

The Dusk and Dawn of Computer Architectures

Hung-Wei Tseng

Recap: 4Cs of cache misses

- 3Cs:
 - Compulsory, Conflict, Capacity
- Coherency miss:
 - A “block” invalidated because of the sharing among processors.
- True sharing
 - Processor A modifies X, processor B also want to access X.
- False sharing
 - Processor A modifies X, processor B also want to access Y.
However, Y is invalidated because X and Y are in the same block!

Performance comparison

- Comparing implementations of thread_vadd — L and R, please identify which one will be performing better and why

Version L

```
void *threaded_vadd(void *thread_id)
{
    int tid = *(int *)thread_id;
    int i;
    for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)
    {
        c[i] = a[i] + b[i];
    }
    return NULL;
}
```

Version R

```
void *threaded_vadd(void *thread_id)
{
    int tid = *(int *)thread_id;
    int i;
    for(i=tid*(ARRAY_SIZE/NUM_OF_THREADS);i<(tid+1)*(ARRAY_SIZE/NUM_OF_THREADS);i++)
    {
        c[i] = a[i] + b[i];
    }
    return NULL;
}
```

- A. L is better, because the cache miss rate is lower
- B. R is better, because the cache miss rate is lower
- C. L is better, because the instruction count is lower
- D. R is better, because the instruction count is lower
- E. Both are about the same

Main thread

```
for(i = 0 ; i < NUM_OF_THREADS ; i++)
{
    tids[i] = i;
    pthread_create(&thread[i], NULL, threaded_vadd, &tids);
}
for(i = 0 ; i < NUM_OF_THREADS ; i++)
    pthread_join(thread[i], NULL);
```

Again — how many values are possible?

- Consider the given program. You can safely assume the caches are coherent. How many of the following outputs will you see?

- ① (0, 0)
 - ② (0, 1)
 - ③ (1, 0)
 - ④ (1, 1)
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

```
#include <stdio.h>
#include <stdlib.h>
#include <pthread.h>
#include <unistd.h>

volatile int a,b;
volatile int x,y;
volatile int f;
void* modifya(void *z) {
    a=1;
    x=b;
    return NULL;
}
void* modifyb(void *z) {
    b=1;
    y=a;
    return NULL;
}
```

```
int main() {
    int i;
    pthread_t thread[2];
    pthread_create(&thread[0], NULL, modifya, NULL);
    pthread_create(&thread[1], NULL, modifyb, NULL);
    pthread_join(thread[0], NULL);
    pthread_join(thread[1], NULL);
    fprintf(stderr, "(%d, %d)\n", x, y);
    return 0;
}
```

Possible scenarios

Thread 1

a=1;

x=b;

Thread 2

b=1;
y=a;

(1,1)

Thread 1

a=1;
x=b;

Thread 2

b=1;
y=a;

(0,1)

Thread 1

a=1;
x=b;

Thread 2

b=1;
y=a;

(1,0)

Thread 1

x=b;
a=1;

Thread 2

y=a;

OoO Scheduling!

b=1;

(0,0)

Take-aways of parallel programming

- Processor behaviors are non-deterministic
 - You cannot predict which processor is going faster
 - You cannot predict when OS is going to schedule your thread
- Cache coherency only guarantees that everyone would eventually have a coherent view of data, but not when
- Cache consistency is hard to support

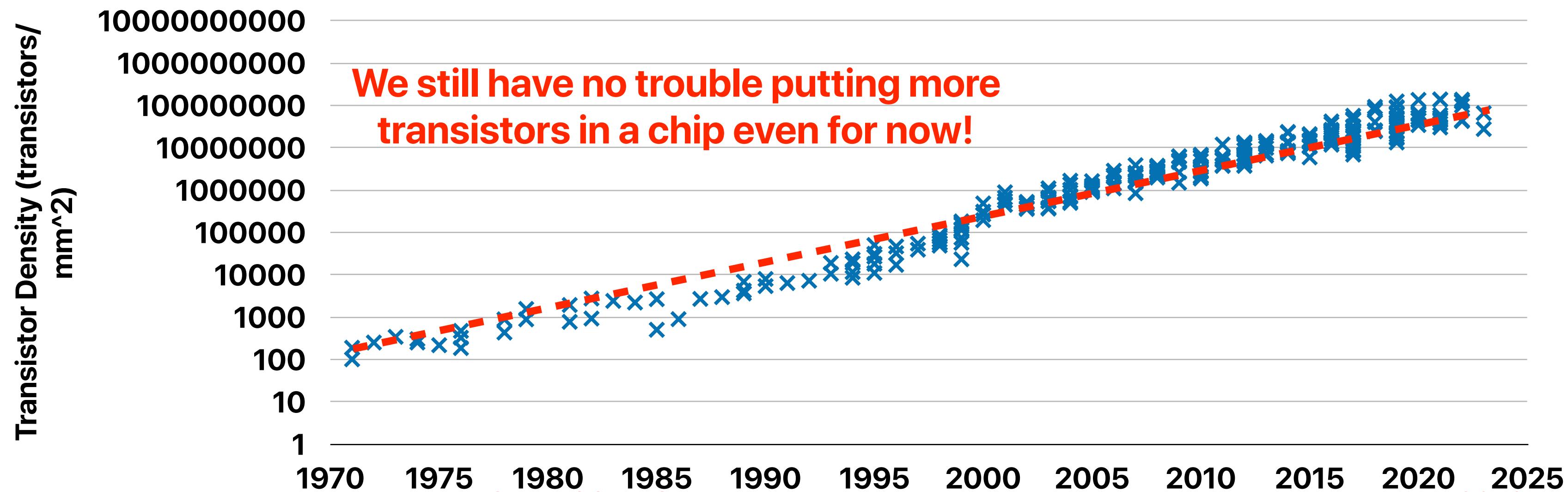
Outline

- Dark silicon and the Golden Age of Computer Architecture

The limiting factor of modern processor performance

Moore's Law⁽¹⁾

- The number of transistors we can build in a fixed area of silicon doubles every 12 ~ 24 months.
- Moore's Law "was" the most important driver for historic CPU performance gains



(1) Moore, G. E. (1965), 'Cramming more components onto integrated circuits', Electronics 38 (8).

Power v.s. Energy

- Power is the direct contributor of “heat”
 - Packaging of the chip
 - Heat dissipation cost
- $\text{Energy} = P * ET$
 - The electricity bill and battery life is related to energy!
 - Lower power does not necessarily means better battery life if the processor slow down the application too much

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor

$$P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$$

- α : average switches per cycle
- C : capacitance
- V : voltage
- f : frequency, usually linear with V
- N : the number of transistors

Static/Leakage Power

- The power consumption due to leakage — transistors do not turn all the way off during no operation
- Becomes the **dominant** factor in the most advanced process technologies.

$$P_{leakage} \sim N \times V \times e^{-V_t}$$

- N : number of transistors
- V : voltage
- V_t : threshold voltage where transistor conducts (begins to switch)

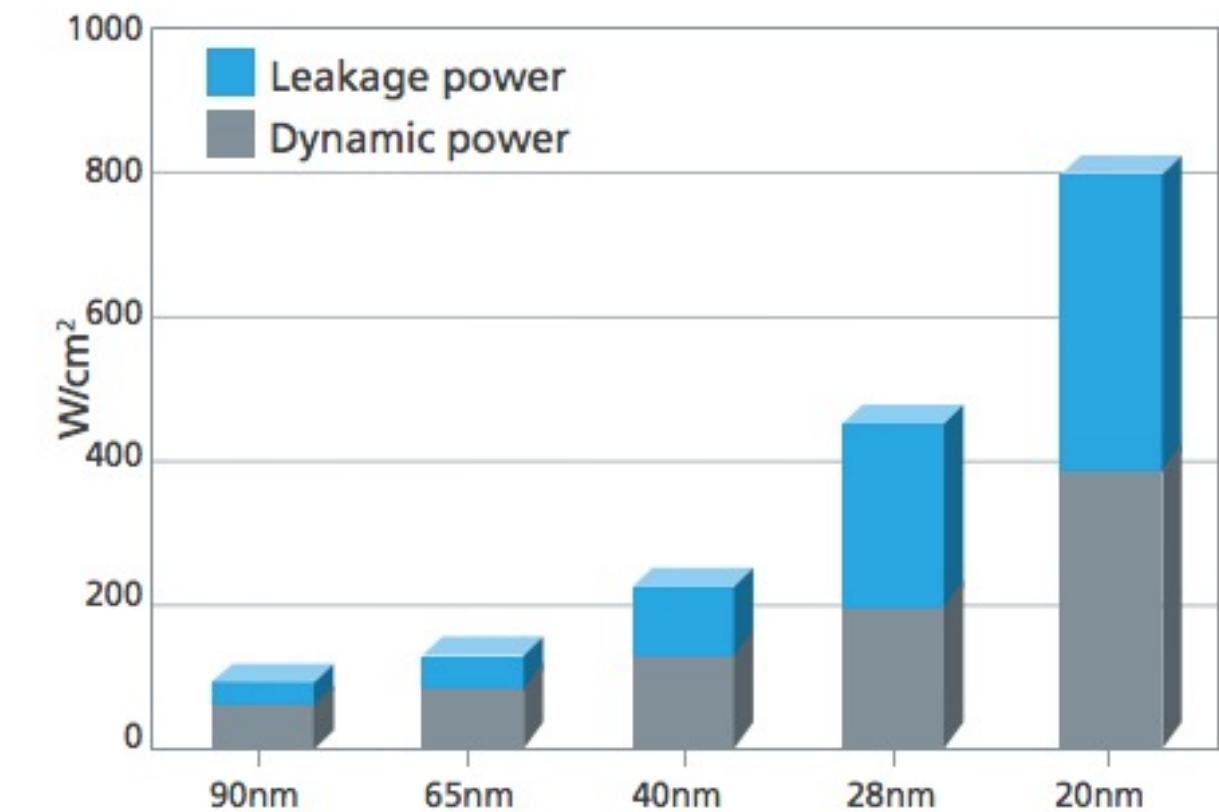


Figure 1: Leakage power becomes a growing problem as demands for more performance and functionality drive chipmakers to nanometer-scale process nodes (Source: IBS).



What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ② The power density of the chip will increase
 - ③ Given the same power budget, we may not able to power on all chip area if we maintain the same clock rate
 - ④ Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area

A. 0
B. 1
C. 2
D. 3
E. 4

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ② The power density of the chip will increase
 - ③ Given the same power budget, we may not be able to power on all chip area if we maintain the same clock rate
 - ④ Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area

A. 0

B. 1

C. 2

D. 3

E. 4

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ② The power density of the chip will increase
 - ③ Given the same power budget, we may not be able to power on all chip area if we maintain the same clock rate
 - ④ Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area

A. 0

B. 1

C. 2

D. 3

E. 4

Power consumption & power density

- The power consumption due to the switching of transistor states
- Dynamic power per transistor:

$$P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$$

- α : average switches per cycle
- C : capacitance
- V : voltage
- f : frequency, usually linear with V
- N : the number of transistors

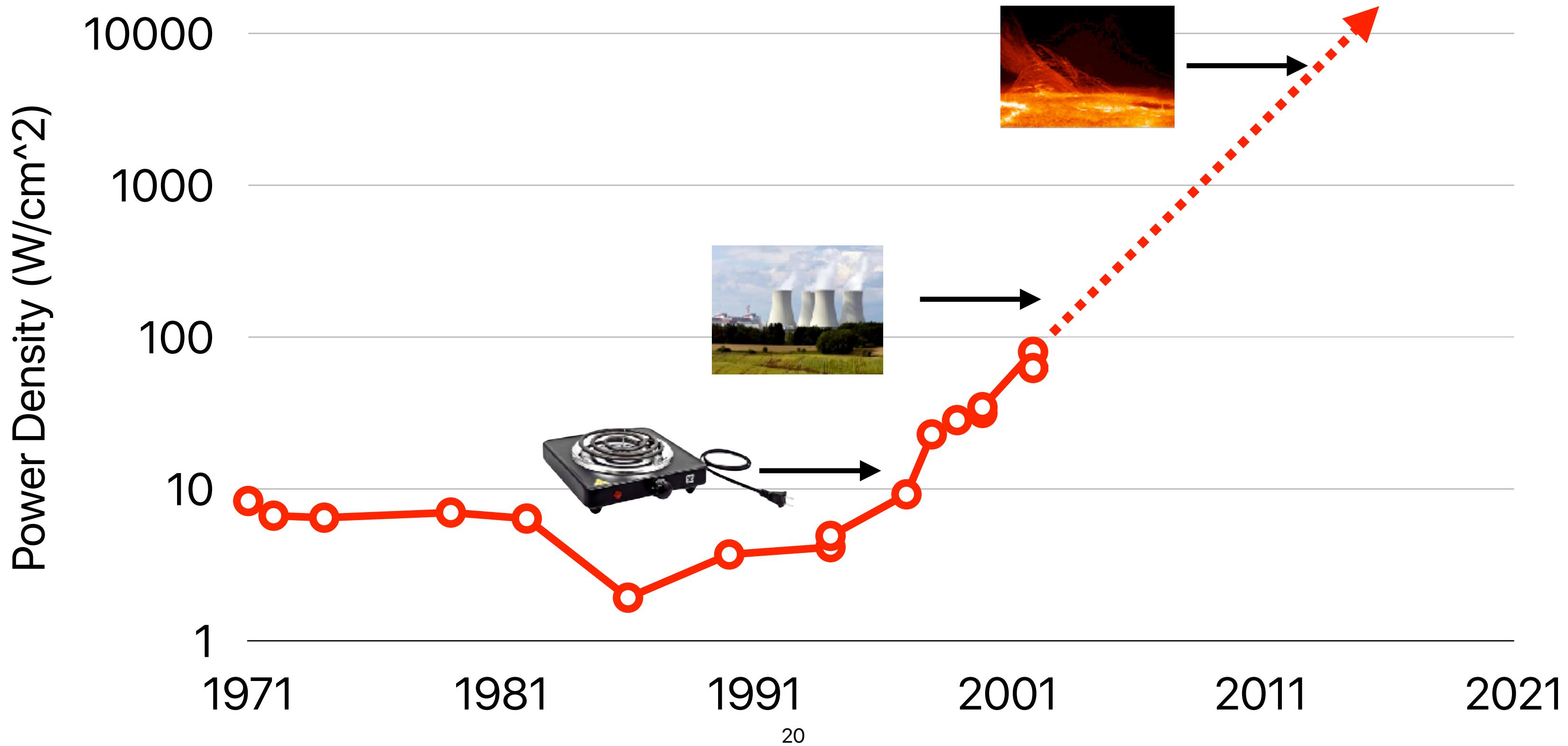
Dennard scaling discontinued — we cannot make voltage lower

- Power density:

$$P_{density} = \frac{P}{area}$$

Moore's Law allows higher frequencies as transistors are smaller
Moore's Law makes this smaller

Power Density of Processors



Power consumption to light on all transistors

Chip

1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1

=49W

Dennardian Scaling

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W

Dennardian Broken

Chip

1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1

On ~ 50W
Off ~ 0W
Dark!

=100W!

Or we have to dim down all transistors

Chip

1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1
1	1	1	1	1	1	1

=49W

Dennardian Scaling

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W

Dennardian Broken

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

Low frequency

~0.5W

=50W!

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ② The power density of the chip will increase
 - ③ Given the same power budget, we may not able to power on all chip area if we maintain the same clock rate
 - ④ Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area

A. 0

B. 1

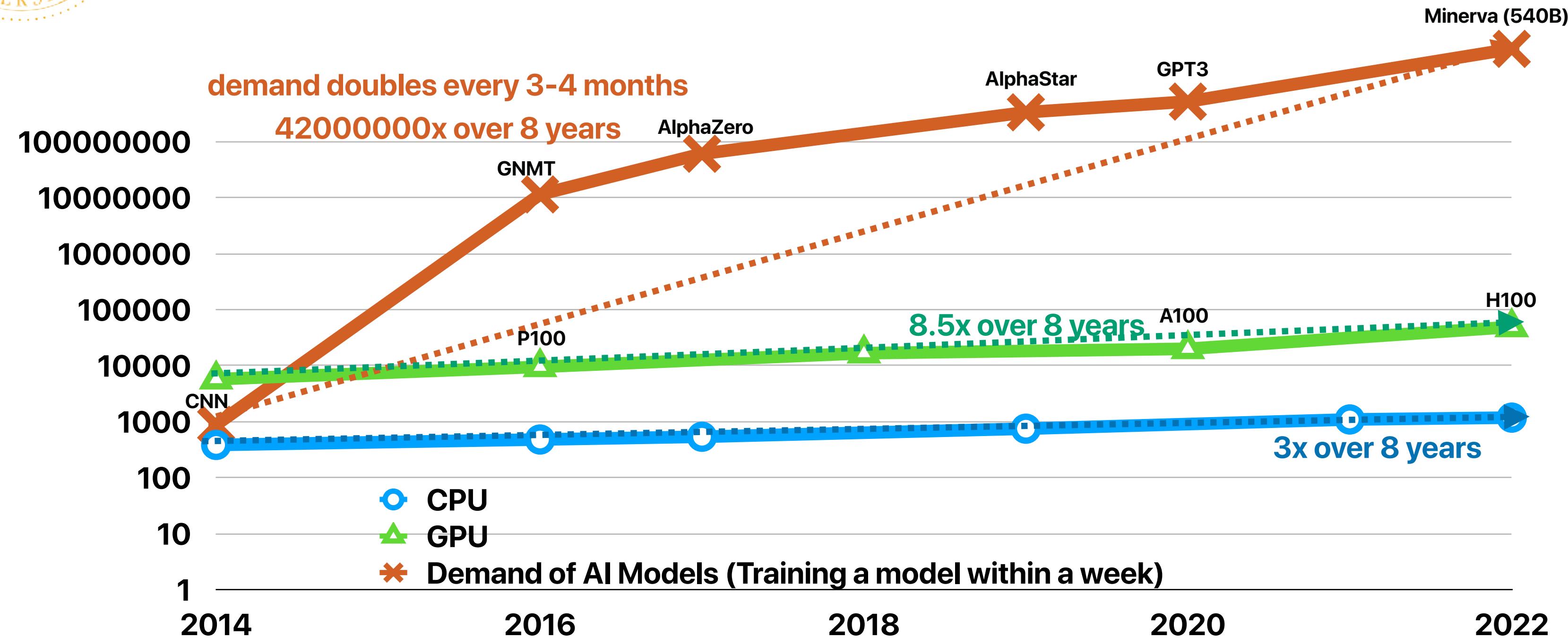
C. 2

D. 3

E. 4



Mis-matching AI/ML demand and general-purpose processing



<https://ourworldindata.org/grapher/artificial-intelligence-training-computation>

Solutions/trends in dark silicon era

Trends in the Dark Silicon Era

- Aggressive dynamic voltage/frequency scaling
- Throughout oriented — slower, but more
- Just let it dark — activate part of circuits, but not all
- From general-purpose to domain-specific — ASIC

Aggressive dynamic frequency scaling

Modern processor's frequency

Socket(s):	1	
NUMA node(s):	1	
Vendor ID:	GenuineIntel	
CPU family:	6	
Model:	151	i7-12700K
Model name:	12th Gen Intel(R) Core(TM) i7-12700KF	
Stepping:	2	Intel 7
CPU MHz:	1226.409	\$450.00 - \$460.00
CPU max MHz:	5000.0000	
CPU min MHz:	800.0000	PC/Client/Tablet, Workstation
Boost MPS:	7219.20	

CPU Specifications

Total Cores	12
# of Performance-cores	8
# of Efficient-cores	4
Total Threads	20
Max Turbo Frequency	5.00 GHz
Intel® Turbo Boost Max Technology 3.0 Frequency ¹	5.00 GHz
Performance-core Max Turbo Frequency	4.90 GHz
Efficient-core Max Turbo Frequency	3.80 GHz
Performance-core Base Frequency	3.50 GHz
Efficient-core Base Frequency	2.70 GHz
Cache	25 MB Intel® Smart Cache

Modern processor's frequency

Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	48 bits physical, 48 bits virtual
CPU(s):	12
On-line CPU(s) list:	0-11
Thread(s) per core:	2
Core(s) per socket:	6
Socket(s):	1
NUMA node(s):	1
Vendor ID:	AuthenticAMD
CPU family:	25
Model:	80
Model name:	AMD Ryzen 5 5500
Stepping:	0
Frequency boost:	enabled
CPU MHz:	3600.000
CPU max MHz:	3600.0000
CPU min MHz:	1400.0000
Processor	71%

Power consumption & power density

- The power consumption due to the switching of transistor states
- Dynamic power per transistor:

$$P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$$

- α : average switches per cycle
- C : capacitance
- V : voltage
- f : frequency, usually linear with V
- N : the number of transistors

Lower the frequency can reduce active power

- Power density:

$$P_{density} = \frac{P}{area}$$

Moore's Law allows higher frequencies as transistors are smaller
Moore's Law makes this smaller



Aggressive frequency scaling

- Regarding dynamic frequency scaling, how many of the following statements are correct?
 - ① Lowering the frequency helps extending the battery life
 - ② Lowering the frequency helps reducing the heat generation
 - ③ Lowering the frequency helps reducing the electricity bill
 - ④ A CPU operating at 25% of the peak frequency can still consume more than 50% of the peak power

A. 0
B. 1
C. 2
D. 3
E. 4

Power & Energy

- Regarding dynamic frequency scaling, how many of the following statements are correct?
 - ① Lowering the frequency helps extending the battery life
 - ② Lowering the frequency helps reducing the heat generation
 - ③ Lowering the frequency helps reducing the electricity bill
 - ④ A CPU operating at 25% of the peak frequency can still consume more than 50% of the peak power

A. 0

B. 1

C. 2

D. 3

E. 4

Recap: Demo — changing the max frequency and performance

- Change the maximum frequency of the intel processor — you learned how to do this when we discuss programmer's impact on performance
- LIKWID a profiling tool providing power/energy information
 - likwid-perfctr -g ENERGY [command_line]
 - Let's try blockmm and popcorn and see what's happening!

Power & Energy

- Regarding dynamic frequency scaling, how many of the following statements are correct?
 - ① Lowering the frequency helps extending the battery life
 - ✓ ② Lowering the frequency helps reducing the heat generation
 - ③ Lowering the frequency helps reducing the electricity bill
 - ✓ ④ A CPU operating at 25% of the peak frequency can still consume more than 50% of the peak power
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

Slower, but more

Transistor counts

Microarchitecture	Transistor Count	Issue-width	Year
Alder Lake	325 M	5x ALU, 7x Memory	2021
Coffee Lake	217 M	4x ALU, 4x Memory	2017
Sandy Bridge	290 M	3x ALU, 3x Memory	2011
Nehalem	182.75 M	3x ALU, 3x Memory	2008



How many transistors per core on Coffee Lake?



The Coffee Lake processor has 217 million transistors per core. It is manufactured using Intel's second 14 nm process. The Coffee Lake processors introduced i5 and i7 CPUs featuring six cores (along with hyper-threading in the case of the later models) and no hyperthreading.



The transistor count per core on Coffee Lake is lower than that of some other modern processors, such as the Ryzen 5 5600X which has 390 million transistors per core. However, Coffee Lake still offers good performance, thanks to its high clock speeds and efficient power delivery.

Here is a table of the transistor counts per core for some other modern processors:

Processor	Transistors per core
Coffee Lake	217 million
Ryzen 5 5600X	390 million
Core i7 12700K	475.8 million

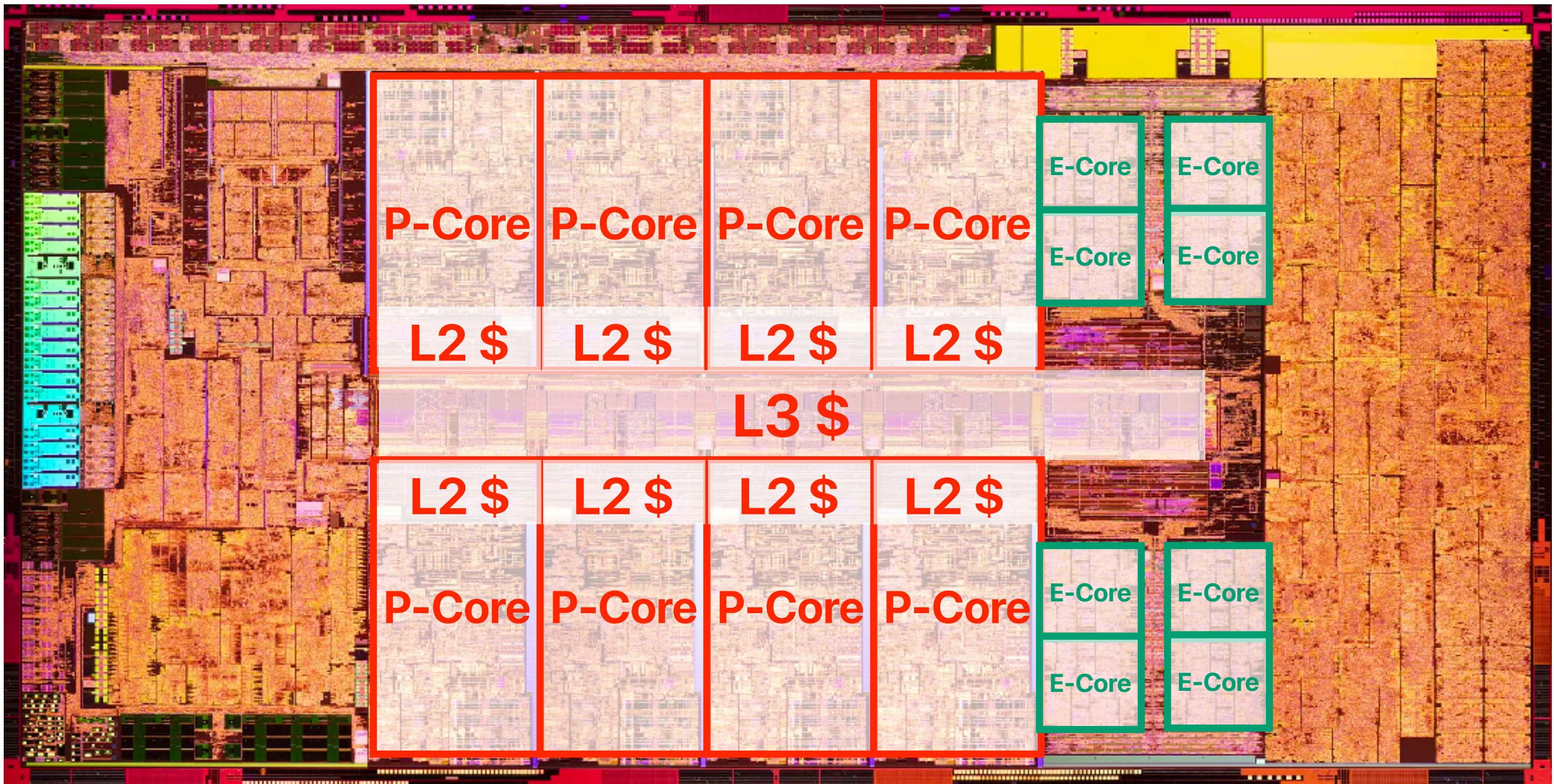
2x 3-issue ALUs Nehalem

Nehalem Alder Lake Nehalem
6-issue 12-issue 6-issue

1x 5-issue ALUs Alder Lake

Based on https://en.wikipedia.org/wiki/Transistor_count

Intel Alder Lake

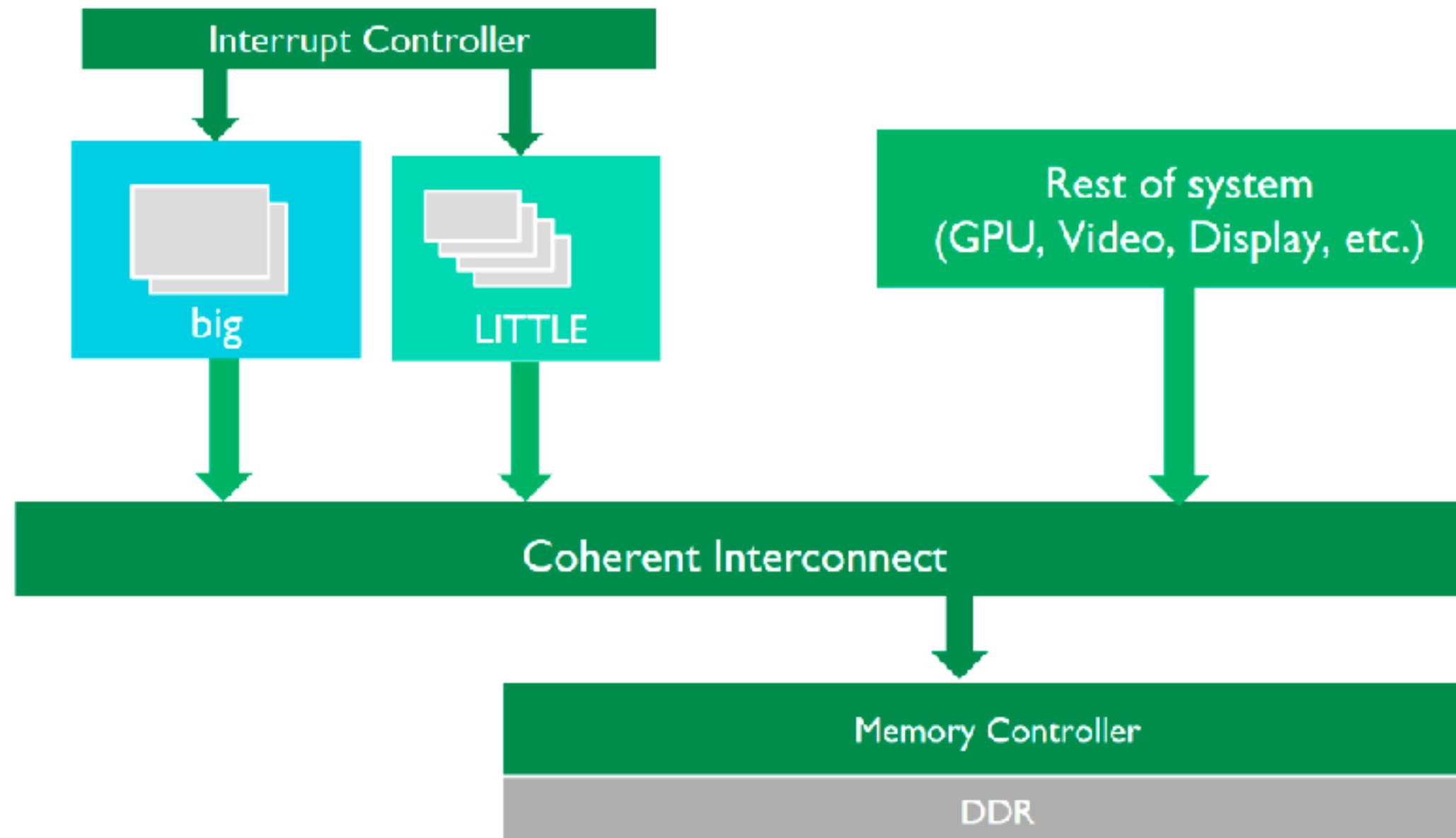


Single ISA heterogeneous CMP in Intel Processors

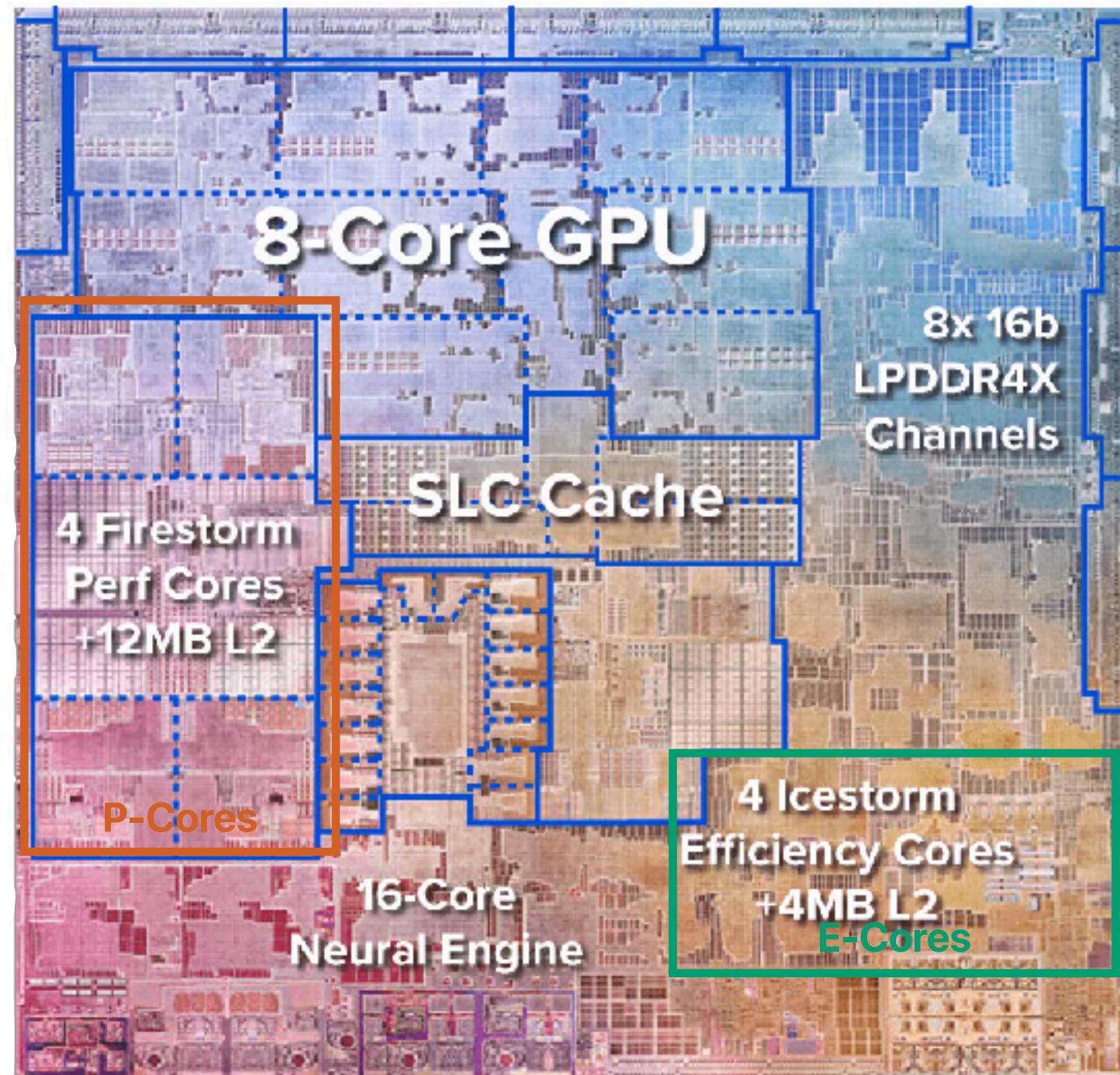


Single ISA heterogeneous CMP in ARM's big.LITTLE architecture

big.LITTLE system



Single ISA heterogeneous CMP in Apple's M1



Demo

- We can use taskset command in linux to control the task allocation.
- Core i7 13700 is a processor with big.Little architecture

Essentials

Product Collection	13th Generation Intel® Core™ i7 Processors
Code Name	Products formerly Raptor Lake
Vertical Segment	Desktop
Processor Number <small>?</small>	i7-13700
Lithography <small>?</small>	Intel 7
Recommended Customer Price <small>?</small>	\$384.00 - \$394.00
Use Conditions <small>?</small>	PC/Client/Tablet, Workstation

CPU Specifications

Total Cores <small>?</small>	16
# of Performance-cores	8
# of Efficient-cores	8
Total Threads <small>?</small>	24



Single ISA heterogeneous CMP (big.Little)

- Regarding “Single-ISA Heterogeneous Multi-Core Architectures”, how many of the following statements is/are correct?
 - ① You need to recompile and optimize the binary for each core architecture to exploit the thread-level parallelism in this architecture
 - ② For a program with limited thread-level parallelism, single ISA heterogeneous CMP would deliver better than or at least the same level of performance as homogeneous CMP built with older-generation cores
 - ③ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with older-generation cores
 - ④ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with newer-generation cores

A. 0
B. 1
C. 2
D. 3
E. 4

Single ISA heterogeneous CMP

- Regarding “Single-ISA Heterogeneous Multi-Core Architectures”, how many of the following statements is/are correct?
 - ① You need to recompile and optimize the binary for each core architecture to exploit the thread-level parallelism in this architecture
 - ② For a program with limited thread-level parallelism, single ISA heterogeneous CMP would deliver better than or at least the same level of performance as homogeneous CMP built with older-generation cores
 - ③ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with older-generation cores
 - ④ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with newer-generation cores
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4
- Corollary #4: Exploiting more parallelism from a program is the key to performance gain in modern architectures

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

More cores per chip, slower per core

	Intel® Xeon® Platinum 849..	Intel® Xeon® Platinum 846..	Intel® Xeon® Gold 6448H ..	Intel® Xeon® Platinum 844..	Intel® Xeon® Gold 6434H ..
Total Cores	60	48	32	16	8
Total Threads	120	96	64	32	16
Max Turbo Frequency	3.50 GHz	3.80 GHz	4.10 GHz	4.00 GHz	4.10 GHz
Processor Base Frequency	1.90 GHz	2.10 GHz	2.40 GHz	2.90 GHz	3.70 GHz
Cache	112.5 MB	105 MB	50 MB	45 MB	22.5 MB
Intel® UPI Speed	16 GT/s	16 GT/s	16 GT/s	16 GT/s	16 GT/s
Max # of UPI Links	4	4	3	4	3
TDP	350 W	330 W	250 W	270 W	195 W

Xeon Phi

Essentials

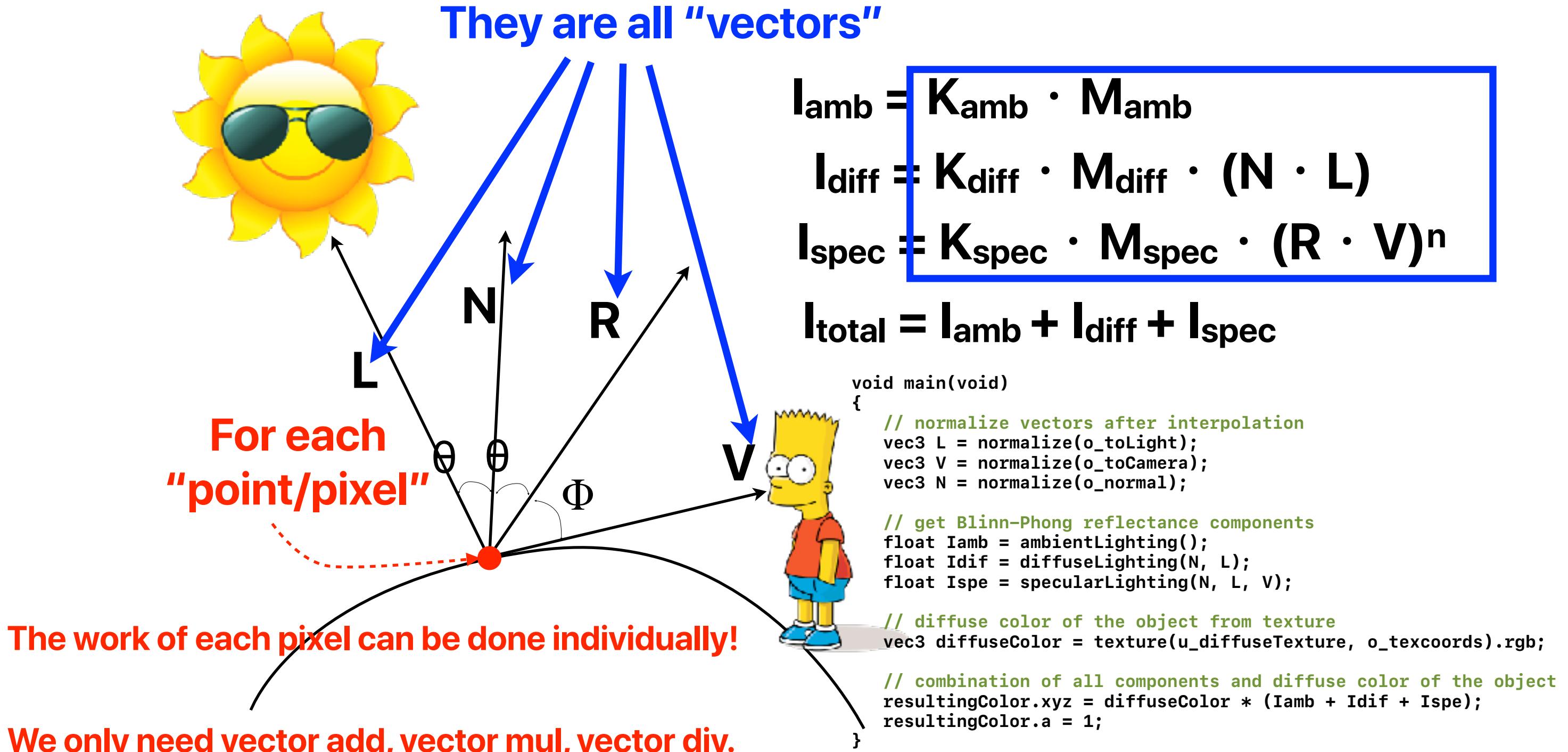
Product Collection	Intel® Xeon Phi™ 72x5 Processor Family
Code Name	Products formerly Knights Mill
Vertical Segment	Server
Processor Number	7295
Off Roadmap	No
Status	Launched
Launch Date ?	Q4'17
Lithography ?	14 nm

Performance

# of Cores ?	72
# of Threads ?	72
Processor Base Frequency ?	1.50 GHz
Max Turbo Frequency ?	1.60 GHz
Cache ?	36 MB L2 Cache
TDP ?	320 W

Slower, but more — GPU

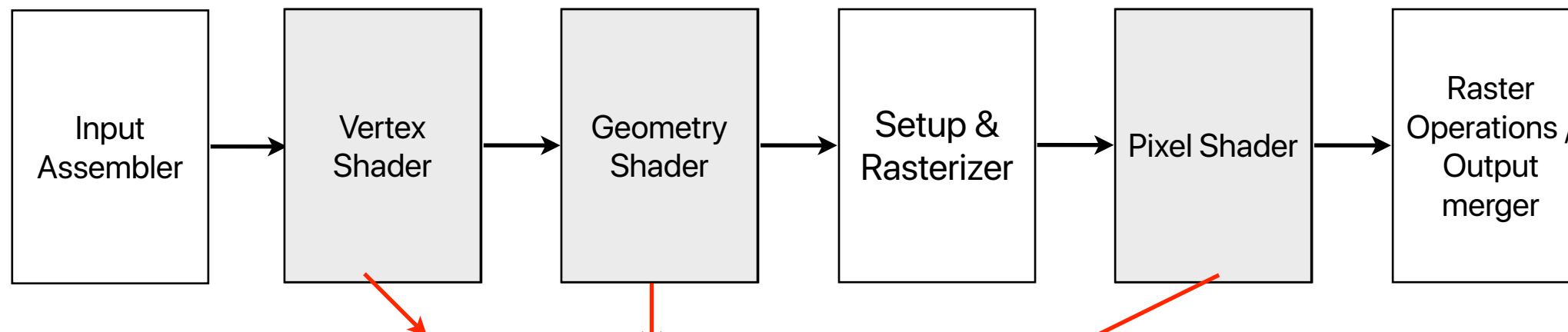
Basic concept of shading



GPU (Graphics Processing Unit)

- Originally for displaying images
- HD video: 1920×1080 pixels * 60 frames per second
 - Therefore, GPU is not latency-oriented by design!
 - Even for 120 frames, you still have 8ms latency to get everything done!
- Graphics processing pipeline

1 GHz can give you 8000000 cycles!!!

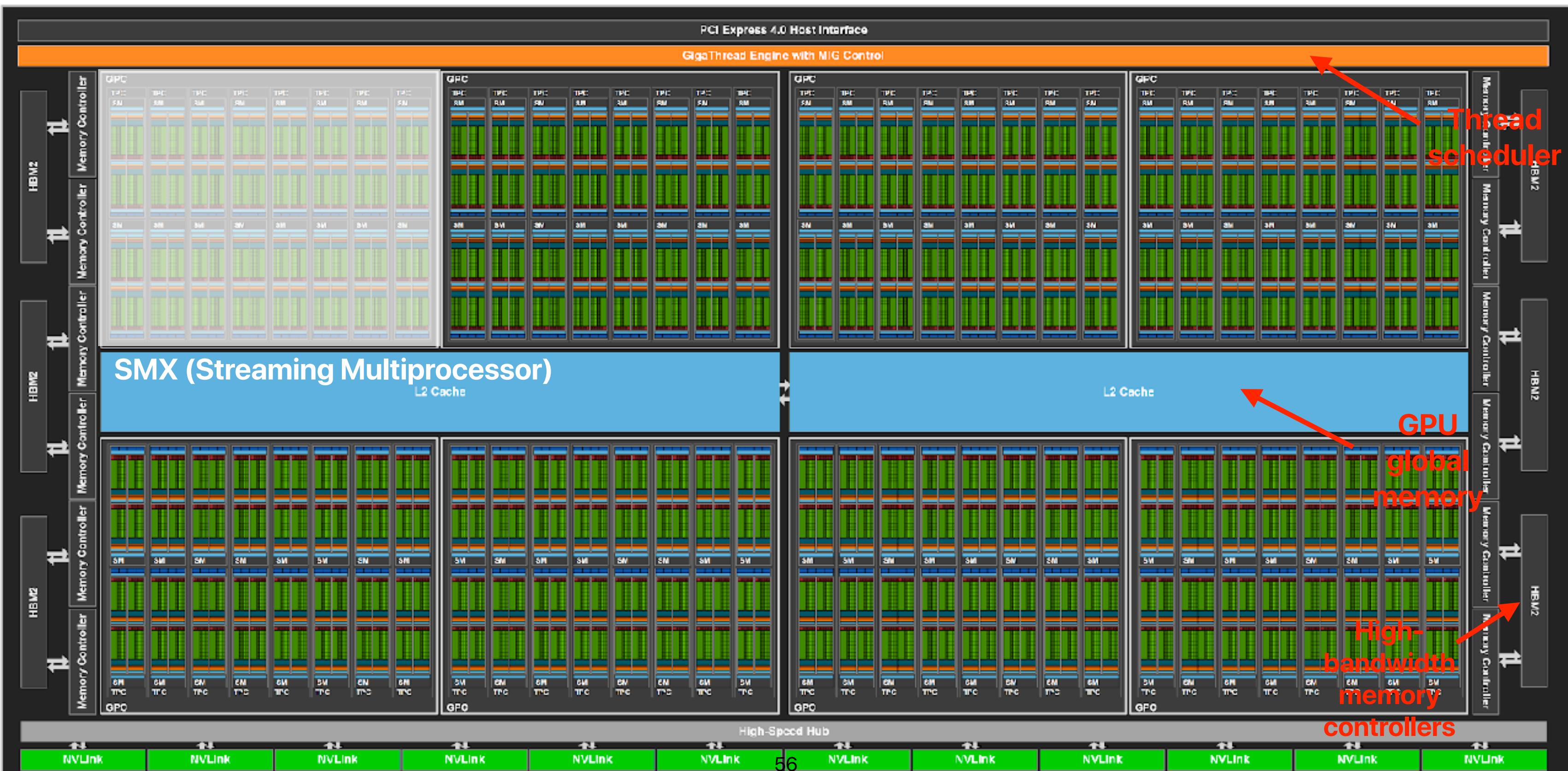


These shaders need to be “programmable” to apply different rendering effects/algorithms
(Phong shading, Gouraud shading, and etc...)

What's the “appropriate” GPU architecture

- Lots of ALUs to process pixels in parallel — 2M pixels in HD resolution, very regular workloads
 - Vector processing model
- Simple operations
 - The ALUs only supports very few instructions
 - Almost no branches
- Deadline driven and throughput-oriented rather than latency oriented
 - High-bandwidth but also “higher-latency” memory
 - ALUs can be slower

GPU Architecture



Inside an SM



A total of $16 \times 4 = 64$ FP32 cores
A total of $16 \times 4 = 64$ INT32 cores
A total of $16 \times 4 = 16$ FP64 cores

- All of these can only perform the same operation at the same time, but each of these is named as a “thread” in CUDA
- You can only use either FP32, FP64, INT32 and “Tensor Cores” at the same time

Static/Leakage Power

- The power consumption due to leakage — transistors do not turn all the way off during no operation
- Becomes the **dominant** factor in the most advanced process technologies.

$P_{leakage} \sim N \times V \times e^{-V/V_t}$ **How about static power?**

- N : number of transistors
- V : voltage
- V_t : threshold voltage where transistor conducts (begins to switch)

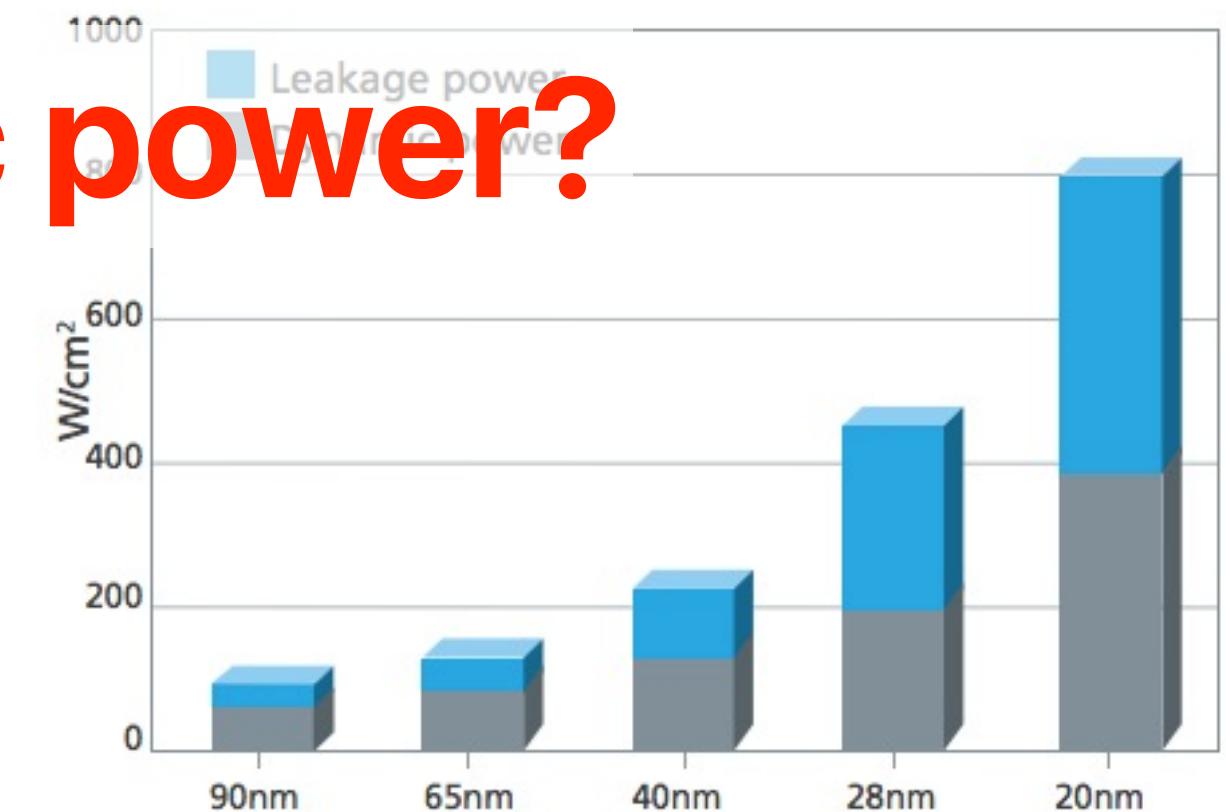


Figure 1: Leakage power becomes a growing problem as demands for more performance and functionality drive chipmakers to nanometer-scale process nodes (Source: IBS).

Just let it dark

Programming in Turing Architecture

Use tensor cores

```
cublasErrCheck(cublasSetMathMode(cublasHandle, CUBLAS_TENSOR_OP_MATH));
```

Make them 16-bit

```
convertFp32ToFp16 <<< (MATRIX_M * MATRIX_K + 255) / 256, 256 >>> (a_fp16, a_fp32,  
MATRIX_M * MATRIX_K);  
convertFp32ToFp16 <<< (MATRIX_K * MATRIX_N + 255) / 256, 256 >>> (b_fp16, b_fp32,  
MATRIX_K * MATRIX_N);
```

```
cublasErrCheck(cublasGemmEx(cublasHandle, CUBLAS_OP_N, CUBLAS_OP_N,  
    MATRIX_M, MATRIX_N, MATRIX_K,  
    &alpha,  
    a_fp16, CUDA_R_16F, MATRIX_M,  
    b_fp16, CUDA_R_16F, MATRIX_K,  
    &beta,  
    c_cublas, CUDA_R_32F, MATRIX_M,  
    CUDA_R_32F, CUBLAS_GEMM_DFALT_TENSOR_OP));
```

call Gemm

Inside an SM



A total of $16 \times 4 = 64$ FP32 cores

A total of $16 \times 4 = 64$ INT32 cores

A total of $16 \times 4 = 16$ FP64 cores

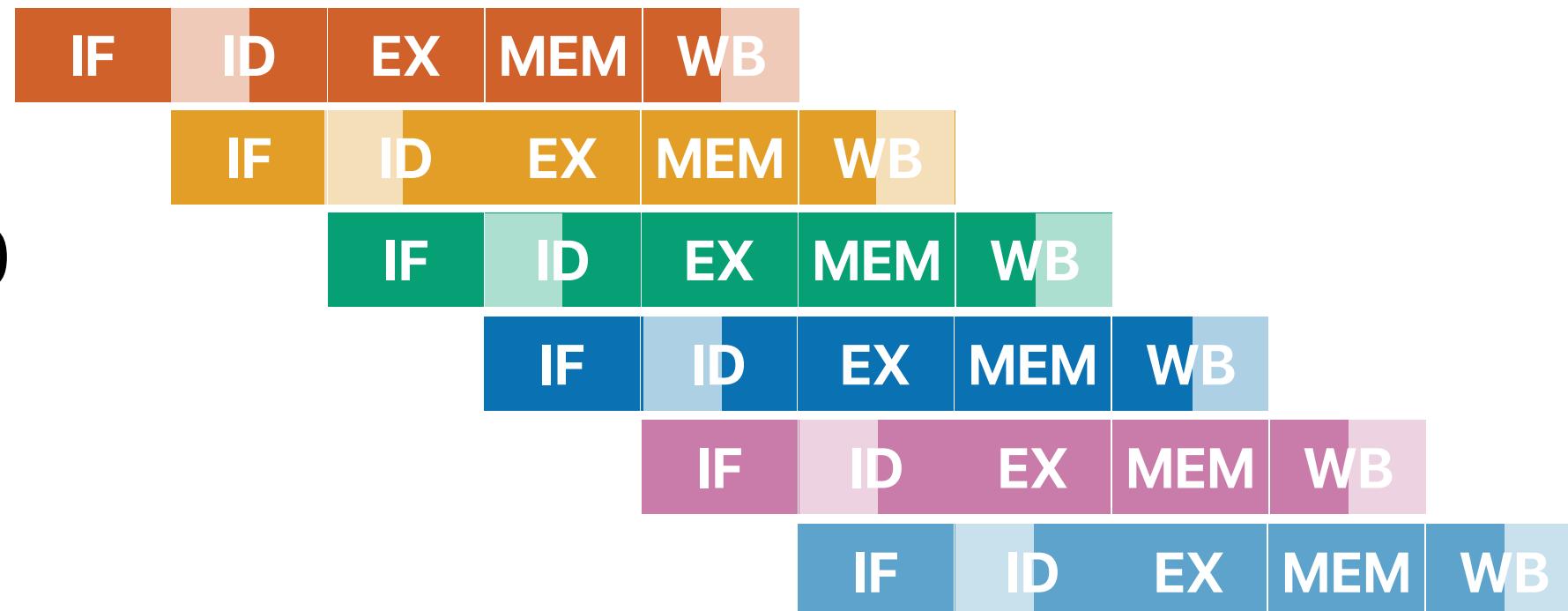
- All of these can only perform the same operation at the same time, but each of these is named as a "thread" in CUDA
- You can only use either FP32, FP64, INT32 or "Tensor Cores" at the same time

The Rise of ASICs/Accelerators — A New “Golden” Age of Architects

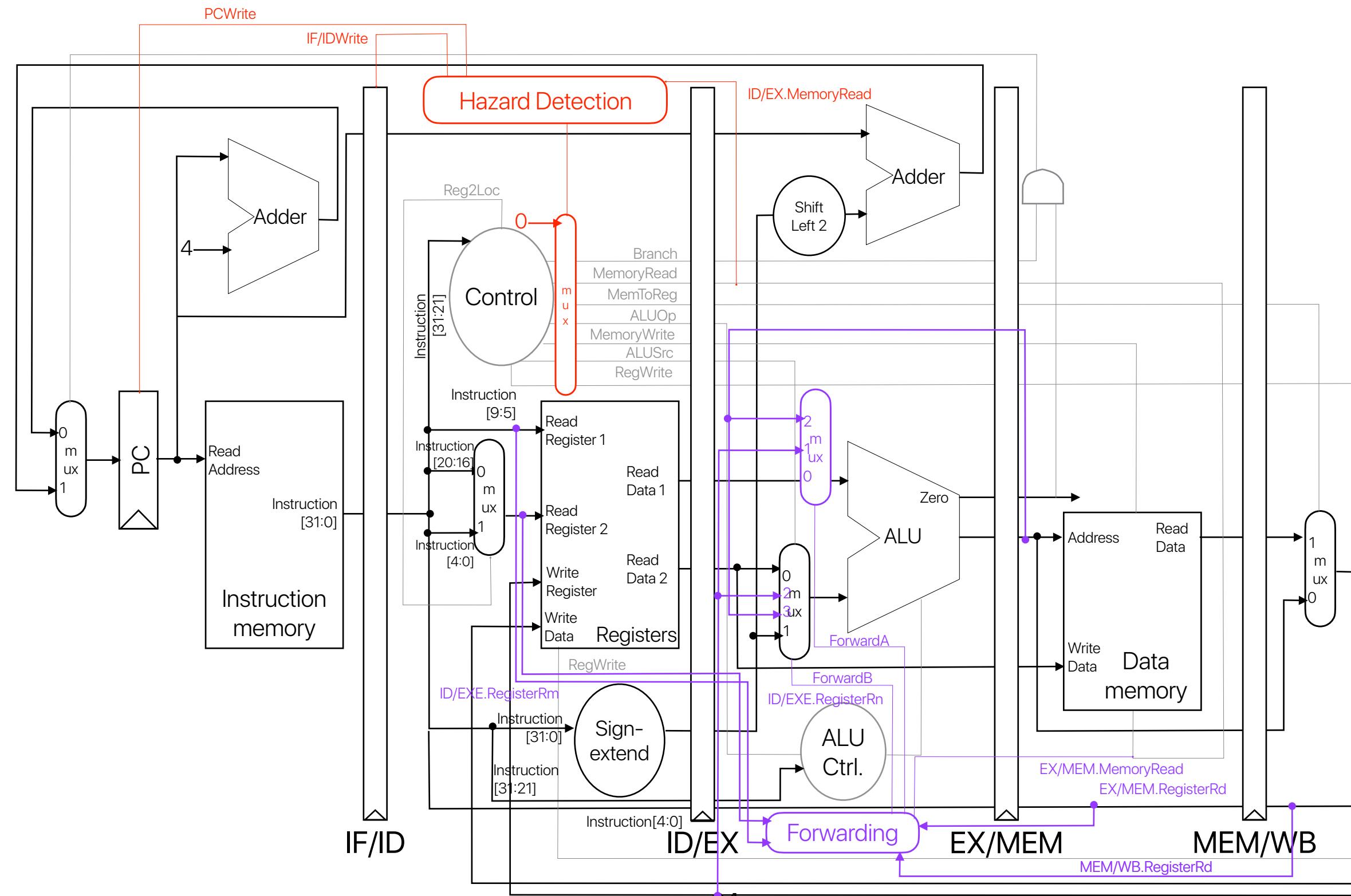
Say, we want to implement $a[i] += a[i+1]*20$

- This is what we need in RISC-V in each iteration

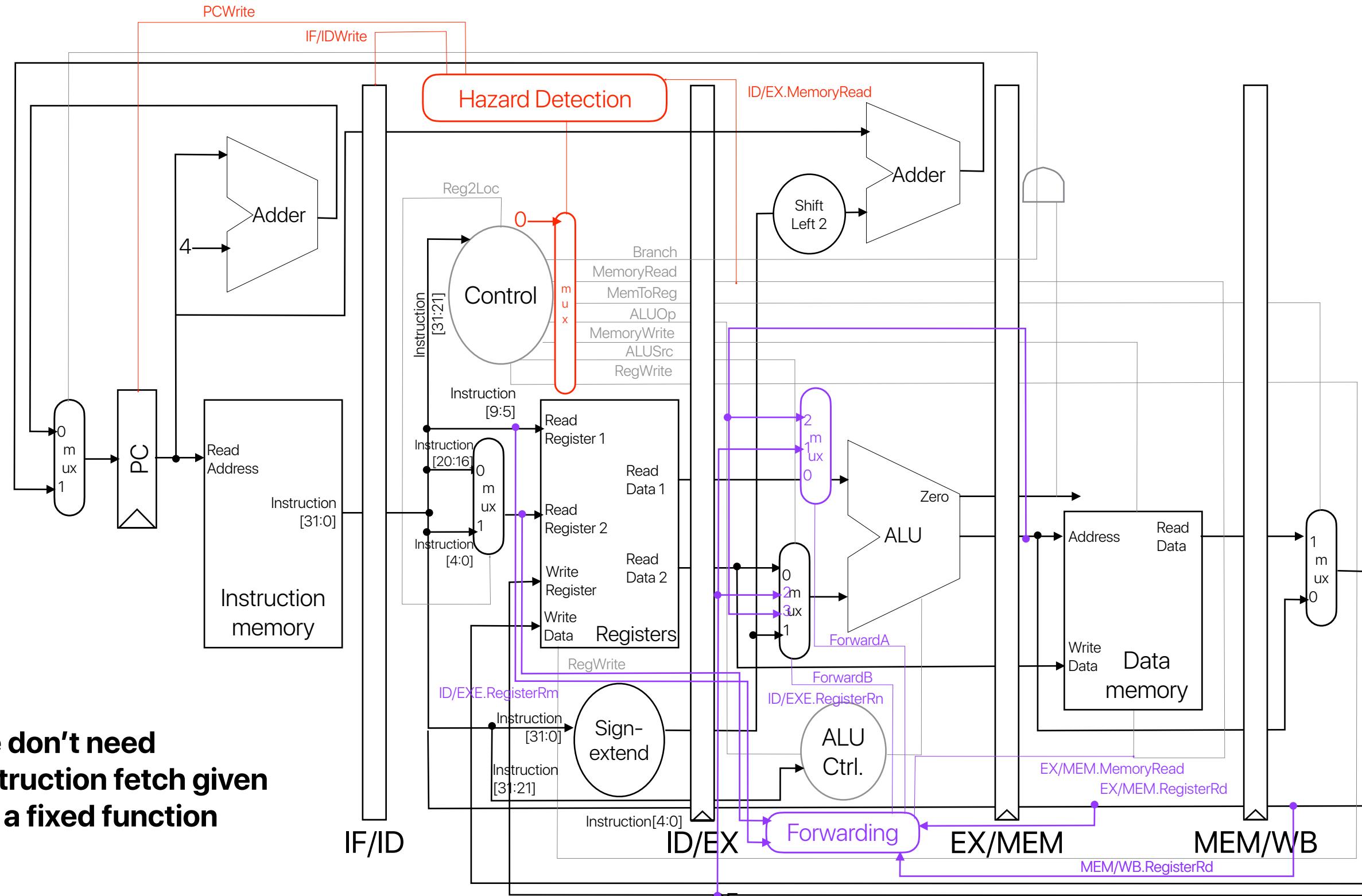
ld	X1,	0(X0)
ld	X2,	8(X0)
add	X3,	X31, #20
mul	X2,	X2, X3
add	X1,	X1, X2
sd	X1,	0(X0)



This is what you need for these instructions



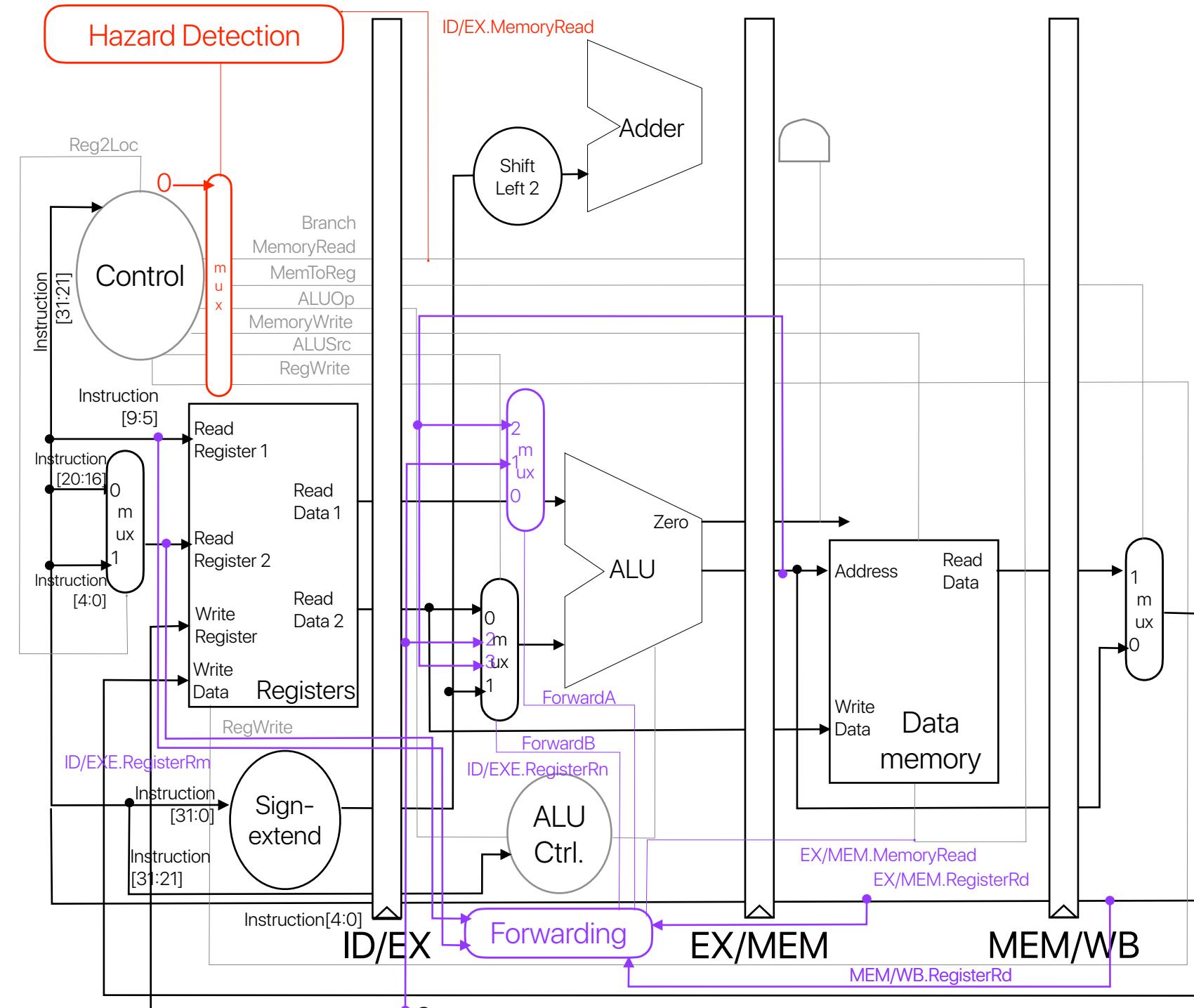
Specialize the circuit



Specialize the circuit

We don't need these many registers, complex control, decode

We don't need instruction fetch given it's a fixed function

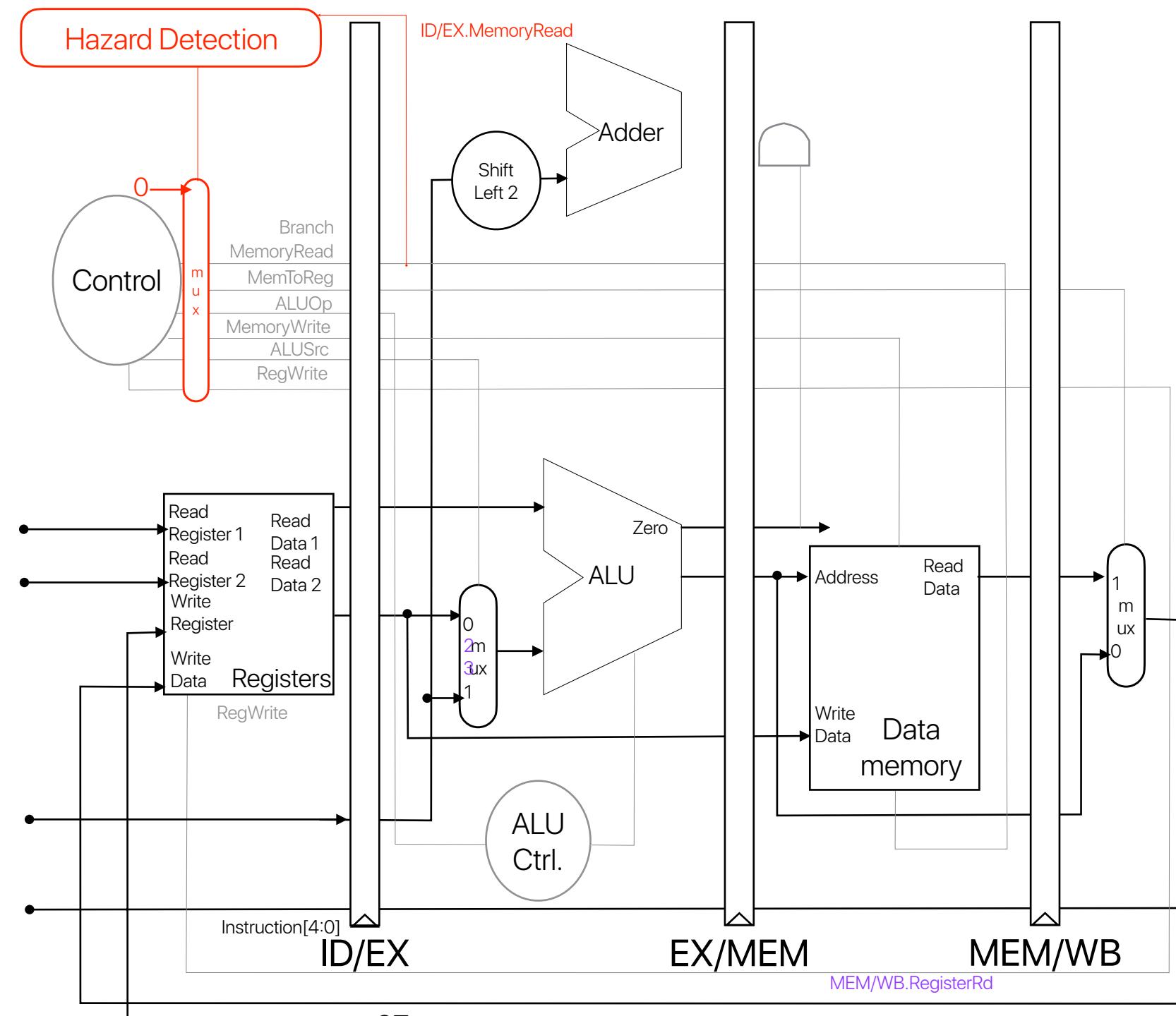


Specialize the circuit

We don't need ALUs,
branches, hazard
detections...

We don't need these
many registers, complex
control, decode

We don't need
instruction fetch given
it's a fixed function

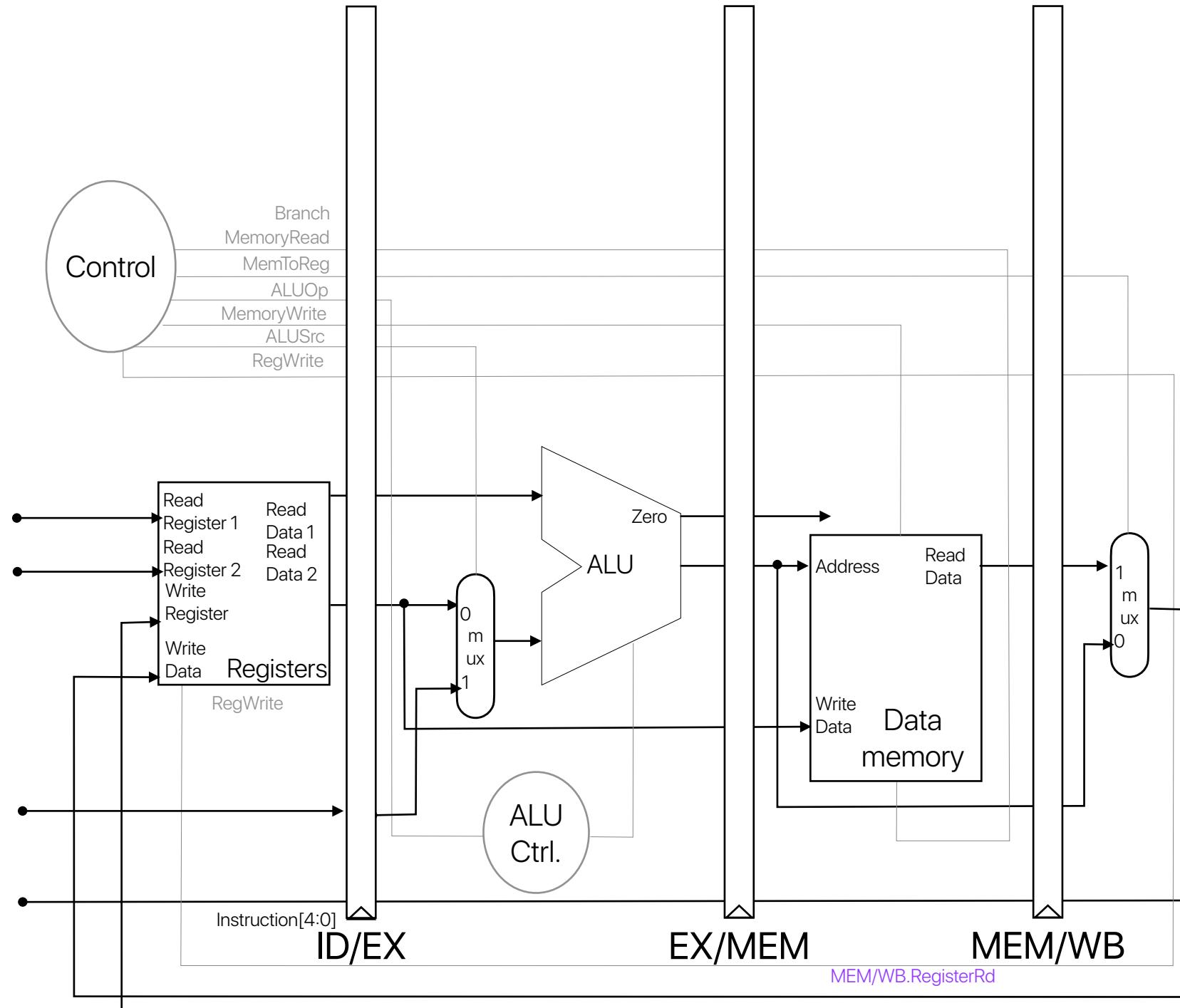


Specialize the circuit

We don't need big ALUs,
branches, hazard
detections...

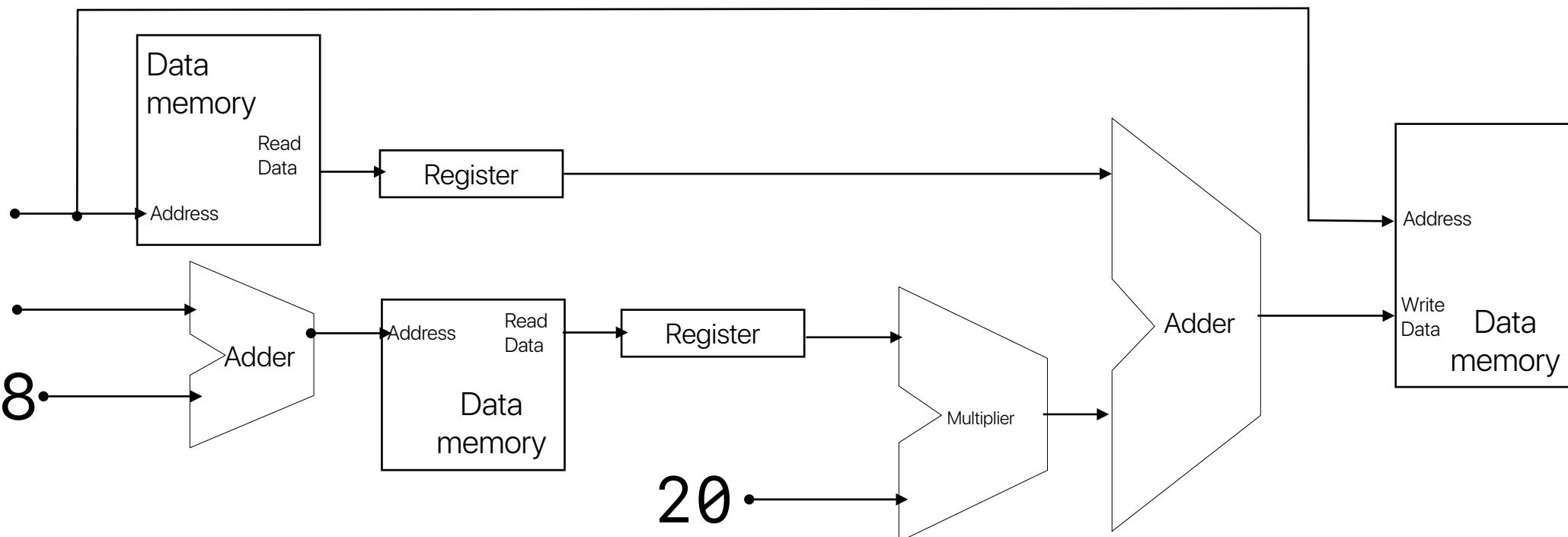
We don't need these
many registers, complex
control, decode

We don't need
instruction fetch given
it's a fixed function



Rearranging the datapath

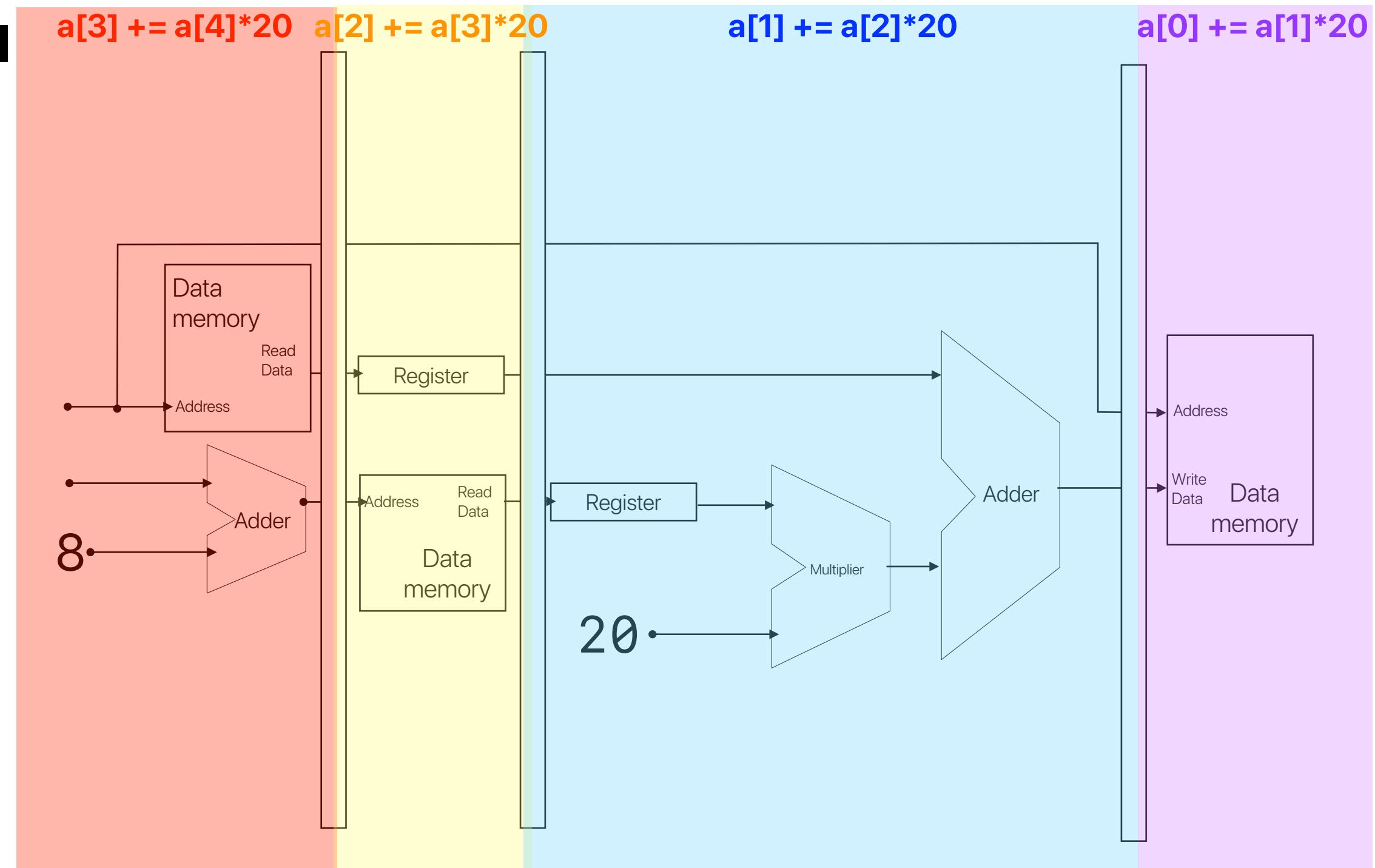
```
ld    X1, 0(X0)
ld    X2, 8(X0)
add   X3, X31, #20
mul   X2, X2, X3
add   X1, X1, X2
sd    X1, 0(X0)
```



The pipeline for $a[i] += a[i+1]*20$

**Each stage can still
be as fast as the
pipelined
processor**

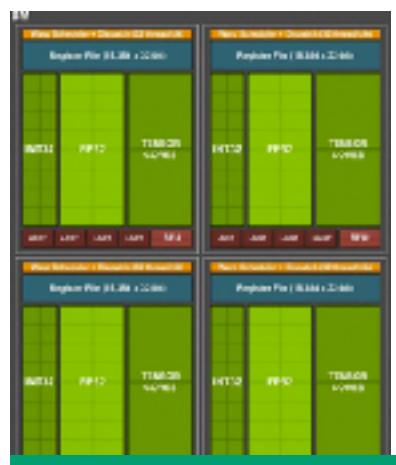
**But each stage is
now working on
what the original 6
instructions would
do**



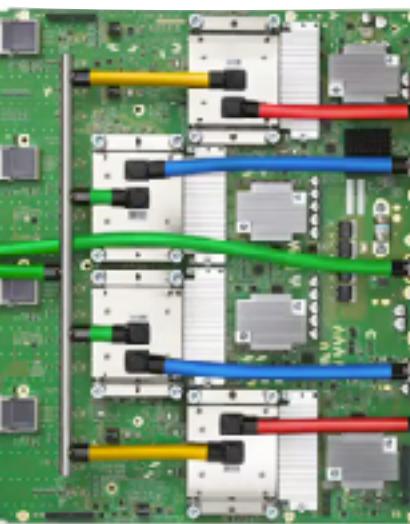


The “landscape” of modern computers

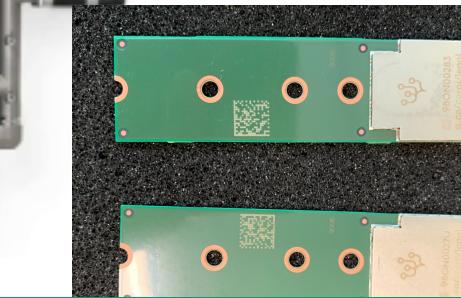
Google Datacenter TPUs



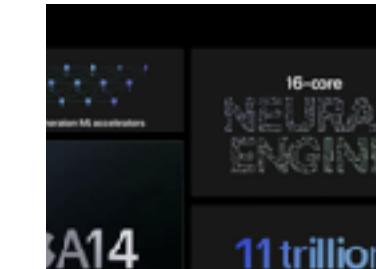
NVIDIA Tensor Core Units



Google Edge TPUs



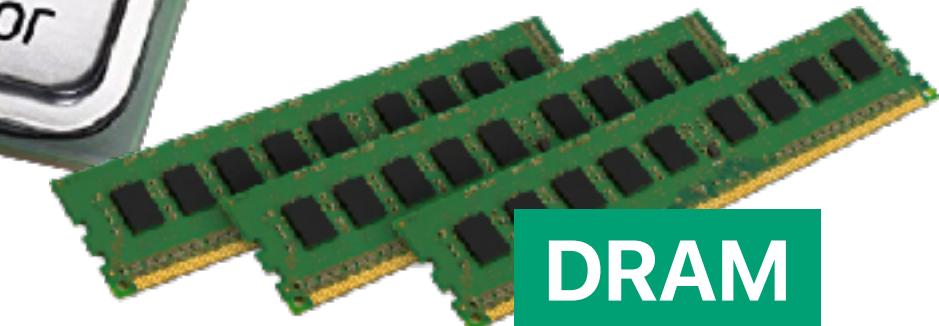
AI/ML Accelerators



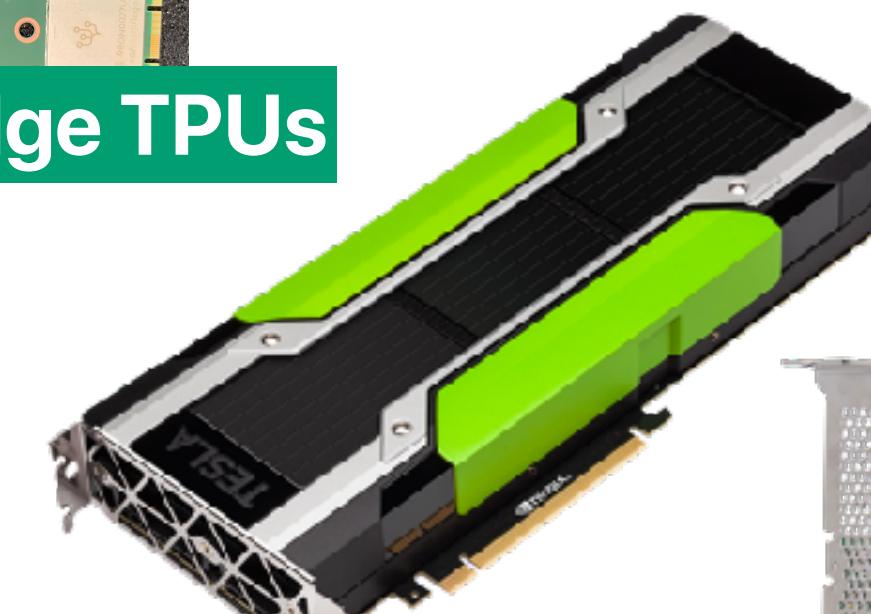
Apple Neural Engines



CPU



DRAM



GPU

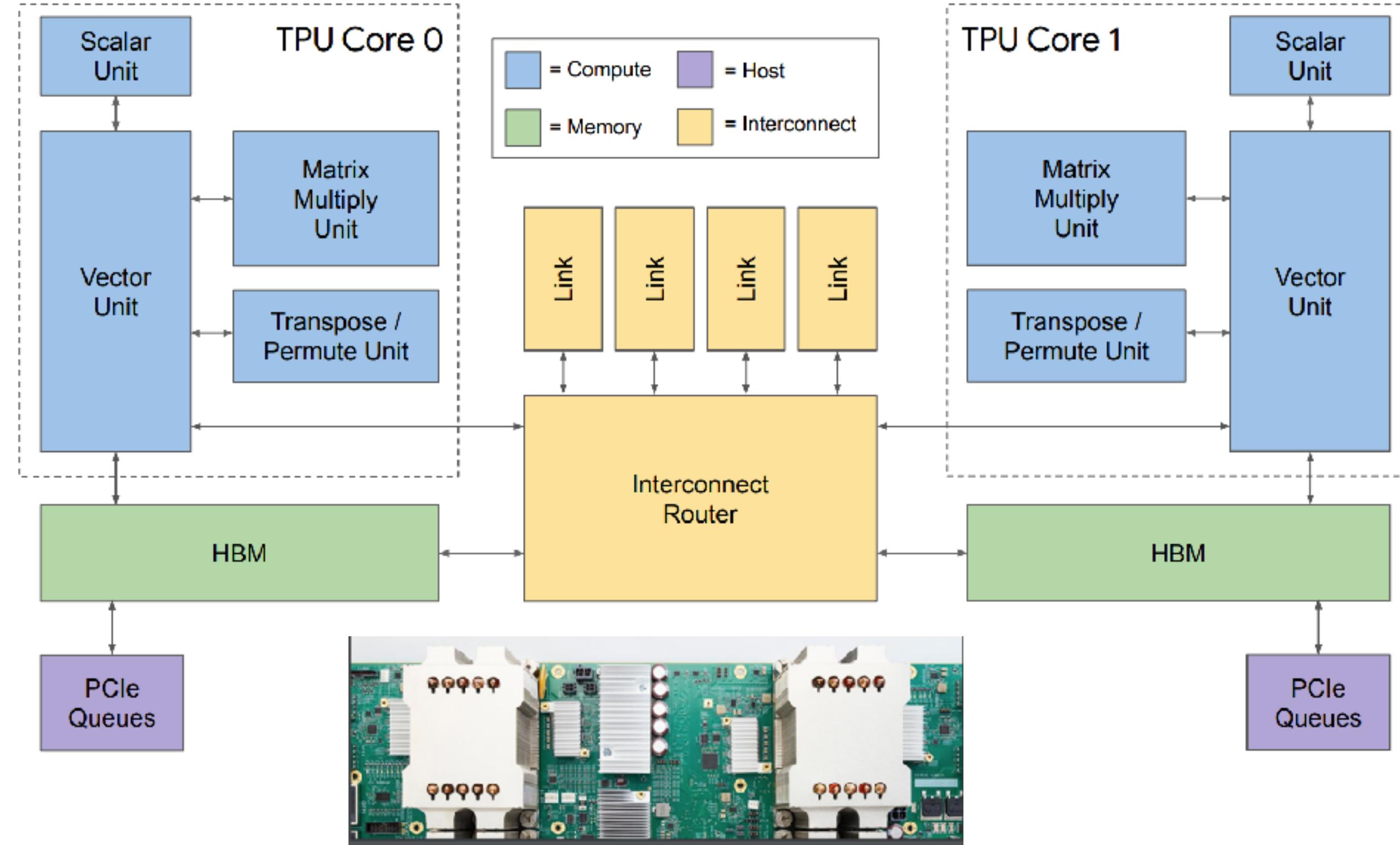


SSD

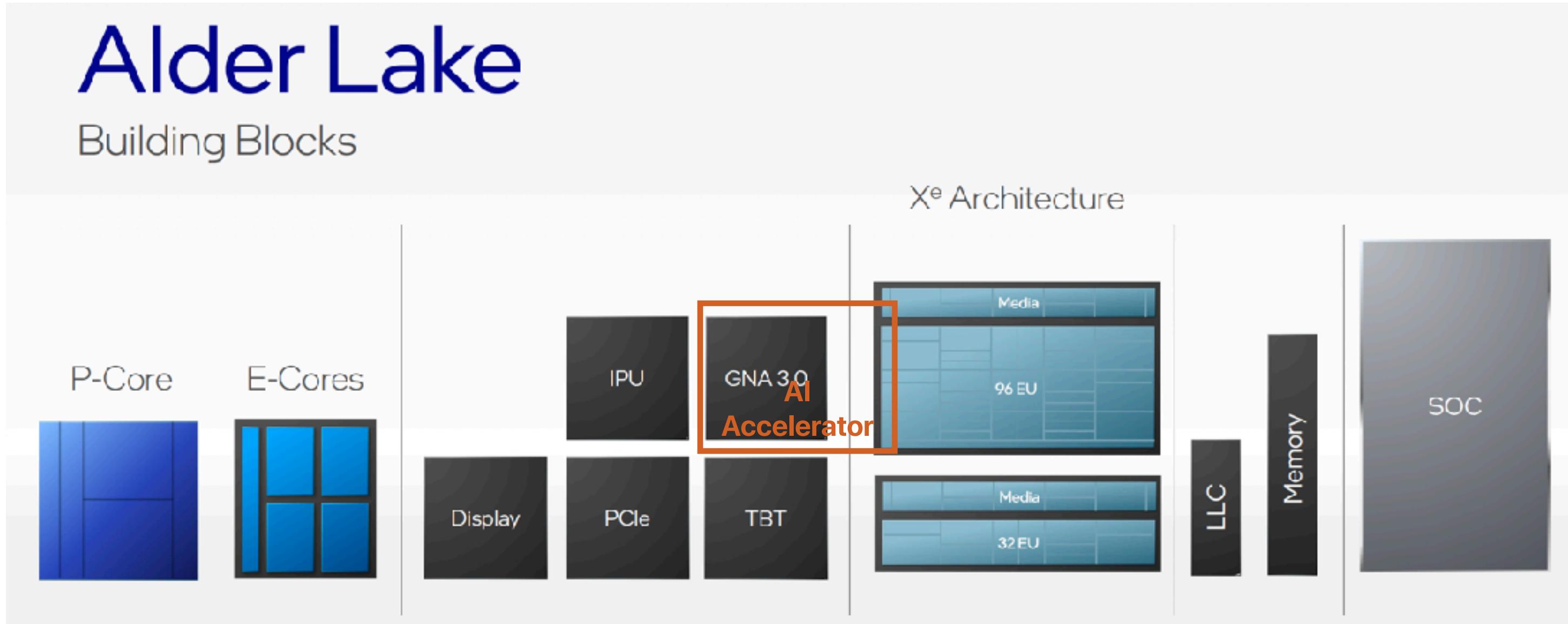


NIC

Tensor Processing Unit



Modern processors also contain “accelerators”



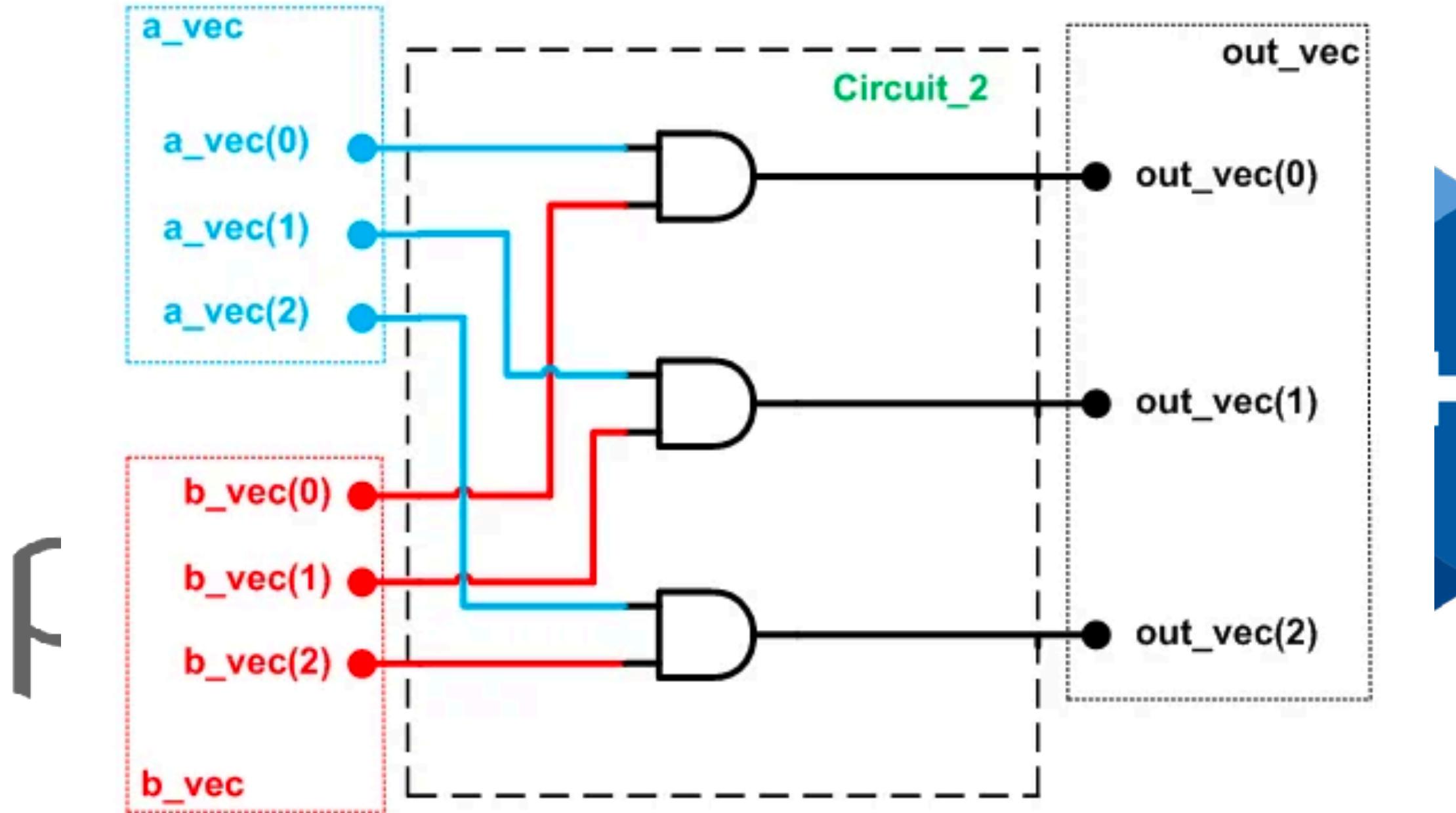


Conclusion: A New *Golden Age*

- End of Dennard Scaling and Moore's Law
⇒ architecture innovation to improve performance/cost/energy
- Security ⇒ architecture innovation too
- Domain Specific Languages ⇒ Domain Specific Architectures
- Free, open architectures and open source implementations
⇒ everyone can innovate and contribute
- Cloud FPGAs ⇒ all can design and deploy custom "HW"
- Agile HW development ⇒ all can afford to make (small) chips
- Like 1980s, great time for architects in academia & in industry!

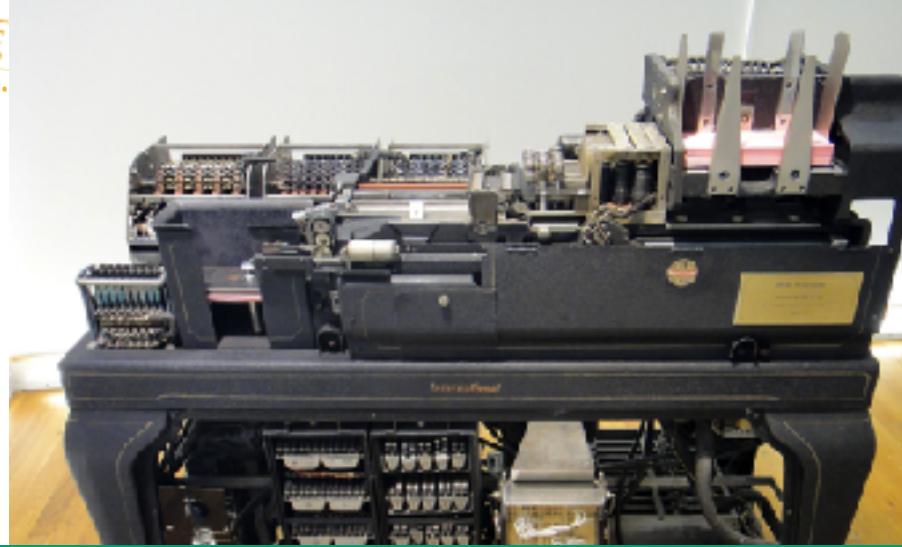


Is building domain specific accelerators the only avenue?

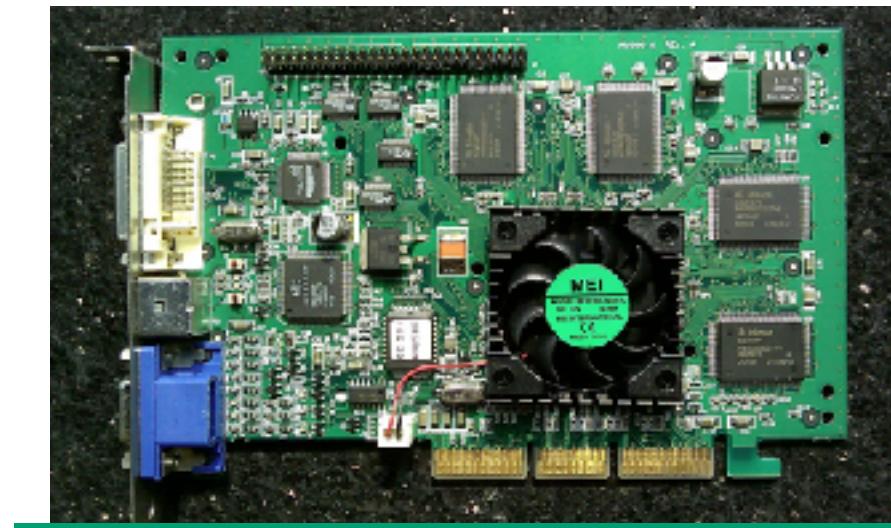




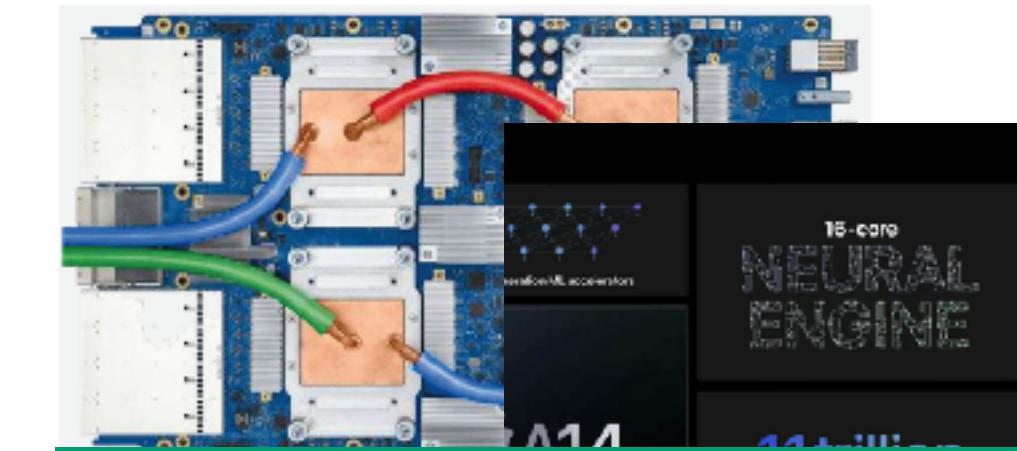
The evolutionary cycle of computing



A combination of special-purposed logical units

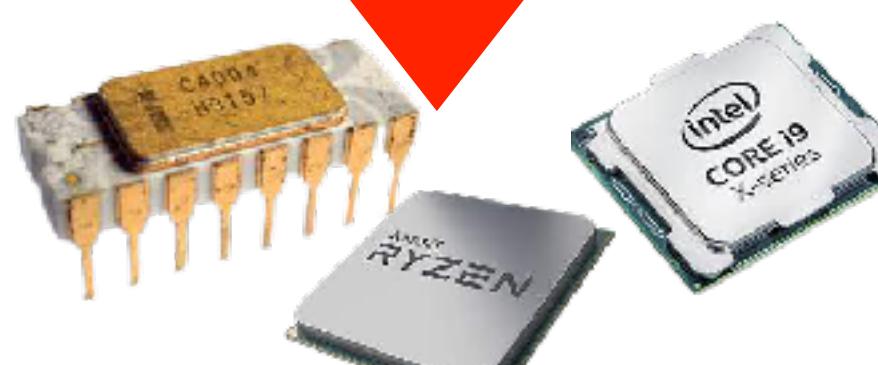


Graphics "Accelerator"



AI/ML/NN/RayTracing
"Accelerators"

Special-purpose
General-purpose



General-purpose microprocessors



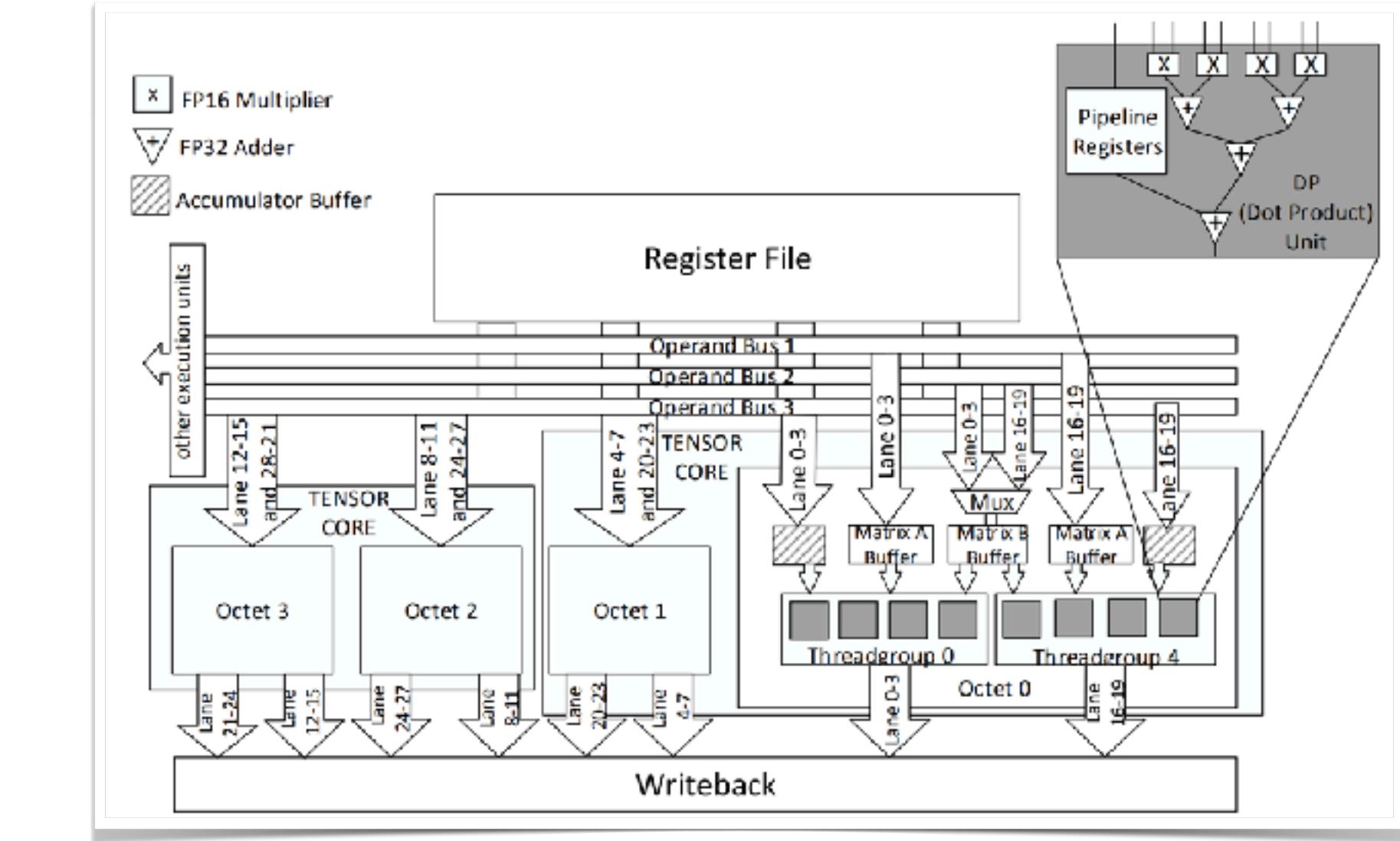
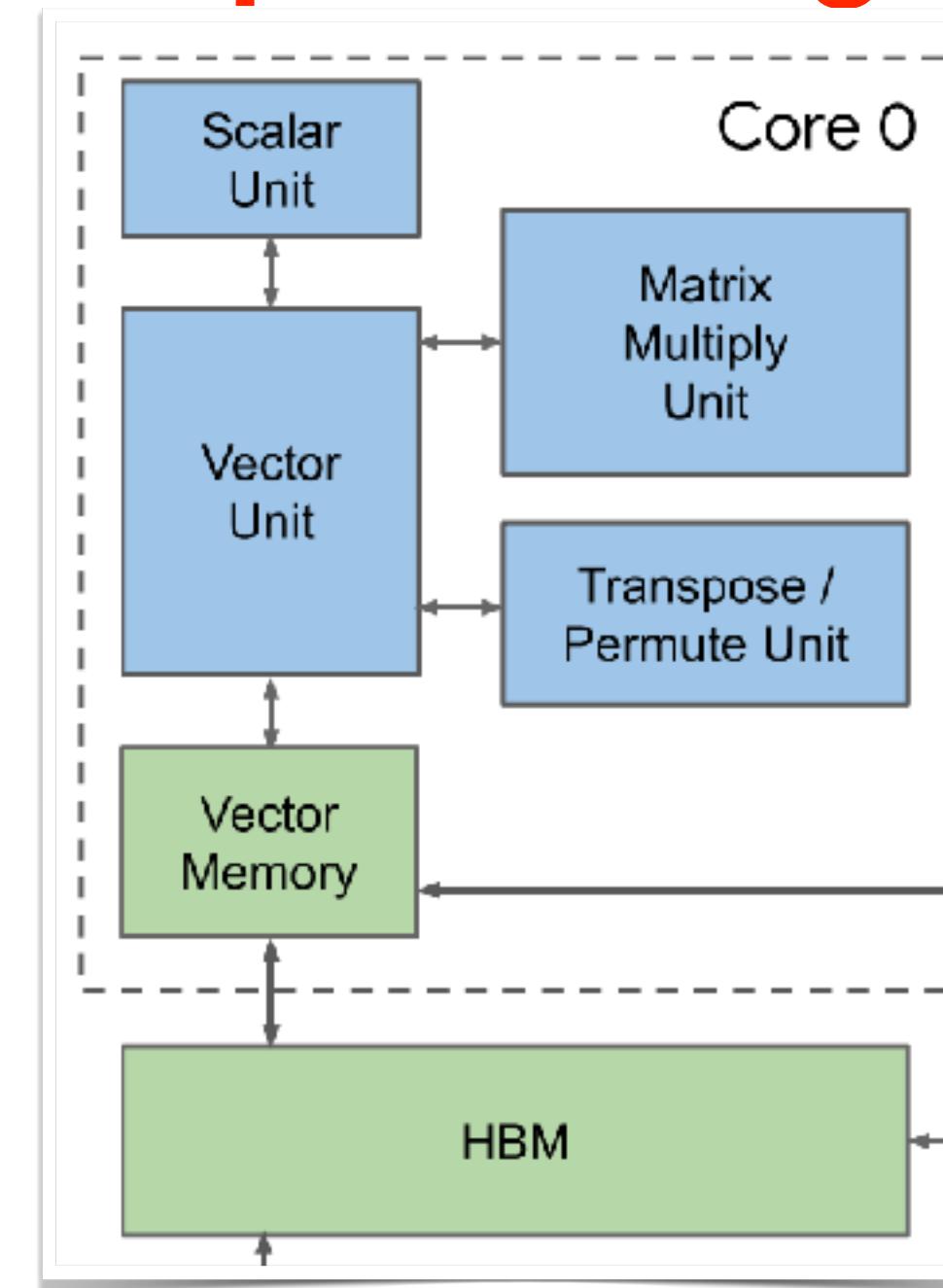
General Purpose Computing On GPUs

General Purpose Computing On Modern Accelerators ?





Matrix processing — the core of AI/ML accelerators



T. Norrie et al., "The Design Process for Google's Training Chips: TPUv2 and TPUv3," in IEEE Micro, vol. 41, no. 2, pp. 56-63

M. Raihan, N. Goli and T. Aamodt, "Modeling Deep Learning Accelerator Enabled GPUs, ISPASS 2019



Who is my most loyal customer?

```
SELECT A.customer, B.brand, SUM(A.Quantity * B.Price) AS value FROM A INNER JOIN B WHERE ON  
A.ProductID = B.ProductID GROUP BY A.customer, B.brand;
```

Customer	ProductID	Brand	Quant.
Abe	i9-12900	Intel	1
Abe	i7-12700	Intel	1
Abe	RTX3080	NVIDIA	1
Abe	660p	Intel	2
Bob	RyZen 5800G	AMD	1
Bob	980Pro	Samsung	1
Cindy	RTX3080	NVIDIA	1
Cindy	i7-12700	Intel	2
Cindy	660p	Intel	2
Diana	RyZen 5800G	AMD	1
Diana	RX5000	AMD	1
Diana	980Pro	Samsung	1

B

Brand	ProductID	Price
Intel	i9-12900	600
Intel	i7-12700	400
Intel	660p	100
AMD	RX5000	500
AMD	RyZen 5800G	200
NVIDIA	RTX3080	1200
Samsung	980Pro	100

Customer	Intel	AMD	NVIDIA	Samsung
Abe	1200	0	1200	0
Bob	0	400	0	100
Cindy	1000	0	1200	0
Diana	0	900	0	100

Matrix multiplications

Customer/ ProductID	i9-1290	i7-1270	660	RX500	RyZen 5800G	RTX308	980Pr o
	0	0	p	0	0	0	0
Abe	1	1	2	0	0	1	0
Bob	0	0	0	0	1	0	1
Cindy	0	2	2	0	0	1	0
Diana	0	0	0	1	1	0	1

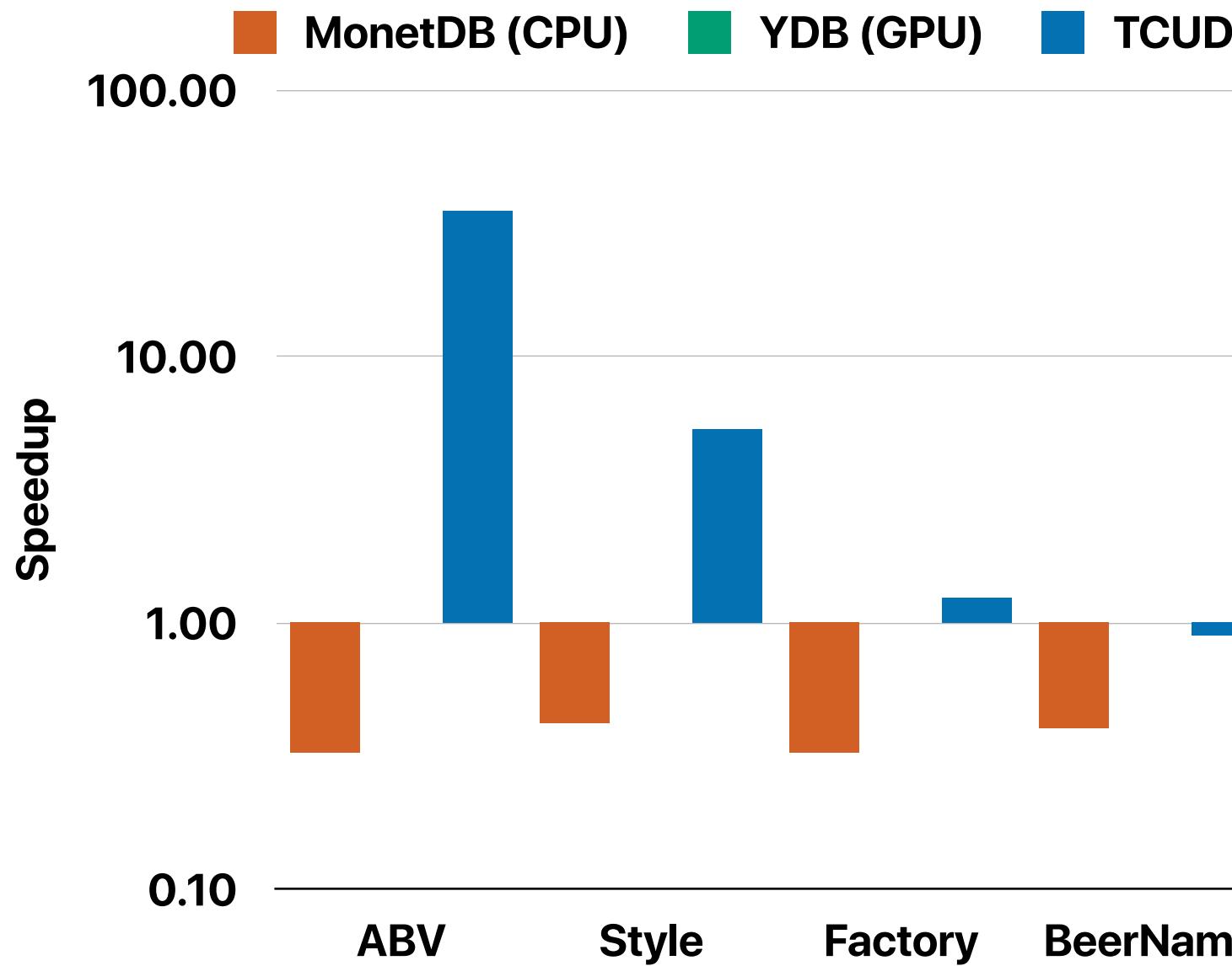
ProductID/ Brand	Intel	AMD	NVIDIA	Samsung
i9-12900	600	0	0	0
i7-12700	400	0	0	0
660p	100	0	0	0
RX5000	0	500	0	0
RyZen 5800G	0	400	0	0
RTX3080	0	0	1200	0
980Pro	0	0	0	100

Memory copy

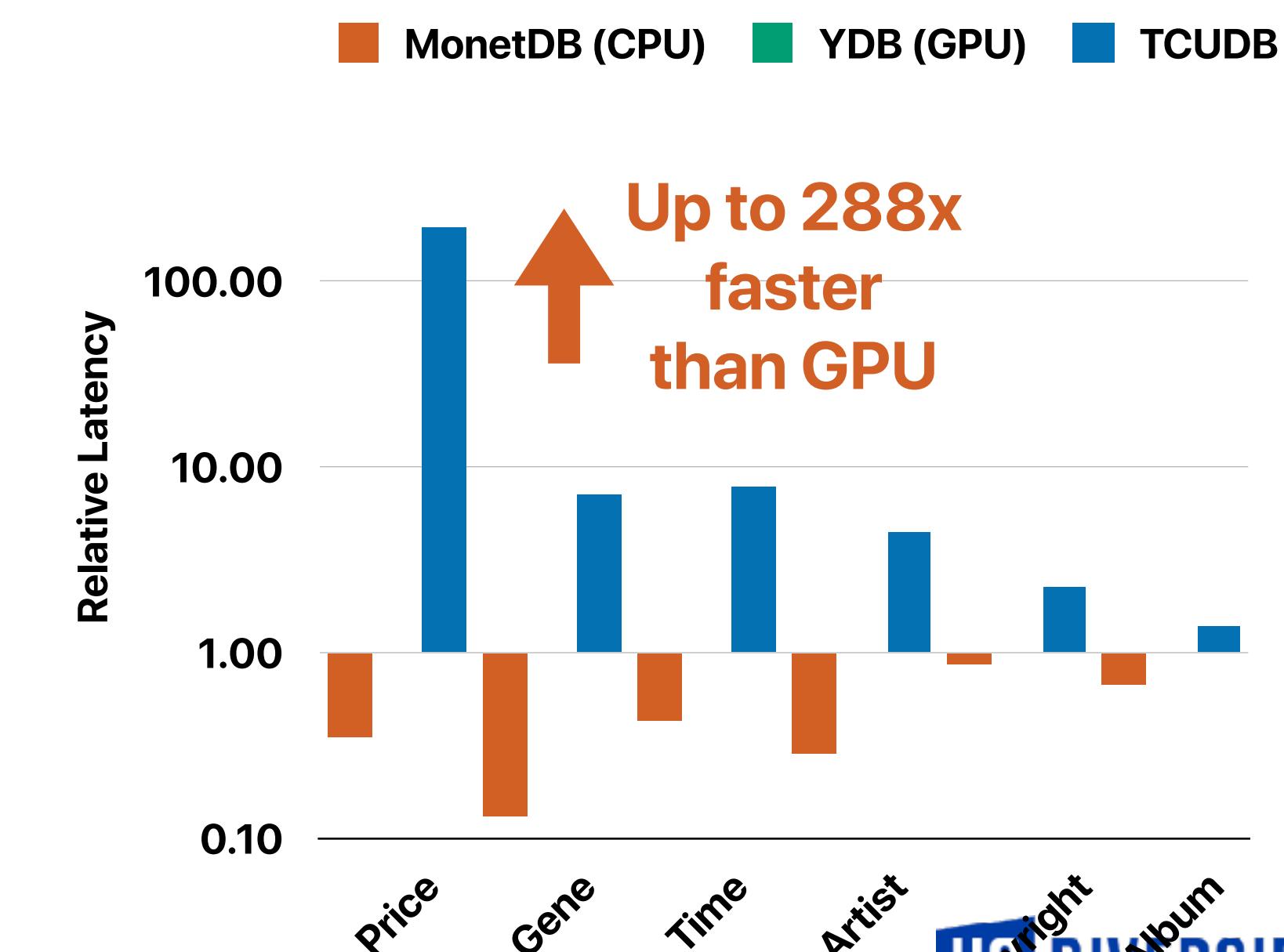


Entity matching

- `SELECT TABLE_A.ID, TABLE_A.BEER_NAME, TABLE_B.ID, TABLE_B.BEER_NAME FROM TABLE_A, TABLE_B WHERE TABLE_A.ABV = TABLE_B.ABV;`



- `SELECT TABLE_A.ID, TABLE_A.SONG, TABLE_B.ID, TABLE_B.SONG FROM TABLE_A, TABLE_B WHERE TABLE_A.ARTIST = TABLE_B.ARTIST;`



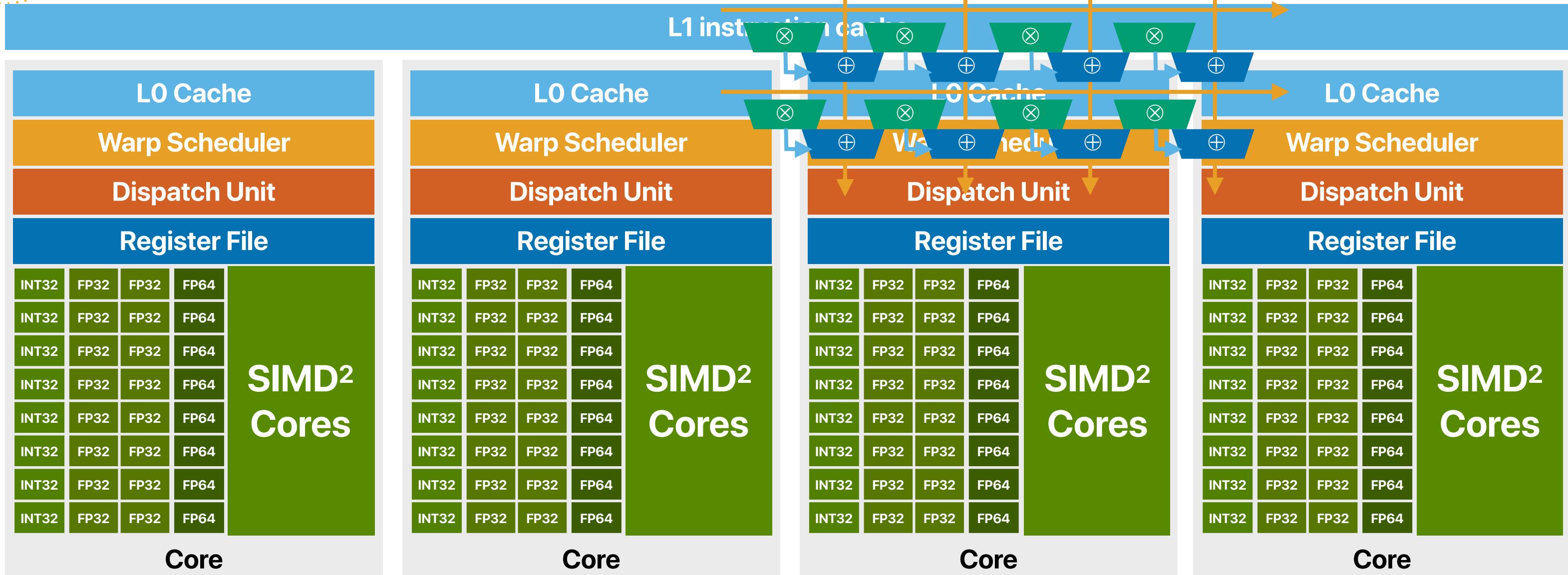


SIMD² instructions

	1st OP	2nd OP
Matrix Multiplications	+	×
All pair shortest path	min	+
Critical path	max	+
Minimum reliability paths	min	×
Maximum reliability paths	max	×
Minimum spanning tree	min	max
Maximum capacity paths	max	min
Transitive and reflexive closure	or	and
L2 Distance	+	$ a-b ^2$

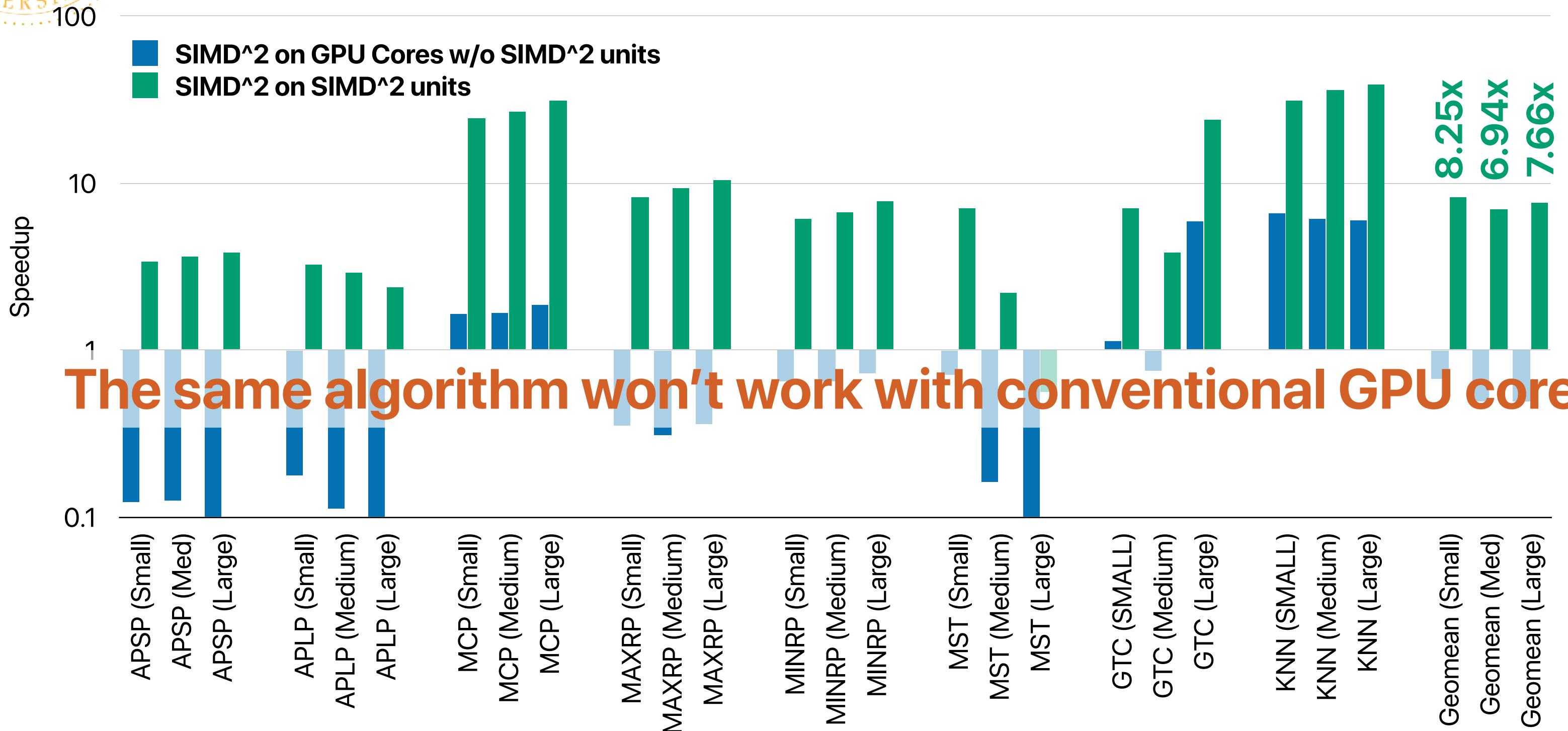


Integration of SIMD2





SIMD² performance



Final words

Conclusion

- Computer architecture is now more important than you could ever imagine
- Being a “programmer” is easy. You need to know architecture a lot to be a “performance programmer”
 - Branch prediction
 - Cache
- Multicore era — to get your multithreaded program correct and perform well, you need to take care of coherence and consistency
- We’re now in the “dark silicon era”
 - Single-core isn’t getting any faster
 - Multi-core doesn’t scale anymore
 - We will see more and more ASICs
 - You need to write more “system-level” programs to use these new ASICs.

One more thing...

Sample Final

Format of the online final

- Multiple choices (15 questions)
 - They're like your clicker/midterm multiple choices questions
 - Cumulative, don't forget your midterm and midterm review
- Free answer and open-ended questions * 2 problem sets, 8 questions in total
 - Explain your answer clearly and "correctly". If you include "incorrect" information in your answer, it's going to hurt your grade
 - May not have a standard answer. You need to understand the concepts to provide a good answer
 - You should review your assignments carefully

Multiple choices

How many dependencies do we have?

- How many pairs of data dependences are there in the following x86 instructions?

```
movl    (%rdi), %eax  
xorl    (%rsi), %eax  
movl    %eax, (%rdi)  
xorl    (%rsi), %eax  
movl    %eax, (%rsi)  
xorl    %eax, (%rdi)
```

- A. 1
- B. 2
- C. 3
- D. 4
- E. 5

The effect of code optimization

- By reordering which pair of the following instruction stream can we eliminate all stalls without affecting the correctness of the code?
 - ① `movl (%rdi), %ecx`
 - ② `addl %ecx, %eax`
 - ③ `addq $4, %rdi`
 - ④ `cmpq %rdx, %rdi`
 - ⑤ `jne .L3`
 - ⑥ `ret`
 - A. (1) & (2)
 - B. (2) & (3)
 - C. (3) & (4)
 - D. (4) & (5)
 - E. None of the pairs can be reordered

CMP advantages

- How many of the following are advantages of CMP over traditional superscalar processor
 - ① CMP can provide better energy-efficiency within the same area
 - ② CMP can deliver better instruction throughput within the same die area (chip size)
 - ③ CMP can achieve better ILP for each running thread
 - ④ CMP can improve the performance of a single-threaded application without modifying code
- A. 0
B. 1
C. 2
D. 3
E. 4

What about “linked list”

- Assume the current PC is already at instruction (1) and this linked list has only three nodes. This processor can fetch and issue 2 instructions per cycle, with exactly the same register renaming hardware and pipeline as we showed previously.

Which of the following C state of the code snippet determines the performance?

- A. do {
- B. number_of_nodes++;
- C. current = current->next;
- D. } while (current != NULL);

Amdahl's Law on Multicore Architectures

- Regarding Amdahl's Law on multicore architectures, how many of the following statements is/are correct?
 - ① If we have unlimited parallelism, the performance of each parallel piece does not matter as long as the performance slowdown in each piece is bounded
 - ② With unlimited amount of parallel hardware units, single-core performance does not matter anymore
 - ③ With unlimited amount of parallel hardware units, the maximum speedup will be bounded by the fraction of parallel parts
 - ④ With unlimited amount of parallel hardware units, the effect of scheduling and data exchange overhead is minor

A. 0
B. 1
C. 2
D. 3
E. 4

Virtual indexed, physical tagged cache limits the cache size

- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the system use 4K pages.
 - A. 32B blocks, 2-way
 - B. 32B blocks, 4-way
 - C. 64B blocks, 4-way
 - D. 64B blocks, 8-way

Power & Energy

- Regarding dynamic frequency scaling, how many of the following statements are correct?
 - ① Lowering the frequency helps extending the battery life
 - ② Lowering the frequency helps reducing the heat generation
 - ③ Lowering the frequency helps reducing the electricity bill
 - ④ A CPU operating at 25% of the peak frequency can still consume more than 50% of the peak power

A. 0

B. 1

C. 2

D. 3

E. 4

Why is D better than C?

- How many of the following statements explains the main reason why B outperforms C with compiler optimizations
 - D has lower dynamic instruction count than C
 - D has significantly lower branch mis-prediction rate than C
 - D has significantly fewer branch instructions than C
 - D can incur fewer memory accesses than C

A. 0

```
inline int __popcount(uint64_t x) {  
    int c = 0;  
    int table[16] = {0, 1, 1, 2, 1,  
                    2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4};  
    while(x) {  
        c += table[(x & 0xF)];  
        x = x >> 4;  
    }  
    return c;  
}
```

B. 1

C. 2

D. 3

E. 4



```
inline int __popcount(uint64_t x) {  
    int c = 0;  
    int table[16] = {0, 1, 1, 2, 1,  
                    2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4};  
    for (uint64_t i = 0; i < 16; i++) {  
        c += table[(x & 0xF)];  
        x = x >> 4;  
    }  
    return c;  
}
```

Demo revisited

- Why the performance is better when option is not “0”
 - ① The amount of dynamic instructions needs to execute is a lot smaller
 - ② The amount of branch instructions to execute is smaller
 - ③ The amount of branch mis-predictions is smaller
 - ④ The amount of data accesses is smaller
- A. 0 **if(option)**
 std::sort(data, data + arraySize);
- B. 1 **for (unsigned i = 0; i < 100000; ++i) {**
- C. 2 int **threshold** = std::rand();
- D. 3 **for (unsigned i = 0; i < arraySize; ++i) {**
 if (data[i] >= threshold)
 sum ++;
- E. 4 **}**
 }

Why can't we proceed without stalls/no-ops?

- How many of the following statements are true regarding why we have to stall for each branch in the current pipeline processor
 - ① The target address when branch is taken is not available for instruction fetch stage of the next cycle
 - ② The target address when branch is not-taken is not available for instruction fetch stage of the next cycle
 - ③ The branch outcome cannot be decided until the comparison result of ALU is not out
 - ④ The next instruction needs the branch instruction to write back its result
- A. 0
B. 1
C. 2
D. 3
E. 4

What if the code look like this?

- D-L1 Cache configuration of NVIDIA Tegra X1
 - Size 64KB, 4-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++)
    c[i] = a[i]; //load a and then store to c
for(i = 0; i < 512; i++)
    c[i] += b[i]; //load b, load c, add, and then store to c
```

What's the data cache miss rate for this code?

- A. 6.25%
- B. 56.25%
- C. 66.67%
- D. 68.75%
- E. 100%

What kind(s) of misses can matrix transpose remove?

- By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Block + Transpose

```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```

What data structure is performing better

	Array of objects	object of arrays
	<pre>struct grades { int id; double *homework; double average; };</pre>	<pre>struct grades { int *id; double **homework; double *average; };</pre>
average of each homework	<pre>for(i=0;i<homework_items; i++) { gradesheet[total_number_students].homework[i] = 0.0; for(j=0;j<total_number_students;j++) gradesheet[total_number_students].homework[i] +=gradesheet[j].homework[i]; gradesheet[total_number_students].homework[i] /= (double)total_number_students; }</pre>	<pre>for(i = 0;i < homework_items; i++) { gradesheet.homework[i][total_number_students] = 0.0; for(j = 0; j <total_number_students;j++) { gradesheet.homework[i][total_number_students] += gradesheet.homework[i][j]; } gradesheet.homework[i][total_number_students] /= total_number_students; }</pre>

- Considering your workload would like to calculate the average score of **one of the homework for all students**, which data structure would deliver better performance?
 - Array of objects
 - Object of arrays

Cache coherency

- Assuming that we are running the following code on a CMP with a cache coherency protocol, how many of the following outputs are possible? (a is initialized to 0 as assume we will output more than 10 numbers)

thread 1	thread 2
while(1) printf("%d ", a);	while(1) a++;

- ① 0123456789
- ② 1259368101213
- ③ 1111111164100
- ④ 111111111100
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

Performance comparison

- Comparing implementations of thread_vadd — L and R, please identify which one will be performing better and why

Version L

```
void *threaded_vadd(void *thread_id)
{
    int tid = *(int *)thread_id;
    int i;
    for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)
    {
        c[i] = a[i] + b[i];
    }
    return NULL;
}
```

Version R

```
void *threaded_vadd(void *thread_id)
{
    int tid = *(int *)thread_id;
    int i;
    for(i=tid*(ARRAY_SIZE/NUM_OF_THREADS);i<(tid+1)*(ARRAY_SIZE/NUM_OF_THREADS);i++)
    {
        c[i] = a[i] + b[i];
    }
    return NULL;
}
```

- L is better, because the cache miss rate is lower
- R is better, because the cache miss rate is lower
- L is better, because the instruction count is lower
- R is better, because the instruction count is lower
- Both are about the same

FalseSharing

Main thread

```
for(i = 0 ; i < NUM_OF_THREADS ; i++)
{
    tids[i] = i;
    pthread_create(&thread[i], NULL, threaded_vadd, &tids);
}
for(i = 0 ; i < NUM_OF_THREADS ; i++)
    pthread_join(thread[i], NULL);
```

Free-answer questions

Branch prediction performance

```
i = 0;  
do {  
    if( i % 2 != 0) // Branch X, taken if i % 2 == 0  
        a[i] *= 2;  
    a[i] += i;  
} while ( ++i < 100)// Branch Y
```

- What's the branch prediction accuracy if the processor uses?
 - A 2-bit local predictor
 - A (4, 2) global history predictor
- Can you rewrite the code to generate the same result but eliminate branch X?

Register renaming

- For the following C code:

```
do {  
    number_of_nodes++;  
    current = current->next;  
} while ( current != NULL );
```

- Assume we have a dual-fetch, dual-issue, out-of-order pipeline where
 - INT ALU takes 1 cycle
 - MEM pipeline: 4 cycles in total
 - BR takes 1 cycle to resolve
- If the loop is taken twice, how many cycles it takes to issue all instructions?
- If the loop is taken 100 times, what's the average CPI?

Cache simulation

- The processor has a 8KB, 256B blocked, 2-way L1 cache. Consider the following code:

```
for(i=0;i<256;i++) {  
    a[i] = b[i] + c[i];  
    // load a[i] and load b[i], store to c[i]  
    // &a[0] = 0x10000, &b[0] = 0x20000, &c[0] = 0x30000  
}
```

- What's the total miss rate? How many of the misses are compulsory misses? How many of the misses are conflict misses?
- How can you improve the cache performance of the above code through changing hardware?
- How can you improve the performance **without** changing hardware?

Algorithm and performance

- Consider the LCS algorithm on the right hand-side.
 - What kind of operation would determine the execution time?
 - What would be the critical path of execution?
 - Can you parallelize the algorithm?

```
int max(int a, int b);  
  
/* Returns length of LCS for X[0..m-1], Y[0..n-1] */  
int lcs( char *X, char *Y, int m, int n )  
{  
    if (m == 0 || n == 0)  
        return 0;  
    if (X[m-1] == Y[n-1])  
        return 1 + lcs(X, Y, m-1, n-1);  
    else  
        return max(lcs(X, Y, m, n-1), lcs(X, Y, m-1, n));  
}  
  
/* Utility function to get max of 2 integers */  
int max(int a, int b)  
{  
    return (a > b)? a : b;  
}
```

Code and cache miss rate

- Assume my cache has 16KB capacity, 16 byte block size and is 2-way set associative. Integers are 4 bytes. Give the C code for a loop that has a very poor hit rate in this cache but whose hit rate raises to almost 100% if we double the capacity to 32KB.

Open-ended questions

SMT v.s. CMP

- Both CMP & SMT exploit thread-level or task-level parallelism. Assuming both application X and application Y have similar instruction combination, say 60% ALU, 20% load/store, and 20% branches. Consider two processors:

P1: CMP with a 2-issue pipeline on each core. Each core has a private L1 32KB D-cache

P2: SMT with a 4-issue pipeline. 64KB L1 D-cache

Which one do you think is better?

- A. P1
- B. P2

Other open-ended questions

- If you're asked to design a machine learning hardware, what will you do?
- If you're asked to build an Xeon Phi type processor where each core also has many-way SMT, are you going to give the processor more cache or better branch predictor?
- Can we focus on improving the throughput of computing instead of latency? Can you give an example on what type of applications will not work well in this way
- Pros and cons for branch prediction using perceptrons?

Other open-ended questions (cont.)

- What's Amdahl's Law implication on parallelism? How does that guide future software design?
- What's Dark Silicon Problem? What are the new design trends in addressing the problem?
- If the OoO pipeline is highly optimized, do we still care about the ISA design?
- What's volatile? What's inline? Why we need them? The pros and cons?



Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

Announcement

- Assignment #5 due this Thursday
- Final review tomorrow
 - You are encouraged to skim through sample midterm before you come
 - We will take a group photo (prepare your smile on camera!)
- Final Exam
 - The final exam will be held at 9/8 3p-6p
 - Similar to the midterm, but more time and about 1.5x—2x longer
- Last office hours
 - Hung-Wei Tseng: 9/6 8:30p-9:30p (online), 9/7 3:30p-5:30p (in-person)

Computer Science & Engineering

142

つづく

