] >= INT_MAX/2)
        sum++;
}
```

$O(n)$

**if option is set to 1:**  $O(n \log_2 n)$

**otherwise, O(n):**  $O(n)$

# **Definition of “Performance”**



CANONICAL

ubuntu®

Enterprise ▾ Developer ▾ Community ▾ Download ▾

Downloads Overview Cloud IoT Raspberry Pi Server Desktop AMD-Xilinx Alternative downloads

## Download Ubuntu Desktop

The open-source desktop operating system that powers millions of PCs and laptops around the world. Find out more about Ubuntu's features and how we support developers and organisations below.

[Ubuntu Desktop homepage](#)

[Visit the Ubuntu Desktop blog >](#)



# Peer instruction

- Before the lecture — You need to complete the required **reading**
- During the lecture — I'll bring in activities to ENGAGE you in exploring your understanding of the material
  - Popup questions
  - Individual **thinking** — use polls in Zoom to express your opinion
  - Group **discussion**
    - Breakout rooms based on your residential colleges!
    - Use polls in Zoom to express your group's opinion
  - Whole-classroom **discussion** — we would like to hear from you

Read

Think

Discuss

**Let's start with "end-to-end latency"**  
— how long it takes to execute a  
program?

# CPU Performance Equation

$$Performance = \frac{1}{Execution\ Time}$$

$$Execution\ Time = \frac{Instructions}{Program} \times \frac{Cycles}{Instruction} \times \frac{Seconds}{Cycle}$$

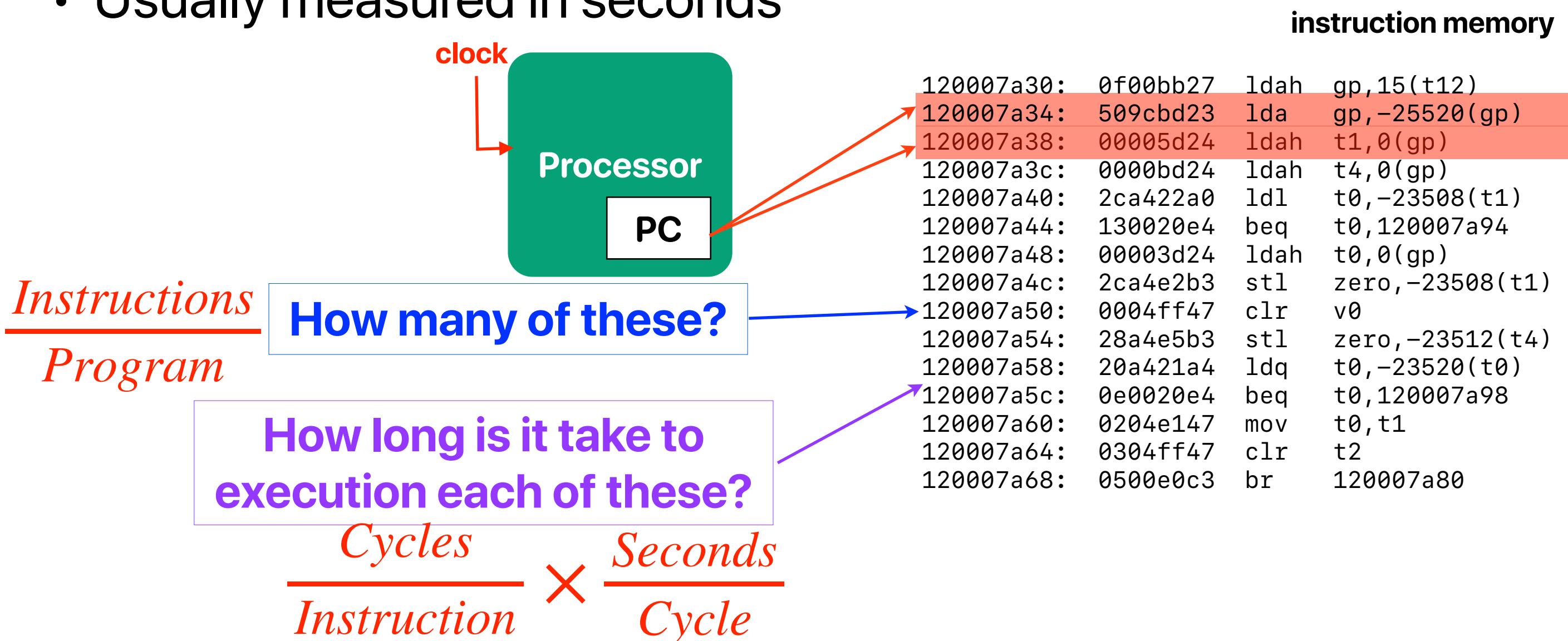
$$ET = IC \times CPI \times CT$$

$$1GHz = 10^9Hz = \frac{1}{10^9}sec\ per\ cycle = 1\ ns\ per\ cycle$$

$\frac{1}{Frequency(i.e.,\ clock\ rate)}$

# Execution Time

- The simplest kind of performance
- Shorter execution time means better performance
- Usually measured in seconds



# Performance Equation (X)

- Assume that we have an application composed with a total of **5000000000** instructions, in which **20%** of them are “Type-A” instructions with an average **CPI of 8** cycles, **20%** of them are “Type-B” instructions with an average **CPI of 4** cycles and **the rest** instructions are “Type-C” instructions with average **CPI of 1** cycle. If the processor runs at **3 GHz**, how long is the execution time?

A. 3.67 sec

B. 5 sec

C. 6.67 sec

D. 15 sec

E. 45 sec

$$ET = (5 \times 10^9) \times (20\% \times 8 + 20\% \times 4 + 60\% \times 1) \times \frac{1}{3 \times 10^{-9}} \text{ sec} = 5 \text{ average CPI}$$

$$ET = IC \times CPI \times CT$$

# Speedup

- The relative performance between two machines, X and Y. Y is  $n$  times faster than X

$$n = \frac{\text{Execution Time}_X}{\text{Execution Time}_Y}$$

- The speedup of Y over X

$$\text{Speedup} = \frac{\text{Execution Time}_X}{\text{Execution Time}_Y}$$

# Speedup of Y over X

- Consider the same program on the following two machines, X and Y. By how much Y is faster than X?

	Clock Rate	Instructions	Percentage of Type-A	CPI of Type-A	Percentage of Type-B	CPI of Type-B	Percentage of Type-C	CPI of Type-C
Machine X	3 GHz	5000000000	20%	8	20%	4	60%	1
Machine Y	5 GHz	5000000000	20%	13	20%	4	60%	1

A. 0.2       $ET_Y = (5 \times 10^9) \times (20\% \times 13 + 20\% \times 4 + 60\% \times 1) \times \frac{1}{5 \times 10^9} sec = 4$

B. 0.25      
$$\begin{aligned} Speedup &= \frac{Execution\ Time_X}{Execution\ Time_Y} \\ &= \frac{5}{4} = 1.25 \end{aligned}$$

D. 1.25

E. No changes

# **What Affects Each Factor in Performance Equation**

# Demo — programmer & performance

A

```
for(i = 0; i < ARRAY_SIZE; i++)  
{  
    for(j = 0; j < ARRAY_SIZE; j++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)  
{  
    for(i = 0; i < ARRAY_SIZE; i++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```

$O(n^2)$

Complexity

$O(n^2)$

Instruction Count?

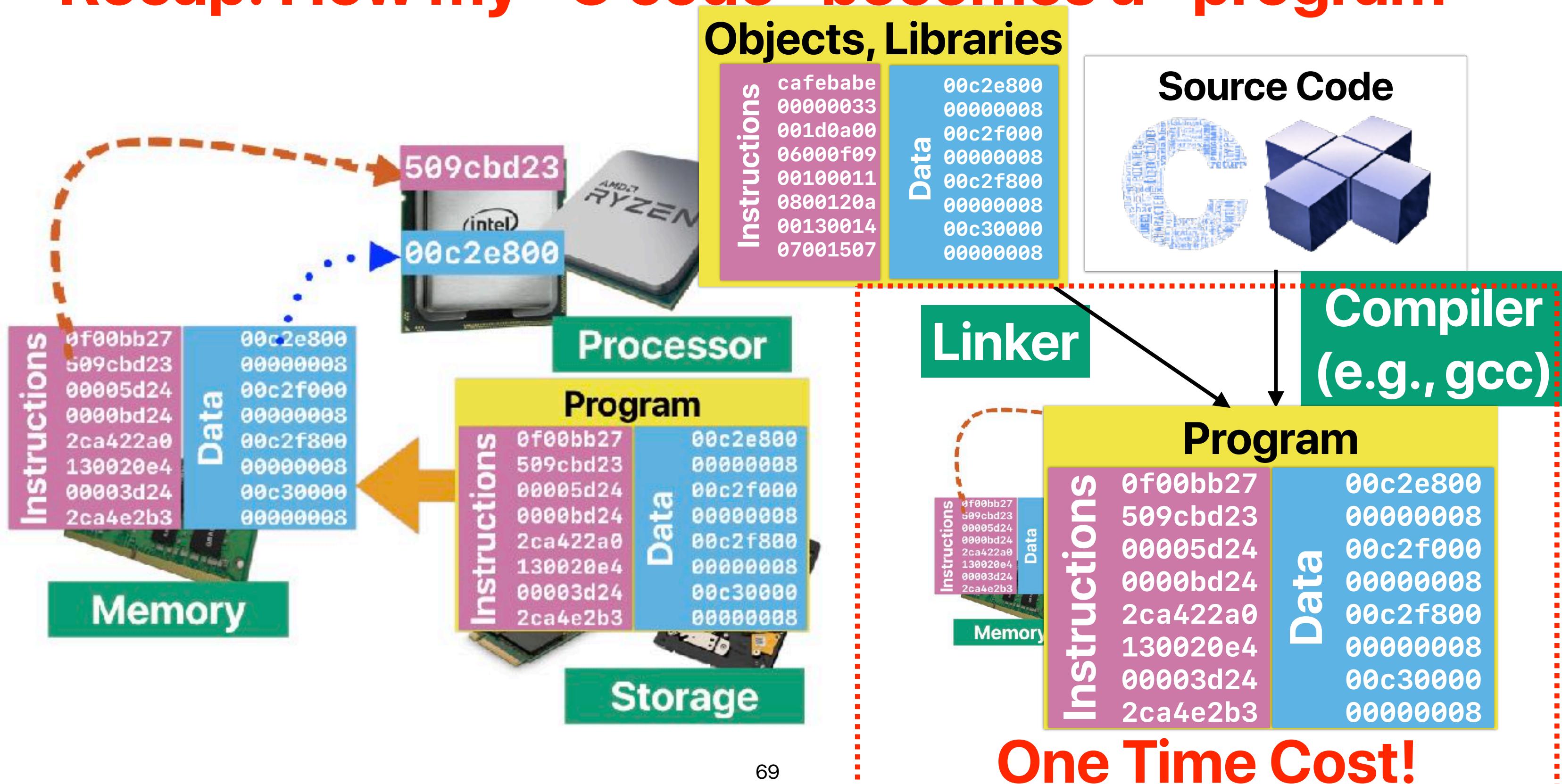
Clock Rate

CPI

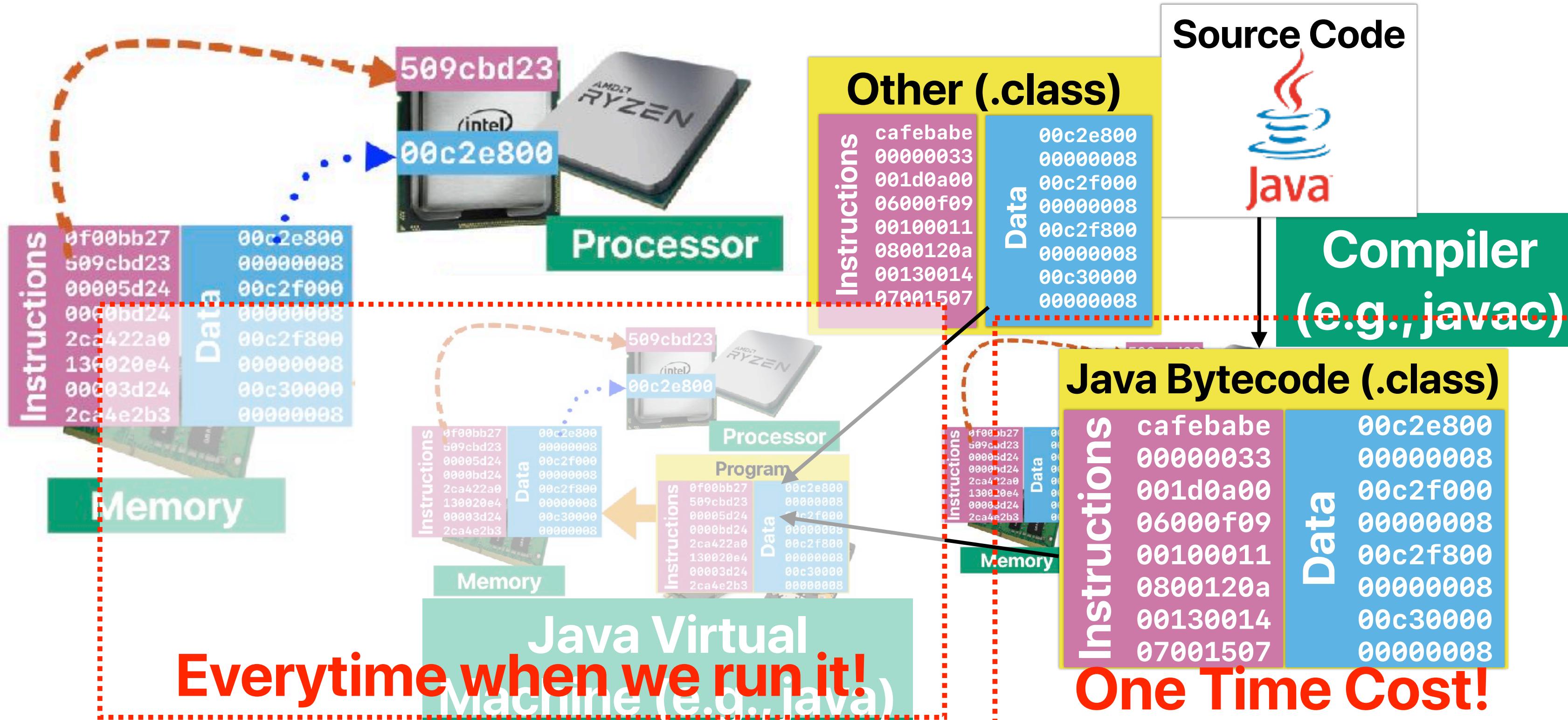
# Use “performance counters” to figure out!

- Modern processors provides performance counters
  - instruction counts
  - cache accesses/misses
  - branch instructions/mis-predictions
- How to get their values?
  - You may use “perf stat” in linux
  - You may use Instruments —> Time Profiler on a Mac
  - Intel’s vtune — only works on Windows w/ intel processors
  - You can also create your own functions to obtain counter values

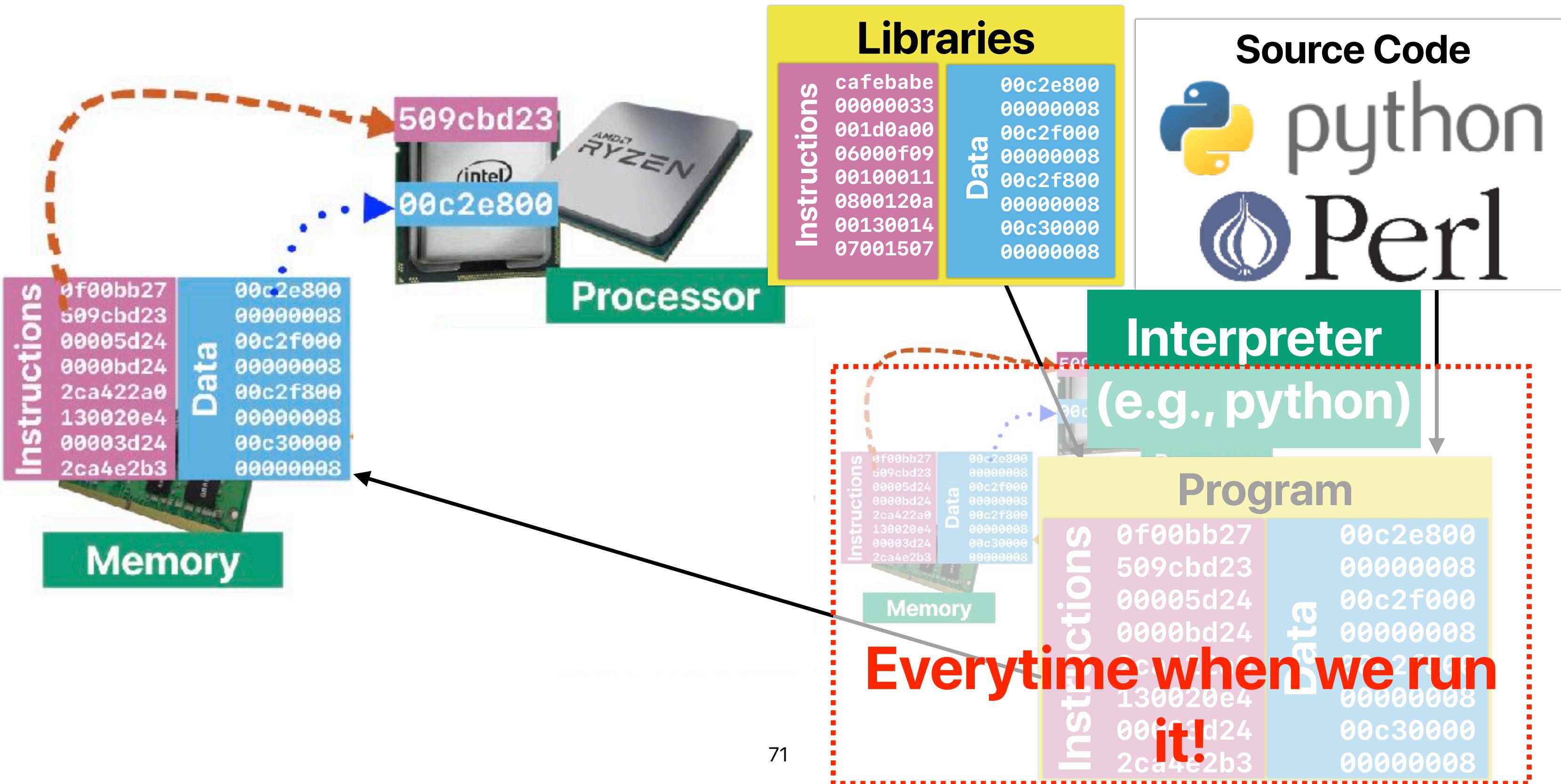
# Recap: How my “C code” becomes a “program”



# Recap: How my “Java code” becomes a “program”



# Recap: How my “Python code” becomes a “program”



# Revisited the demo with compiler optimizations!

- gcc has different optimization levels.
  - -O0 — no optimizations
  - -O3 — typically the best-performing optimization

A

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)
{
    for(i = 0; i < ARRAY_SIZE; i++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

# **Instruction Set Architecture (ISA) & Performance**

# Recap: ISA — the interface b/w processor/software

- Operations
  - Arithmetic/Logical, memory access, control-flow (e.g., branch, function calls)
  - Operands
    - Types of operands — register, constant, memory addresses
    - Sizes of operands — byte, 16-bit, 32-bit, 64-bit
- Memory space
  - The size of memory that programs can use
  - The addressing of each memory locations
  - The modes to represent those addresses

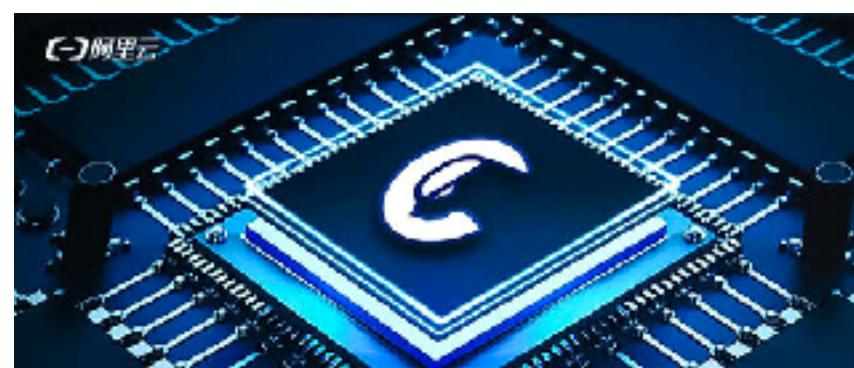
# Popular ISAs



x86

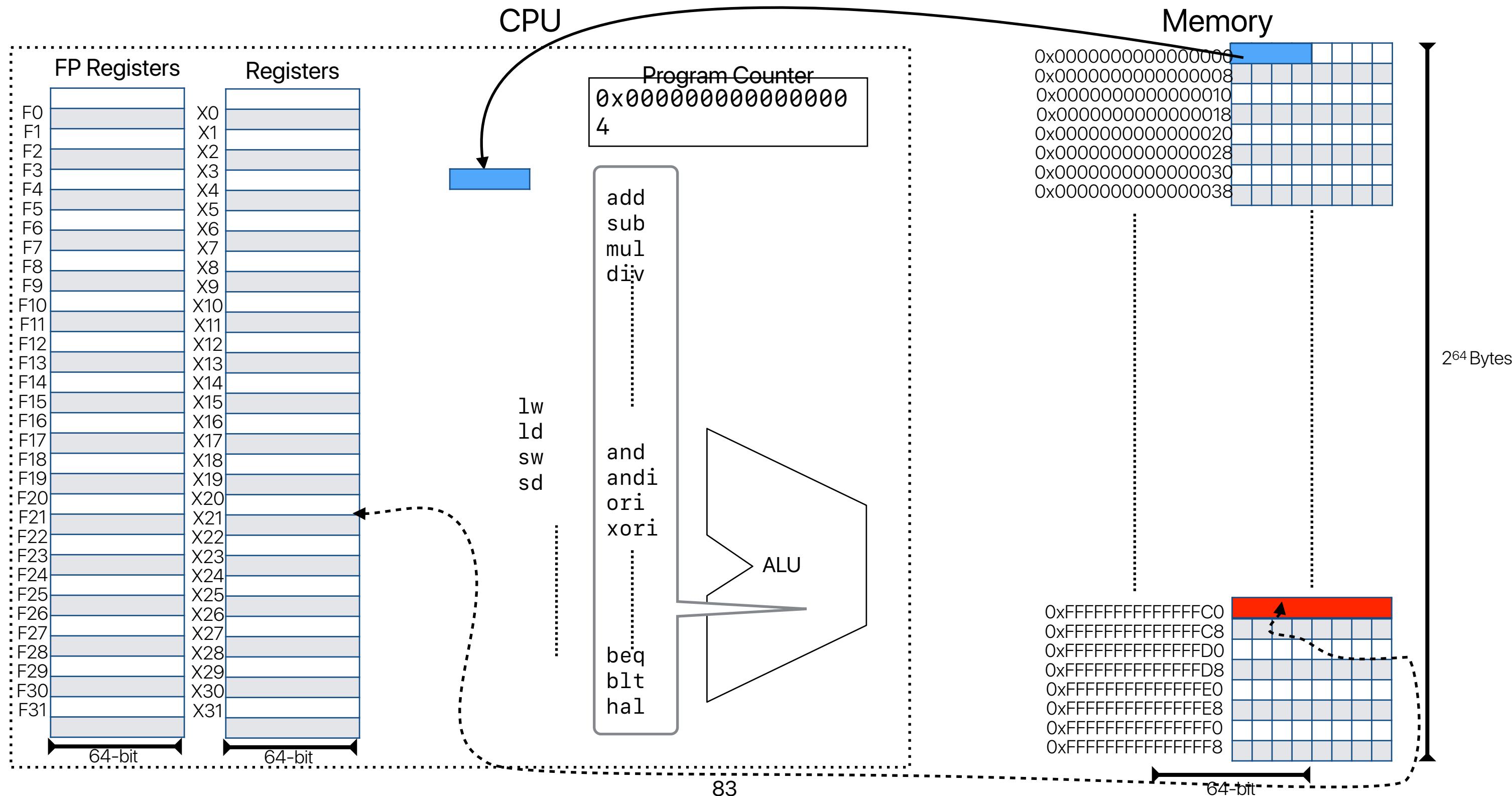


arm



RISC-V

# The abstracted “RISC-V” machine

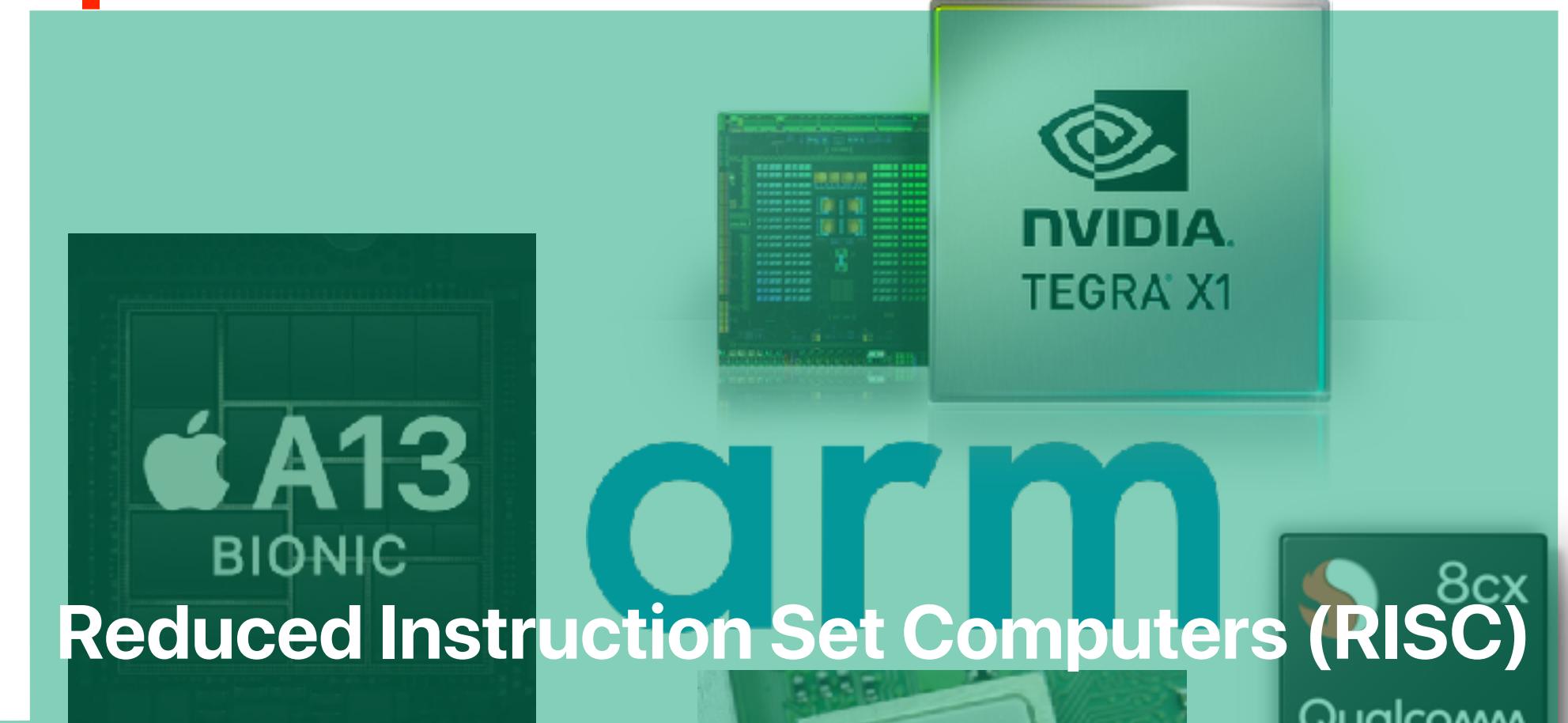


# Subset of RISC-V instructions

Category	Instruction	Usage	Meaning
Arithmetic	add	add x1, x2, x3	$x1 = x2 + x3$
	addi	addi x1, x2, 20	$x1 = x2 + 20$
	sub	sub x1, x2, x3	$x1 = x2 - x3$
Logical	and	and x1, x2, x3	$x1 = x2 \& x3$
	or	or x1, x2, x3	$x1 = x2   x3$
	andi	andi x1, x2, 20	$x1 = x2 \& 20$
	sll	sll x1, x2, 10	$x1 = x2 * 2^{10}$
	srl	srl x1, x2, 10	$x1 = x2 / 2^{10}$
Data Transfer	ld	ld x1, 8(x2)	$x1 = \text{mem}[x2+8]$
	sd	sd x1, 8(x2)	$\text{mem}[x2+8] = x1$
Branch	beq	beq x1, x2, 25	if( $x1 == x2$ ), PC = PC + 100
	bne	bne x1, x2, 25	if( $x1 != x2$ ), PC = PC + 100
Jump	jal	jal 25	\$ra = PC + 4, PC = 100
	jr	jr \$ra	PC = \$ra

The only type of instructions can access memory

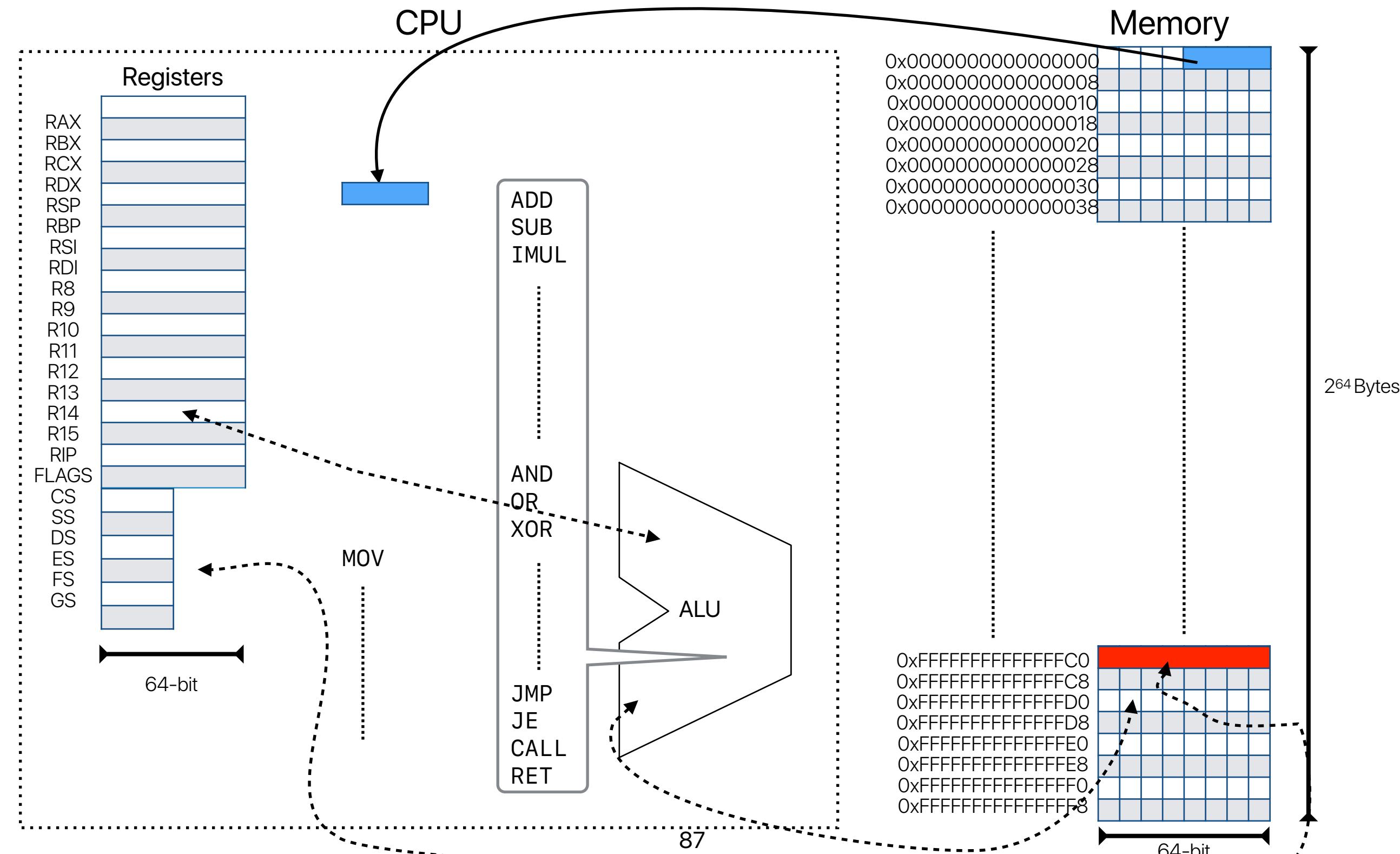
# Popular ISAs



# How many operations: CISC v.s. RISC

- CISC (Complex Instruction Set Computing)
  - Examples: x86, Motorola 68K
  - Provide **many powerful/complex** instructions
    - Many: more than 1503 instructions since 2016
    - Powerful/complex: an instruction can perform both ALU and memory operations
    - Each instruction takes more cycles to execute
- RISC (Reduced Instruction Set Computer)
  - Examples: ARMv8, RISC-V, MIPS (the first RISC instruction, invented by the authors of our textbook)
  - Each instruction only performs simple tasks
  - Easy to decode
  - Each instruction takes less cycles to execute

# The abstracted x86 machine



# RISC-V v.s. x86

	RISC-V	x86
ISA type	Reduced Instruction Set Computers (RISC)	Complex Instruction Set Computers (CISC)
instruction width	32 bits	1 ~ 17 bytes
code size	larger	smaller
registers	32	16
addressing modes	reg+offset	base+offset base+index scaled+index scaled+index+offset
hardware	simple	complex

# **How about complexity?**

# **How about “computational complexity”**

- Algorithm complexity provides a good estimate on the performance if —
  - Every instruction takes exactly the same amount of time
  - Every operation takes exactly the same amount of instructions

**These are unlikely to be true**

# Summary of CPU Performance Equation

$$\text{Performance} = \frac{1}{\text{Execution Time}}$$

$$\text{Execution Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Cycle}}$$

$$ET = IC \times CPI \times CT$$

- IC (Instruction Count)
  - ISA, Compiler, algorithm, programming language, **programmer**
- CPI (Cycles Per Instruction)
  - Machine Implementation, microarchitecture, compiler, application, algorithm, programming language, **programmer**
- Cycle Time (Seconds Per Cycle)
  - Process Technology, microarchitecture, **programmer**

# **Amdahl's Law — and It's Implication in the Multicore Era**

H&P Chapter 1.9

M. D. Hill and M. R. Marty. Amdahl's Law in the Multicore Era. In Computer, vol. 41, no. 7, pp. 33-38, July 2008.

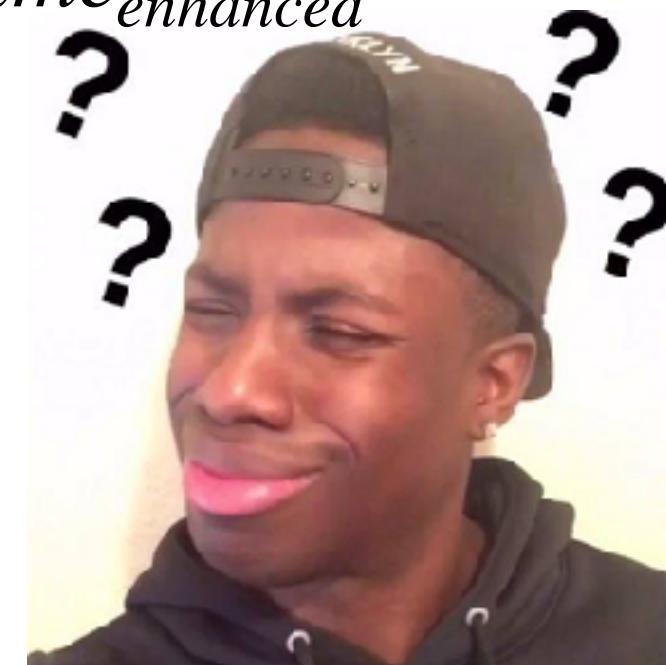
# Amdahl's Law



$$\text{Speedup}_{\text{enhanced}}(f, s) = \frac{1}{(1-f) + \frac{f}{s}}$$

- $f$  — The fraction of time in the original program  
 $s$  — The speedup we can achieve on  $f$

$$\text{Speedup}_{\text{enhanced}} = \frac{\text{Execution Time}_{\text{baseline}}}{\text{Execution Time}_{\text{enhanced}}}$$



# Amdahl's Law

$$Speedup_{enhanced}(f, s) = \frac{1}{(1 - f) + \frac{f}{s}}$$



$$Speedup_{enhanced} = \frac{Execution\ Time_{baseline}}{Execution\ Time_{enhanced}} = \frac{1}{(1 - f) + \frac{f}{s}}$$



# Amdahl's Law on Multiple Optimizations

- We can apply Amdahl's law for multiple optimizations
- These optimizations must be dis-joint!
  - If optimization #1 and optimization #2 are dis-joint:



$$Speedup_{enhanced}(f_{Opt1}, f_{Opt2}, s_{Opt1}, s_{Opt2}) = \frac{1}{(1 - f_{Opt1} - f_{Opt2}) + \frac{f_{Opt1}}{s_{Opt1}} + \frac{f_{Opt2}}{s_{Opt2}}}$$

- If optimization #1 and optimization #2 are not dis-joint:



$$Speedup_{enhanced}(f_{OnlyOpt1}, f_{OnlyOpt2}, f_{BothOpt1Opt2}, s_{OnlyOpt1}, s_{OnlyOpt2}, s_{BothOpt1Opt2}) = \frac{1}{(1 - f_{OnlyOpt1} - f_{OnlyOpt2} - f_{BothOpt1Opt2}) + \frac{f_{BothOpt1Opt2}}{s_{BothOpt1Opt2}} + \frac{f_{OnlyOpt1}}{s_{OnlyOpt1}} + \frac{f_{OnlyOpt2}}{s_{OnlyOpt2}}}$$

# Corollary #1 on Multiple Optimizations

- If we can pick just one thing to work on/optimize



$$Speedup_{max}(f_1, \infty) = \frac{1}{(1 - f_1)}$$

$$Speedup_{max}(f_2, \infty) = \frac{1}{(1 - f_2)}$$

$$Speedup_{max}(f_3, \infty) = \frac{1}{(1 - f_3)}$$

$$Speedup_{max}(f_4, \infty) = \frac{1}{(1 - f_4)}$$

The biggest  $f_x$  would lead to the largest  $Speedup_{max}$ !

## Corollary #2 — make the common case fast!

- When  $f$  is small, optimizations will have little effect.
- Common == **most time consuming** not necessarily the most frequent
- The uncommon case doesn't make much difference
- The common case can change based on inputs, compiler options, optimizations you've applied, etc.

# Identify the most time consuming part

- Compile your program with -pg flag
- Run the program
  - It will generate a gmon.out
  - `gprof your_program gmon.out > your_program.prof`
- It will give you the profiled result in `your_program.prof`

## Corollary #2.1 Don't hurt non-common part too much

- If the program spend 90% in A, 10% in B. Assume that an optimization can accelerate A by 9x, by hurts B by 10x...
- Assume the original execution time is T. The new execution time

$$ET_{new} = \frac{ET_{old} \times 90\%}{9} + ET_{old} \times 10\% \times 10$$

$$ET_{new} = 1.1 \times ET_{old}$$

$$Speedup = \frac{ET_{old}}{ET_{new}} = \frac{ET_{old}}{1.1 \times ET_{old}} = 0.91 \times \dots \text{slowdown!}$$

You may not use Amdahl's Law for this case as Amdahl's Law does NOT  
(1) consider overhead  
(2) bound to slowdown

# Corollary #3 — optimization has a moving target



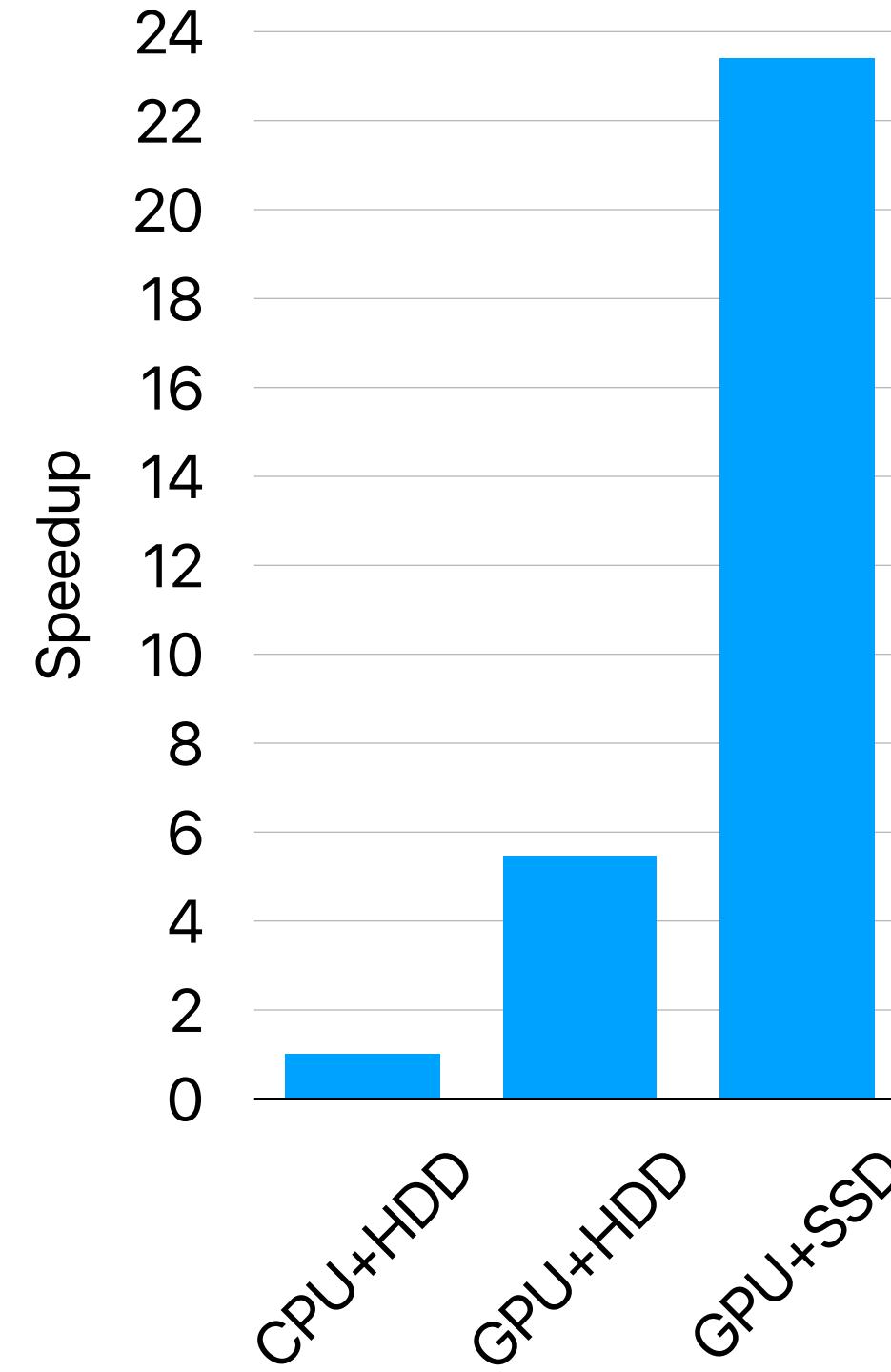
- With optimization, the common becomes uncommon.
- An uncommon case will (hopefully) become the new common case.
- Now you have a new target for optimization — You have to revisit “Amdahl’s Law” every time you applied some optimization

Something else (e.g.,  
data movement)  
matters more now

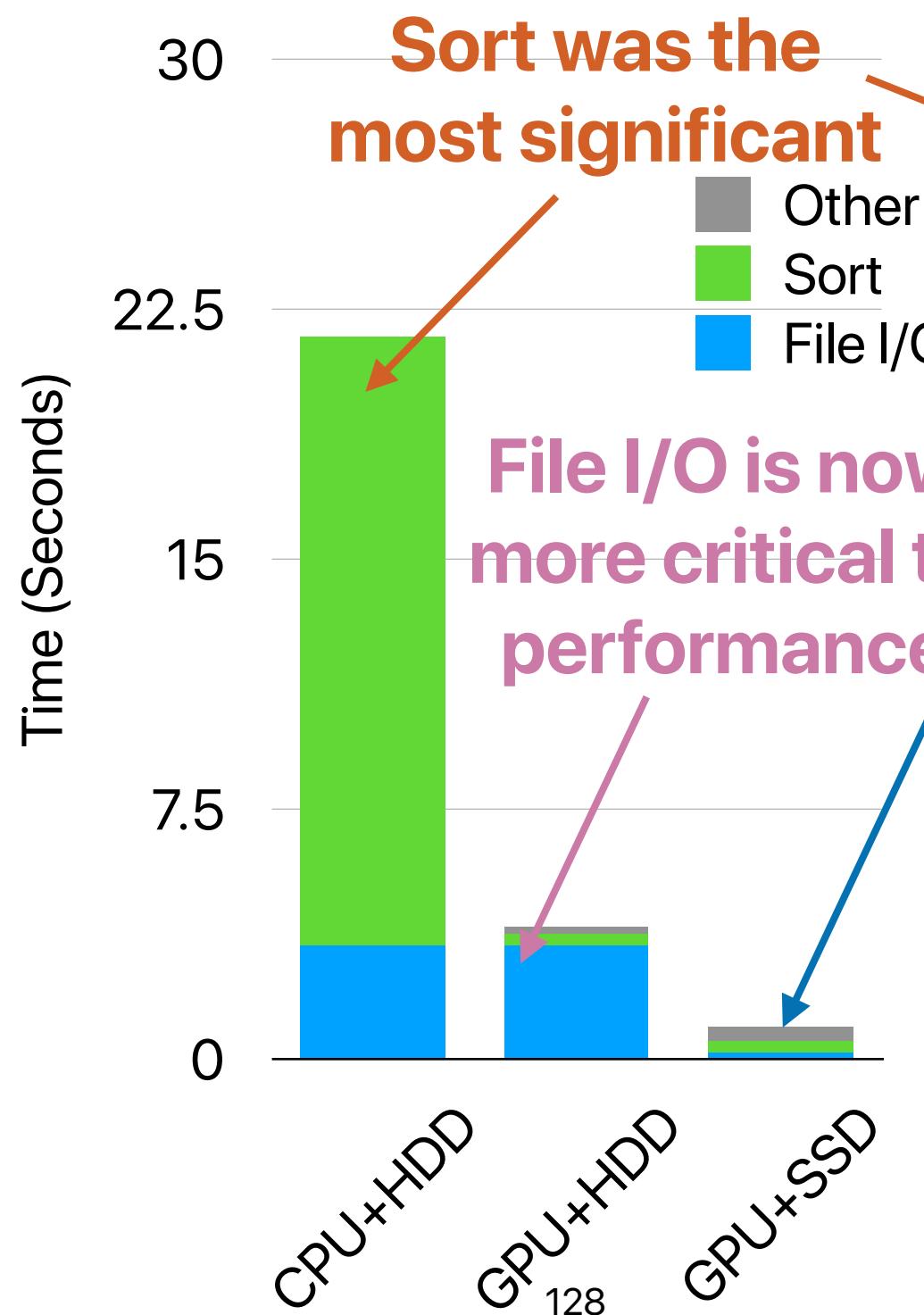
# Demo — sort

Something else (e.g., data movement) matters more

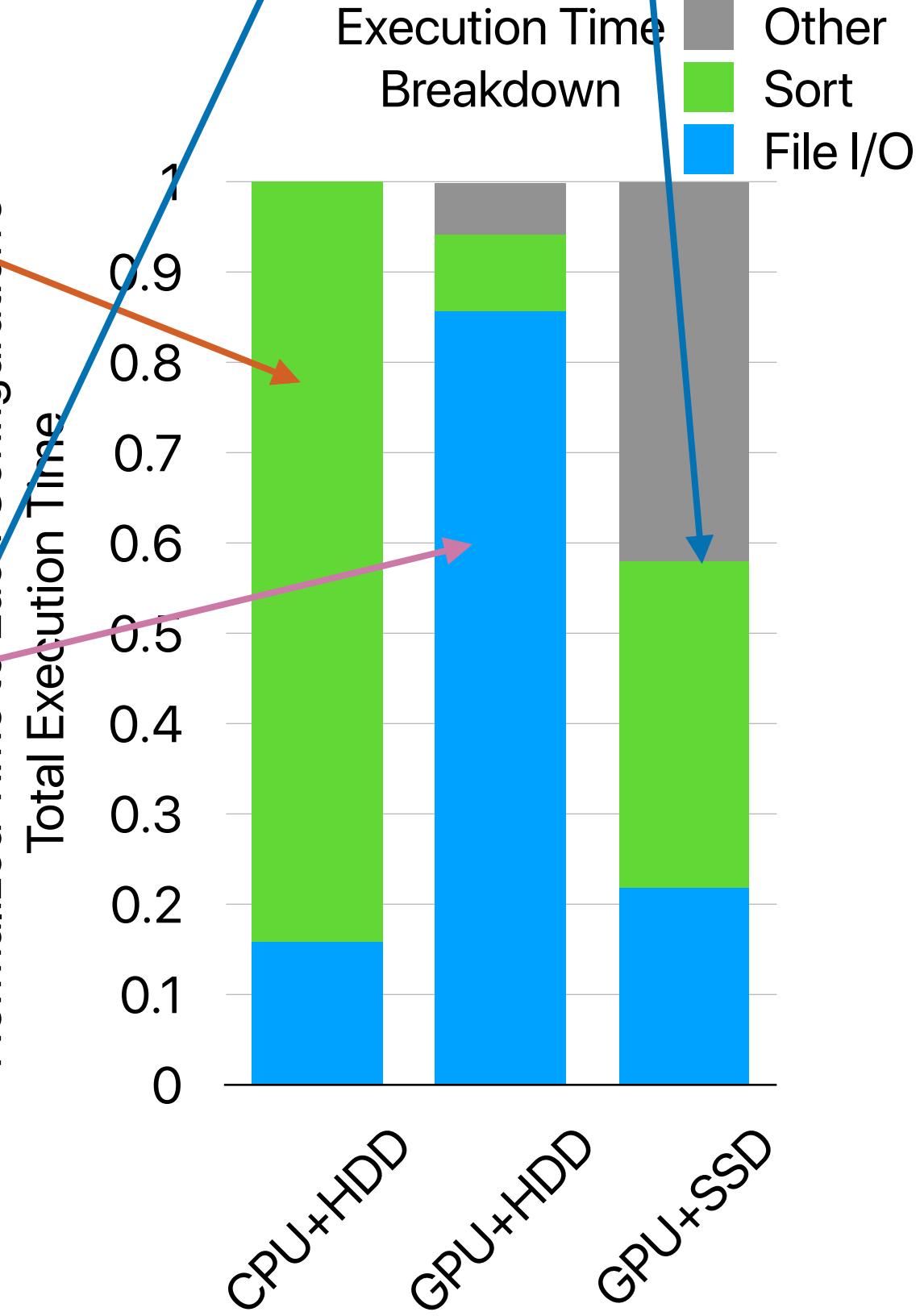
Speedup



Cumulative Execution Time



Normalized Time to Each Configuration's Total Execution Time



# Amdahl's Law on Multicore Architectures

- Symmetric multicore processor with  $n$  cores (if we assume the processor performance scales perfectly)

$$\text{Speedup}_{\text{parallel}}(f_{\text{parallelizable}}, n) = \frac{1}{(1 - f_{\text{parallelizable}}) + \frac{f_{\text{parallelizable}}}{n}}$$

# Demo — merge sort v.s. bitonic sort on GPUs

## Merge Sort

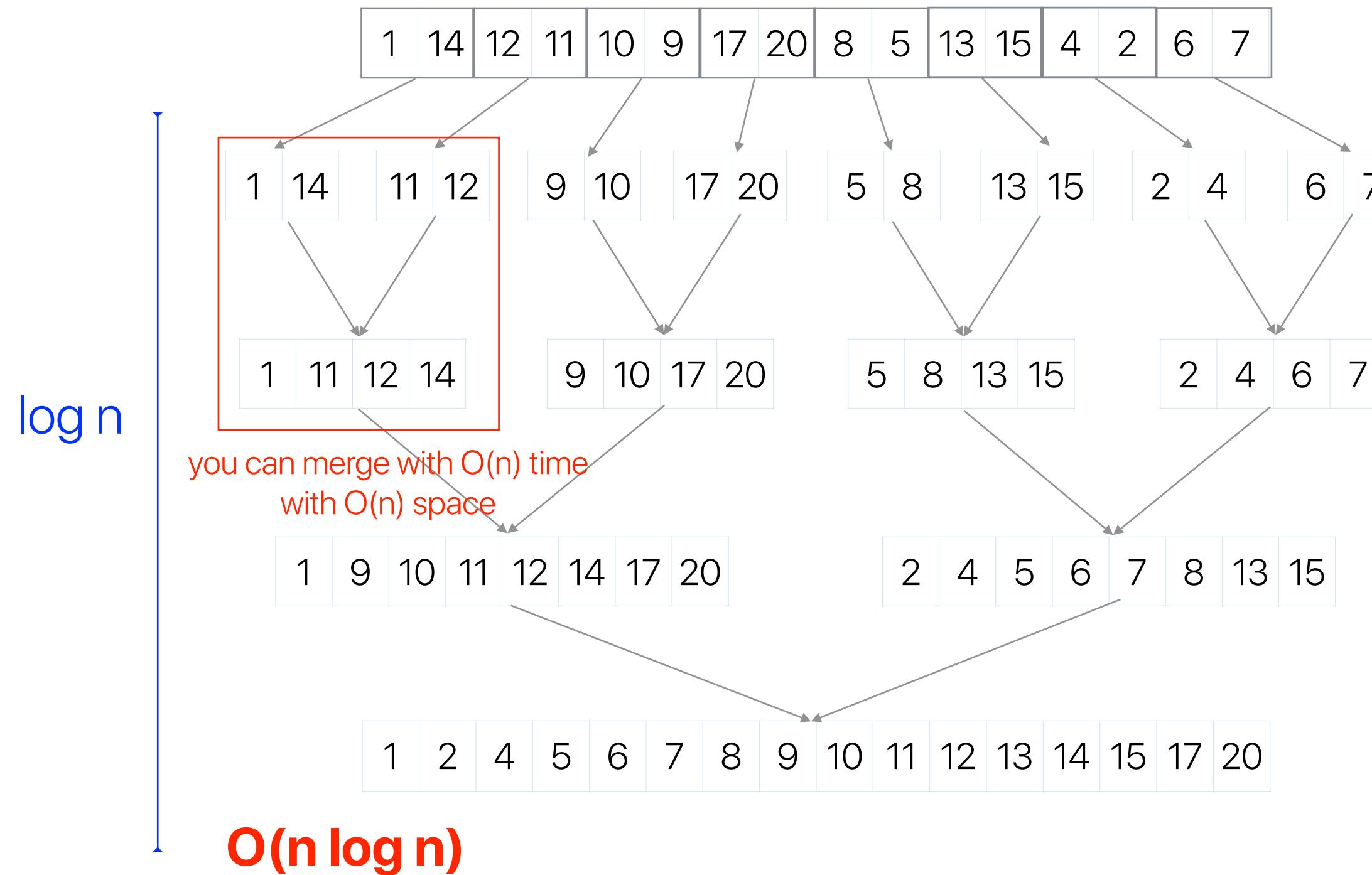
$$O(n \log_2 n)$$

## Bitonic Sort

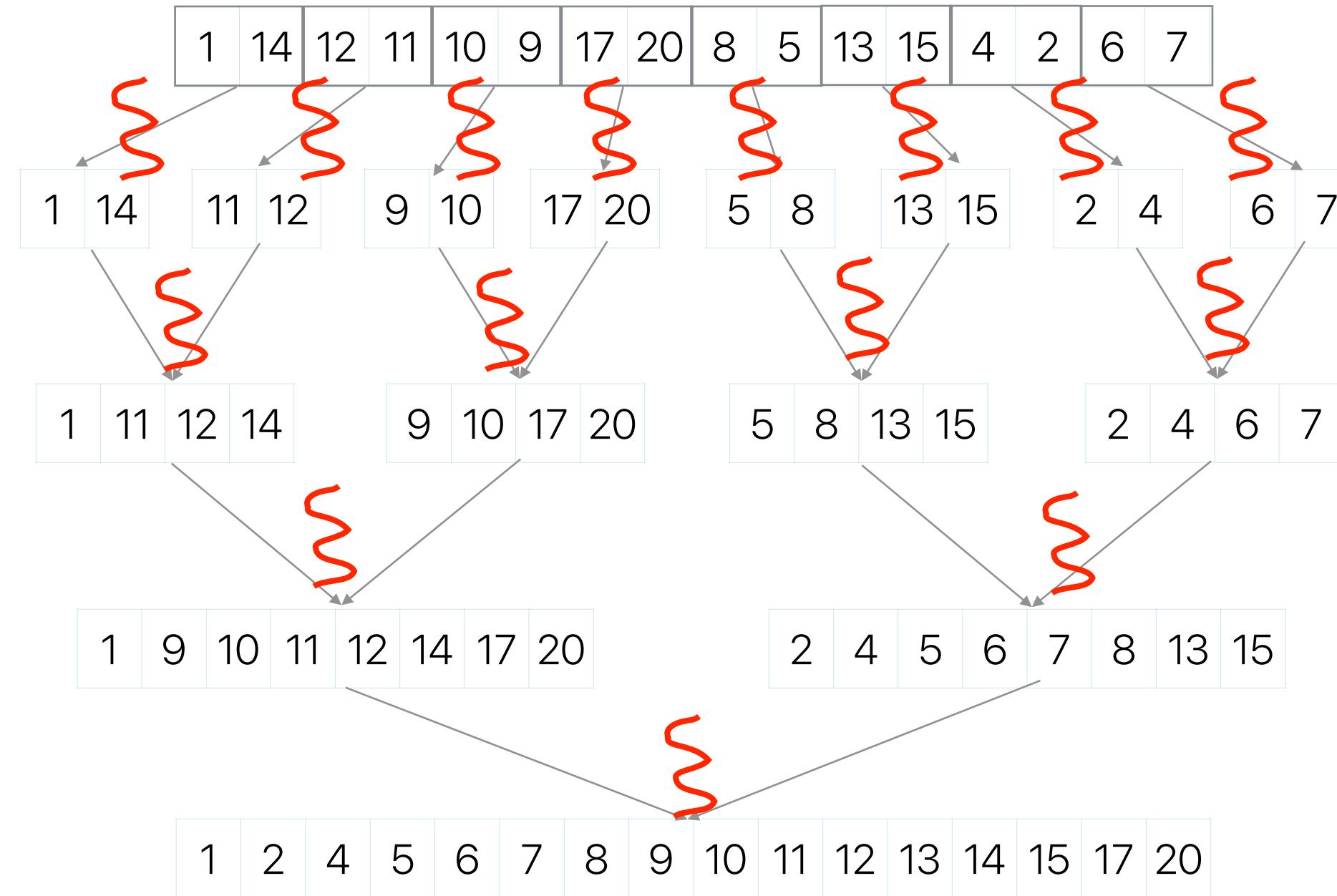
$$O(n \log_2^2 n)$$

```
void BitonicSort() {  
    int i,j,k;  
  
    for (k=2; k<=N; k=2*k) {  
        for (j=k>>1; j>0; j=j>>1) {  
            for (i=0; i<N; i++) {  
                int ij=i^j;  
                if ((ij)>i) {  
                    if ((i&k)==0 && a[i] > a[ij])  
                        exchange(i,ij);  
                    if ((i&k)!=0 && a[i] < a[ij])  
                        exchange(i,ij);  
                }  
            }  
        }  
    }  
}
```

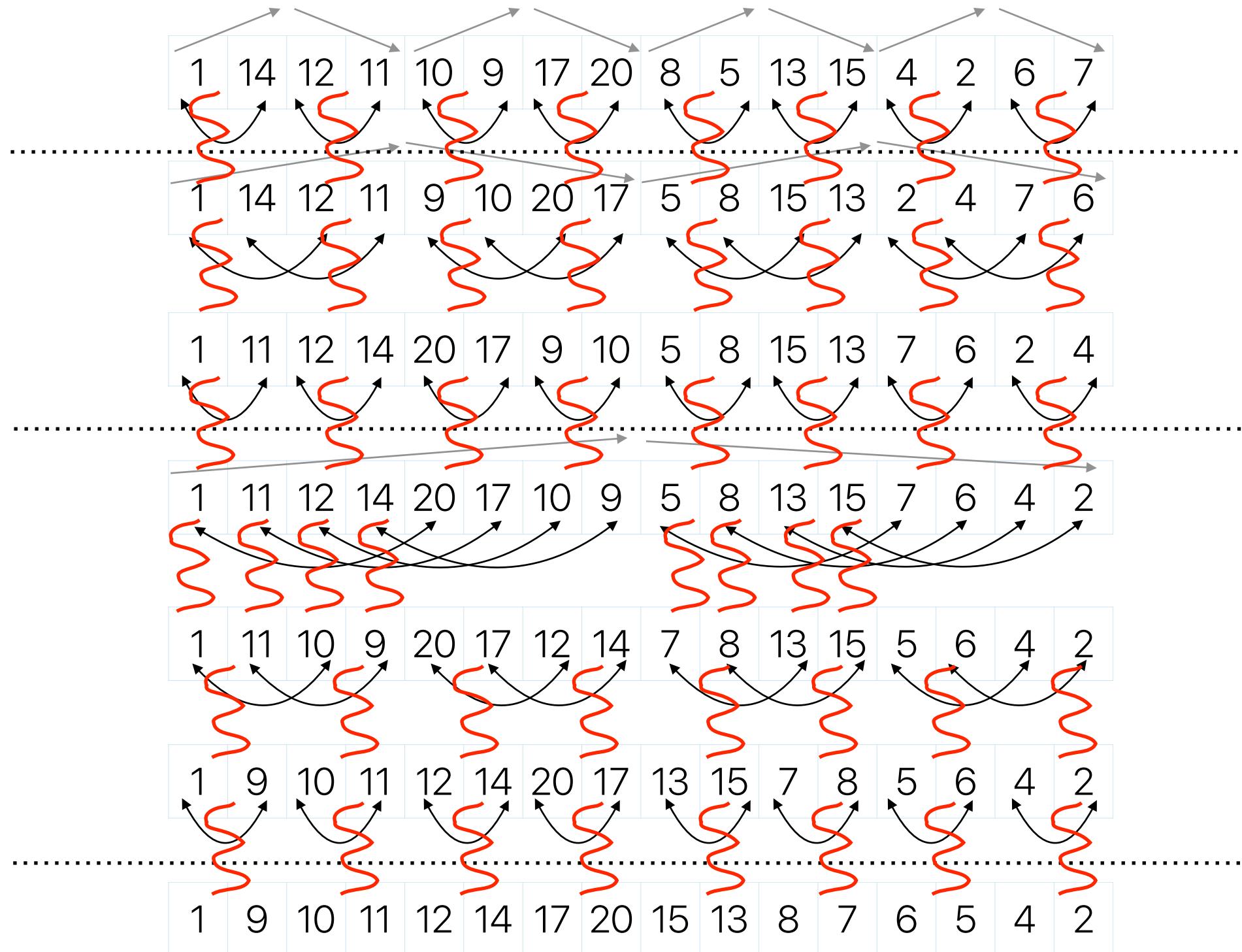
# Merge sort



# Parallel merge sort

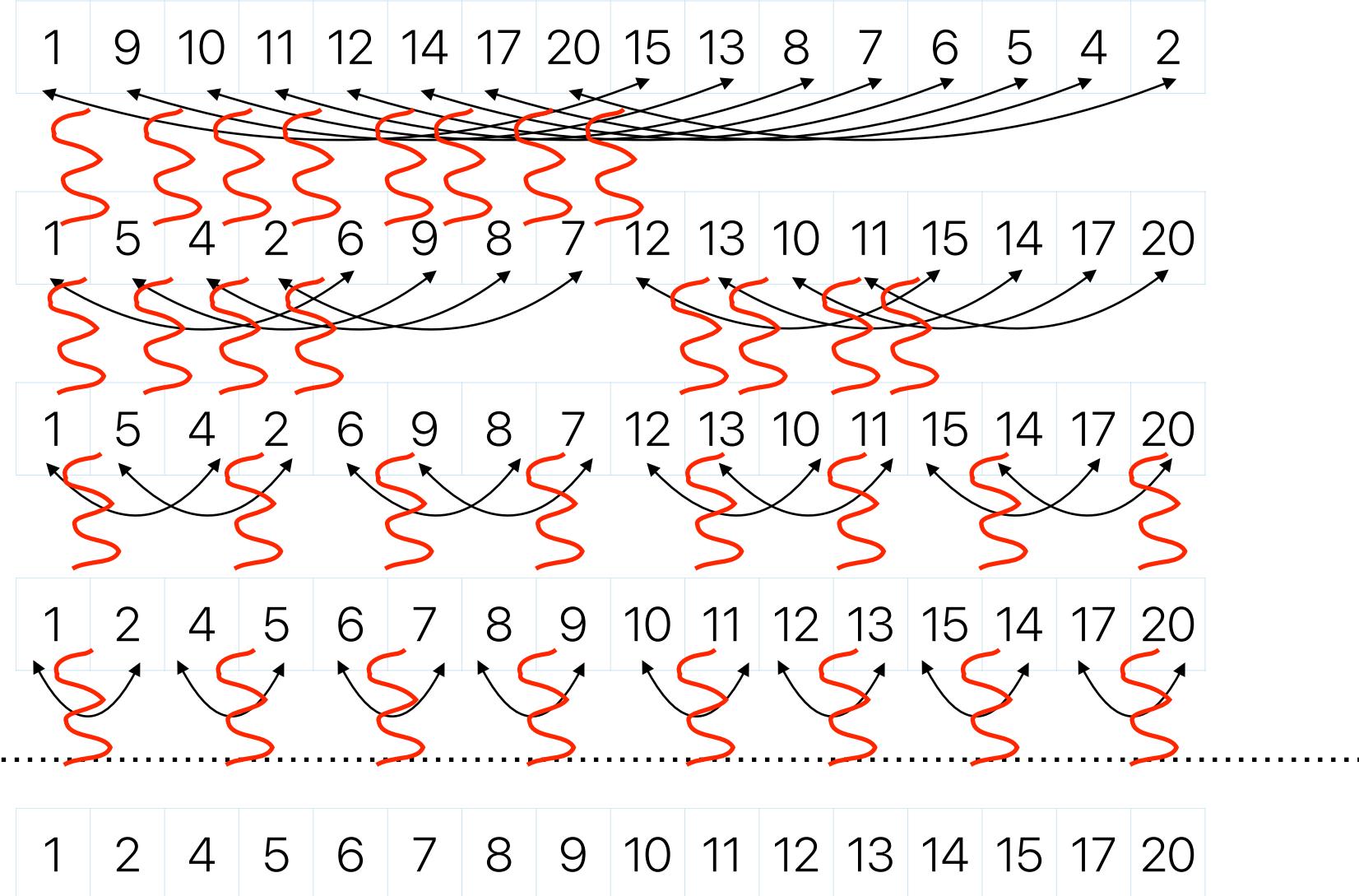


# Bitonic sort



```
void BitonicSort() {  
    int i, j, k;  
  
    for (k=2; k<=N; k=2*k) {  
        for (j=k>>1; j>0; j=j>>1) {  
            for (i=0; i<N; i++) {  
                int ij=i^j;  
                if ((ij)>i) {  
                    if ((i&k)==0 && a[i] > a[ij])  
                        exchange(i,ij);  
                    if ((i&k)!=0 && a[i] < a[ij])  
                        exchange(i,ij);  
                }  
            }  
        }  
    }  
}
```

# Bitonic sort (cont.)



```
void BitonicSort() {  
  
    int i, j, k;  
  
    for (k=2; k<=N; k=2*k) {  
        for (j=k>>1; j>0; j=j>>1) {  
            for (i=0; i<N; i++) {  
                int ij=i^j;  
                if ((ij)>i) {  
                    if ((i&k)==0 && a[i] > a[ij])  
                        exchange(i,ij);  
                    if ((i&k)!=0 && a[i] < a[ij])  
                        exchange(i,ij);  
                }  
            }  
        }  
    }  
}
```

**benefits — in-place merge (no additional space is necessary), very stable comparison patterns**

**$O(n \log^2 n)$  — hard to beat  $n(\log n)$  if you can't parallelize this a lot!**

## Corollary #4

$$\text{Speedup}_{\text{parallel}}(f_{\text{parallelizable}}, \infty) = \frac{1}{(1 - f_{\text{parallelizable}}) + \frac{f_{\text{parallelizable}}}{\infty}}$$

$$\text{Speedup}_{\text{parallel}}(f_{\text{parallelizable}}, \infty) = \frac{1}{(1 - f_{\text{parallelizable}})}$$

- If we can build a processor with unlimited parallelism
  - The complexity doesn't matter as long as the algorithm can utilize all parallelism
  - That's why bitonic sort or MapReduce works!
- **The future trend of software/application design is seeking for more parallelism rather than lower the computational complexity**

**Is it the end of computational  
complexity?**

## Corollary #5

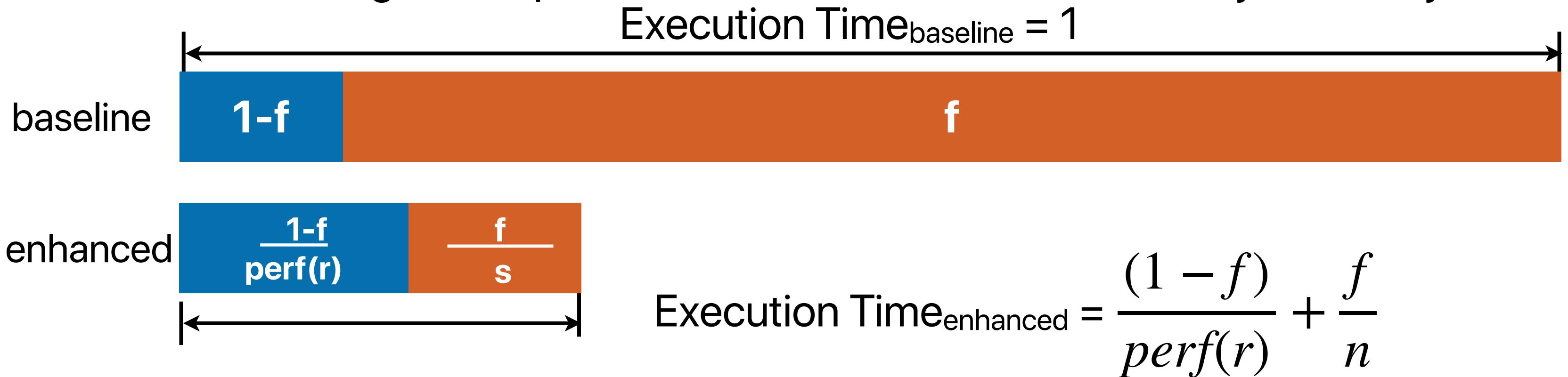
$$\text{Speedup}_{\text{parallel}}(f_{\text{parallelizable}}, \infty) = \frac{1}{(1 - f_{\text{parallelizable}}) + \frac{f_{\text{parallelizable}}}{\infty}}$$

$$\text{Speedup}_{\text{parallel}}(f_{\text{parallelizable}}, \infty) = \frac{1}{(1 - f_{\text{parallelizable}})}$$

- Single-core performance still matters
  - It will eventually dominate the performance
  - If we cannot improve single-core performance further, finding more “parallelizable” parts is more important
  - Algorithm complexity still gives some “insights” regarding the growth of execution time in the same algorithm, though still not accurate

# However, parallelism is not “tax-free”

- Synchronization
- Preparing data
- Addition function calls
- Data exchange if the parallel hardware has its own memory hierarchy



# Lessons learned from Amdahl's Law

$$Speedup_{enhanced}(f, s) = \frac{1}{(1 - f) + \frac{f}{s}}$$

- Corollary #1: Maximum speedup
- Corollary #2: Make the common case fast
  - Common case changes all the time
- Corollary #3: Optimization is a moving target
- Corollary #4: Exploiting more parallelism from a program is the key to performance gain in modern architectures
- Corollary #5: Single-core performance still matters

$$\begin{aligned} Speedup_{max}(f, \infty) &= \frac{1}{(1 - f)} \\ Speedup_{max}(f_1, \infty) &= \frac{1}{(1 - f_1)} \\ Speedup_{max}(f_2, \infty) &= \frac{1}{(1 - f_2)} \\ Speedup_{max}(f_3, \infty) &= \frac{1}{(1 - f_3)} \\ Speedup_{max}(f_4, \infty) &= \frac{1}{(1 - f_4)} \end{aligned}$$

$$\begin{aligned} Speedup_{parallel}(f_{parallelizable}, \infty) &= \frac{1}{(1 - f_{parallelizable})} \\ Speedup_{parallel}(f_{parallelizable}, \infty) &= \frac{1}{(1 - f_{parallelizable})} \end{aligned}$$

**Choose the right metric — Latency  
v.s. Throughput/Bandwidth**

# Latency v.s. Bandwidth/Throughput

- Latency — the amount of time to finish an operation
  - End-to-end execution time of “something”
  - Access time
  - Response time
- Throughput — the amount of work can be done within a given period of time (typically “something” per “timeframe” or the other way around)
  - Bandwidth (MB/Sec, GB/Sec, Mbps, Gbps)
  - IOPs (I/O operations per second)
  - FLOPs (Floating-point operations per second)
  - IPS (Inferences per second)

# RAID — Improving throughput

## MORE SPECS

### Model Code (Capacity)

### General

### Storage



DIMENSION (WxHxD)  
100 X 69.85 X 6.8 (mm)

### Performance

### Power

### Environment

### Warranty

### Performance<sup>2)</sup>

TRIM SUPPORT  
Yes

ENCRYPTION SUPPORT  
AES 256-bit Encryption (Class 0) TCG/Opc  
IEEE1667 (Encrypted drive)

SEQUENTIAL READ  
Up to 550 MB/s

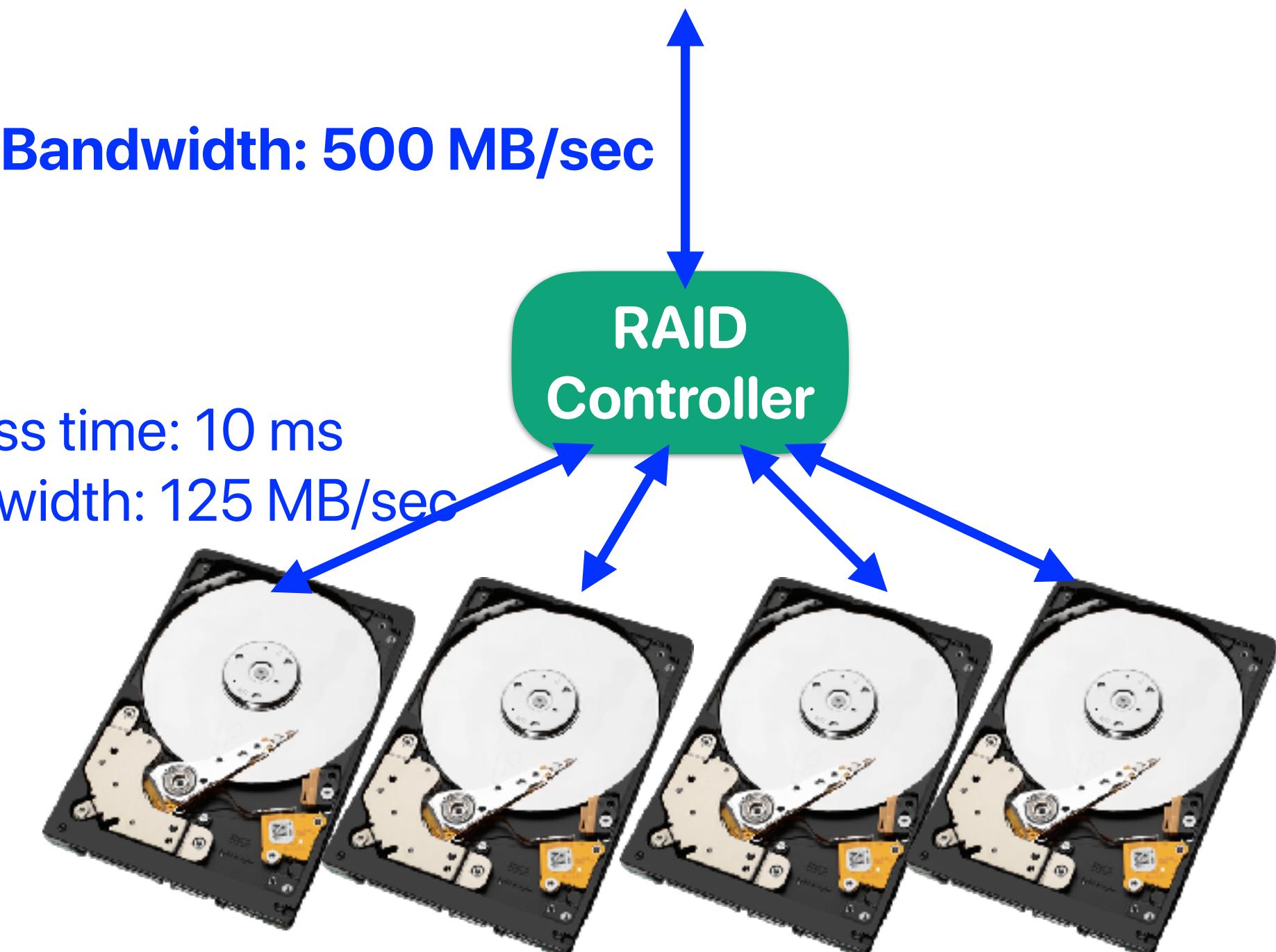
RANDOM WRITE (4KB, QD32)  
Up to 89,000 IOPS

AVERAGE POWER CONSUMPTION  
(SYSTEM LEVEL)<sup>3)</sup>

1,000 GB: Average 2.2 W Maximum 4.0 W  
2,000 GB: Average 3.1 W Maximum 4.2 W  
4,000 GB: Average 3.1 W Maximum 5.4 W  
(Burst mode)

**Aggregated Bandwidth: 500 MB/sec**

Access time: 10 ms  
Bandwidth: 125 MB/sec



# What have we learned?

- Bandwidth does not necessary reflect the performance
- GPUs have better “throughput”, but end2end latency is worse than CPUs if your samples are small

# Latency/Delay v.s. Throughput

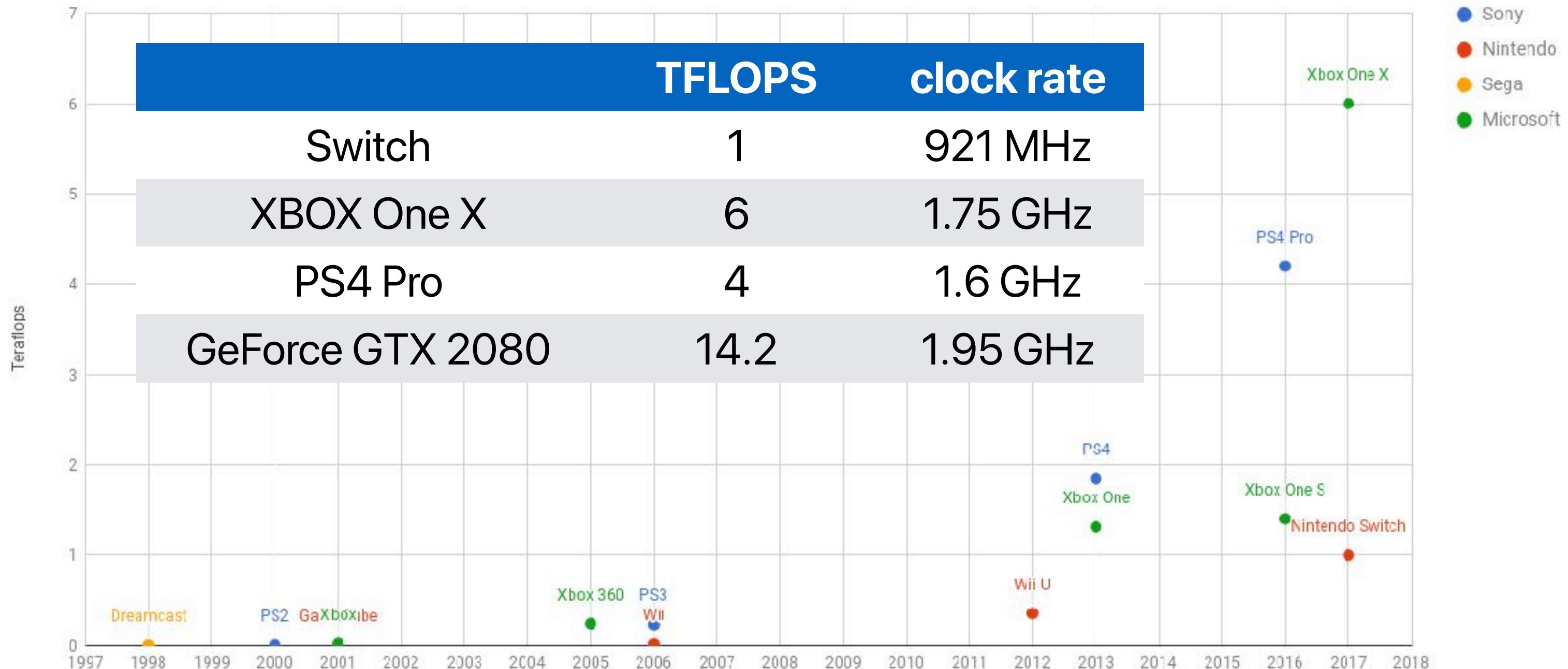
	Toyota Prius	100 Gb Network
bandwidth	290GB/sec	100 Gb/s or 12.5GB/sec
total latency	3.5 hours	2 Peta-byte over 167772 seconds = 1.94 Days
latency in getting the first movie	You see nothing in the first 3.5 hours	100GB/100Gb = 8 secs! You can start watching the first movie in 8 secs!

# “Fair” Comparisons

Andrew Davison. Twelve Ways to Fool the Masses When Giving Performance Results on Parallel Computers. In Humour the Computer, MITP, 1995  
V. Sze, Y. -H. Chen, T. -J. Yang and J. S. Emer. How to Evaluate Deep Neural Network Processors: TOPS/W (Alone) Considered Harmful. In IEEE Solid-State Circuits Magazine, vol. 12, no. 3, pp. 28-41, Summer 2020.

# TFLOPS (Tera FLoating-point Operations Per Second)

Console Teraflops



## Is TFLOPS (Tera FLoating-point Operations Per Second) a good metric?

$$\begin{aligned}TFLOPS &= \frac{\# \text{ of floating point instructions} \times 10^{-12}}{\text{Execution Time}} \\&= \frac{IC \times \% \text{ of floating point instructions} \times 10^{-12}}{IC \times CPI \times CT} \\&= \frac{\% \text{ of floating point instructions} \times 10^{-12}}{CPI \times CT}\end{aligned}$$

**IC is gone!**

- Cannot compare different ISA/compiler
  - What if the compiler can generate code with fewer instructions?
  - What if new architecture has more IC but also lower CPI?
- Does not make sense if the application is not floating point intensive

# TFLOPS (Tera FLoating-point Operations Per Second)

- Cannot compare different ISA/compiler
  - What if the compiler can generate code with fewer instructions?
  - What if new architecture has more IC but also lower CPI?
- Does not make sense if the application is not floating point intensive

	TFLOPS	clock rate
Switch	1	921 MHz
XBOX One X	6	1.75 GHz
PS4 Pro	4	1.6 GHz
GeForce GTX 2080	14.2	1.95 GHz



Artificial Intelligence Computing Leadership from NVIDIA

## CLOUD &amp; DATA CENTER

PRODUCTS ▾

SOLUTIONS ▾

APPS ▾

FOR DEVELOPERS

TECHNOLOGIES ▾

Tesla V100

AI TRAINING

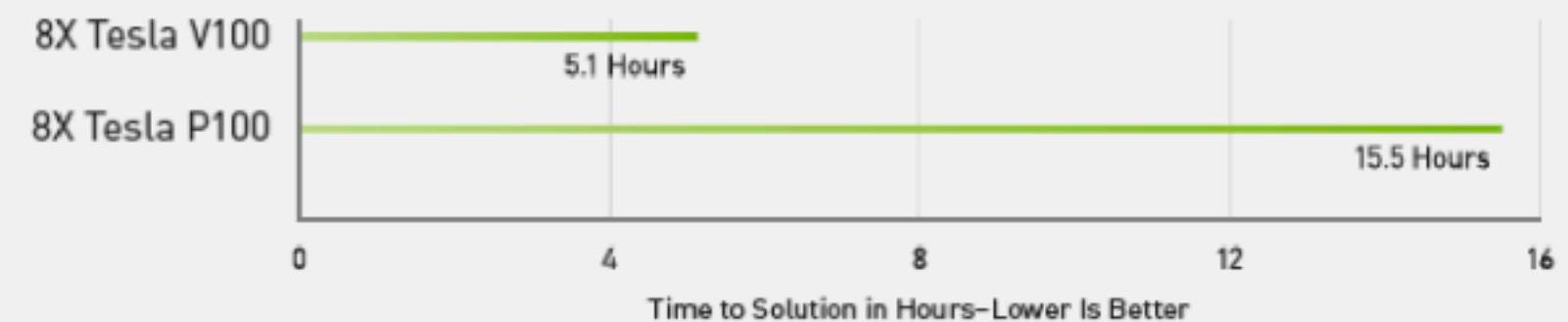
AI INFERENCE

HPC

DATA CENTER GPUs

SPECIFICATIONS

## Deep Learning Training in Less Than a Workday



Server Config: Dual Xeon E5-2699 v4 2.6 GHz | 8X NVIDIA® Tesla® P100 or V100 | ResNet-50 Training on MXNet for 90 Epochs with 1.28M ImageNet Dataset.

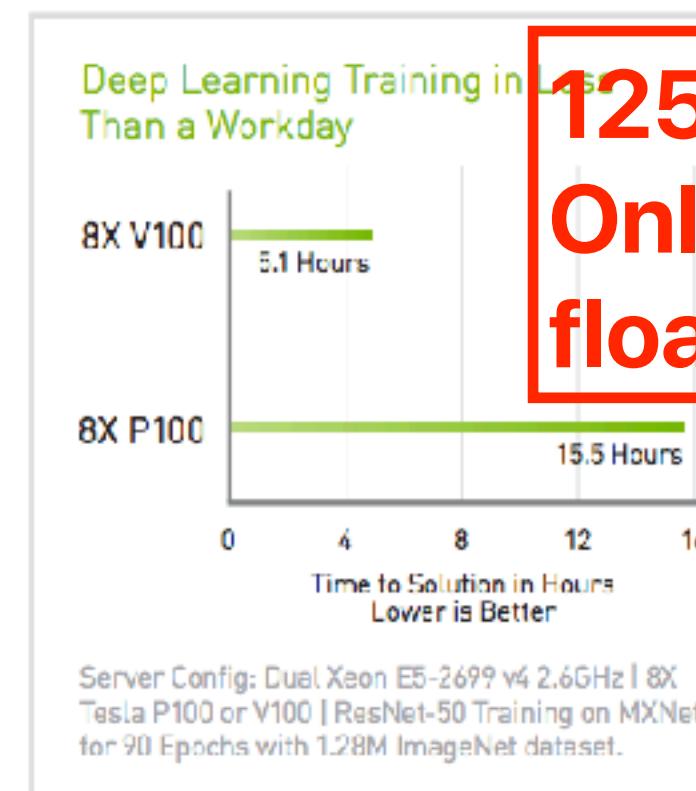
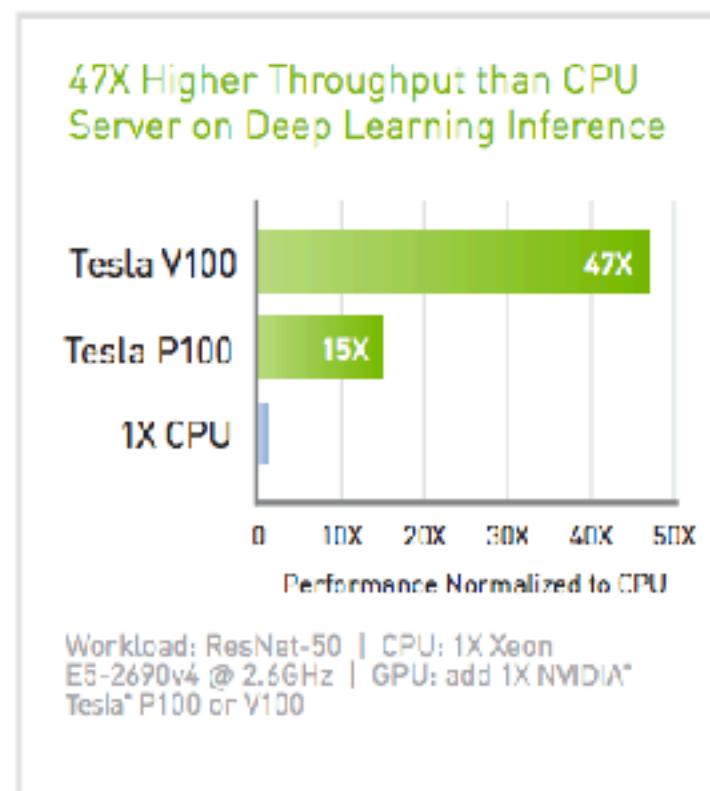
## AI TRAINING

From recognizing speech to training virtual personal assistants and teaching autonomous cars to drive, data scientists are taking on increasingly complex challenges with AI. Solving these kinds of problems requires training deep learning models that are exponentially growing in complexity, in a practical amount of time.

With 640 **Tensor Cores**, Tesla V100 is the world's first GPU to break the 100 teraFLOPS (TFLOPS) barrier of deep learning performance. The next generation of **NVIDIA NVLink™** connects multiple V100 GPUs at up to 300 GB/s to create the world's most powerful computing servers. AI models that would consume weeks of computing resources on previous systems can now be trained in a few days. With this dramatic reduction in training time, a whole new world of problems will now be solvable with AI.

# The Most Advanced Data Center GPU Ever Built.

NVIDIA® Tesla® V100 is the world's most advanced data center GPU ever built to accelerate AI, HPC, and graphics. Powered by NVIDIA Volta, the latest GPU architecture, Tesla V100 offers the performance of up to 100 CPUs in a single GPU—enabling data scientists, researchers, and engineers to tackle challenges that were once thought impossible.



**125 TFLOPS  
Only @ 16-bit floating point**

## SPECIFICATIONS

<b>Tesla V100 PCIe</b>	<b>Tesla V100 SXM2</b>
GPU Architecture	<b>NVIDIA Volta</b>
NVIDIA Tensor Cores	<b>640</b>
NVIDIA CUDA® Cores	<b>5,120</b>

Double-Precision Performance	<b>7 TFLOPS</b>	<b>7.8 TFLOPS</b>
Single-Precision Performance	<b>14 TFLOPS</b>	<b>15.7 TFLOPS</b>
Tensor Performance	<b>112 TFLOPS</b>	<b>125 TFLOPS</b>
GPU Memory	<b>32GB /16GB HBM2</b>	
Memory Bandwidth	<b>900GB/sec</b>	
ECC	<b>Yes</b>	
Interconnect Bandwidth	<b>32GB/sec</b>	<b>300GB/sec</b>
System Interface	<b>PCIe Gen3</b>	<b>NVIDIA NVLink</b>
Form Factor	<b>PCIe Full Height/Length</b>	<b>SXM2</b>
Max Power	<b>375W</b>	<b>300W</b>

1 GPU Node Replaces Up To 54 CPU Nodes

Node Replacement: HPC Mixed Workload

# They try to tell it's the better AI hardware

<https://blogs.nvidia.com/blog/2017/04/10/ai-drives-rise-accelerated-computing-datacenter/>

	K80 2012	TPU 2015	P40 2016
<b>Inferences/Sec &lt;10ms latency</b>	$^{1/13}X$	1X	2X
<b>Training TOPS</b>	6 FP32	NA	12 FP32
<b>Inference TOPS</b>	6 FP32	90 INT8	48 INT8
<b>On-chip Memory</b>	16 MB	24 MB	11 MB
<b>Power</b>	300W	75W	250W
<b>Bandwidth</b>	320 GB/S	34 GB/S	350 GB/S



Start the presentation to see live content. Still no live content? Install the app or get help at [PollEv.com/app](https://PollEv.com/app)

# Inference per second

$$\frac{\text{Inferences}}{\text{Second}} = \frac{\text{Inferences}}{\text{Operation}} \times \frac{\text{Operations}}{\text{Second}}$$

$$= \frac{\text{Inferences}}{\text{Operation}} \times [\frac{\text{operations}}{\text{cycle}} \times \frac{\text{cycles}}{\text{second}} \times \#_{\_PEs} \times \text{Utilization}_{\_PEs}]$$

	Hardware	Model	Input Data
Operations per inference		v	
Operations per cycle	v		
Cycles per second	v		
Number of PEs	v		
Utilization of PEs	v	v	
Effectual operations out of (total) operations		v	v
Effectual operations plus unexploited ineffectual operations per cycle	v		

# What's wrong with inferences per second?

- There is no standard on how they inference — but these affect!
  - What model?
  - What dataset?
  - Quality?
- That's why Facebook is trying to promote an AI benchmark — MLPerf

- *Pitfall: For NN hardware, Inferences Per Second (IPS) is an inaccurate summary performance metric.*

Our results show that IPS is a poor overall performance summary for NN hardware, as it's simply the inverse of the complexity of the typical inference in the application (e.g., the number, size, and type of NN layers). For example, the TPU runs the 4-layer MLP1 at 360,000 IPS but the 89-layer CNN1 at only 4,700 IPS, so TPU IPS vary by 75X! Thus, using IPS as the single-speed summary is *even more misleading* for NN accelerators than MIPS or FLOPS are for regular processors [23], so IPS should be even more disparaged. To compare NN machines better, we need a benchmark suite written at a high-level to port it to the wide variety of NN architectures. Fathom is a promising new attempt at such a benchmark suite [3].



# Extreme Multitasking Performance

- Dual 4K external monitors
- 1080p device display
- 7 applications

# What's missing in this video clip?

- The ISA of the “competitor”
- Clock rate, CPU architecture, cache size, how many cores
- How big the RAM?
- How fast the disk?

# 12 ways to Fool the Masses When Giving Performance Results on Parallel Computers

- Quote only 32-bit performance results, not 64-bit results.
- Present performance figures for an inner kernel, and then represent these figures as the performance of the entire application.
- Quietly employ assembly code and other low-level language constructs.
- Scale up the problem size with the number of processors, but omit any mention of this fact.
- Quote performance results projected to a full system.
- Compare your results against scalar, unoptimized code on Crays.
- When direct run time comparisons are required, compare with an old code on an obsolete system.
- If MFLOPS rates must be quoted, base the operation count on the parallel implementation, not on the best sequential implementation.
- Quote performance in terms of processor utilization, parallel speedups or MFLOPS per dollar.
- Mutilate the algorithm used in the parallel implementation to match the architecture.
- Measure parallel run times on a dedicated system, but measure conventional run times in a busy environment.
- If all else fails, show pretty pictures and animated videos, and don't talk about performance.