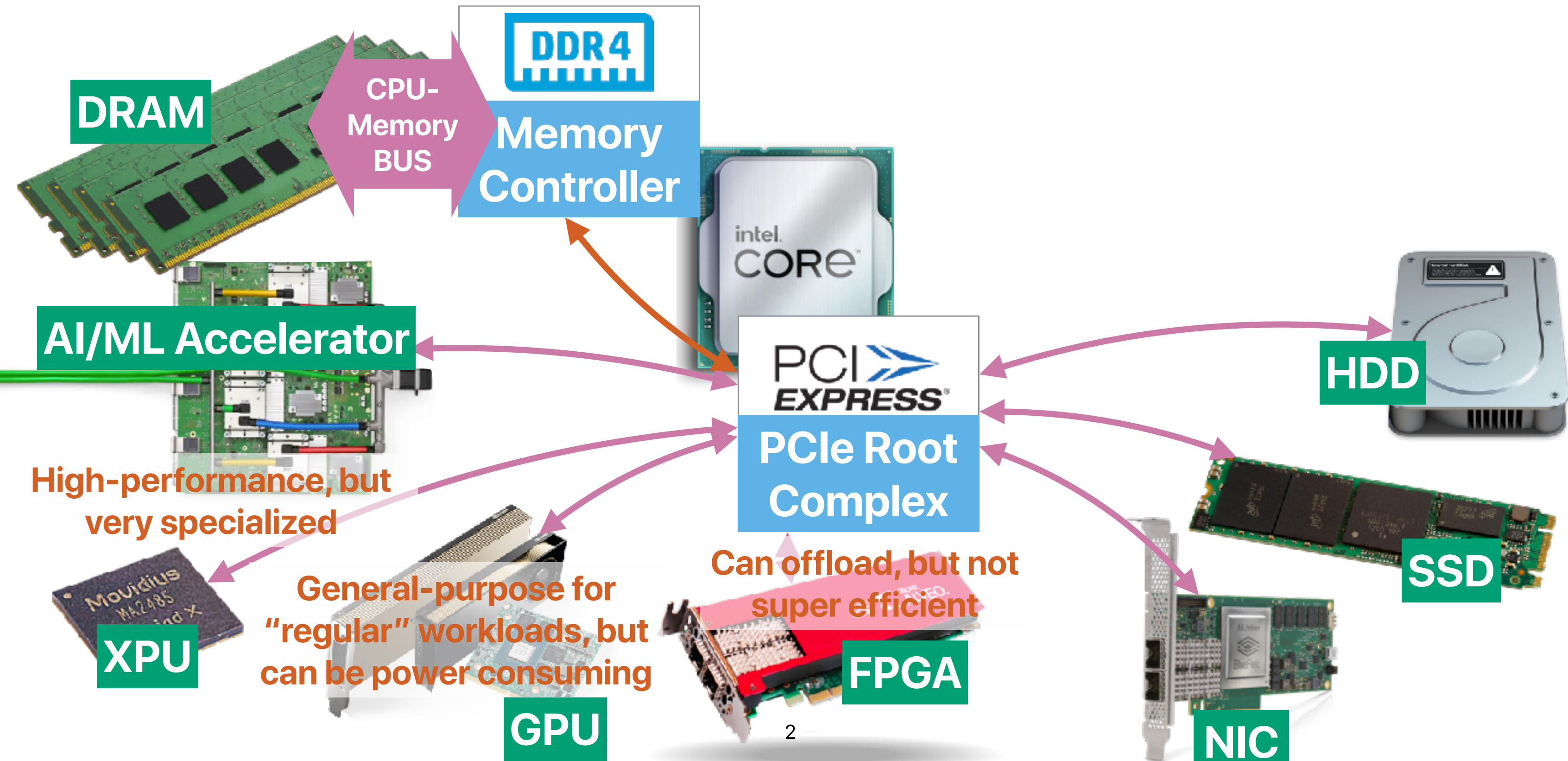


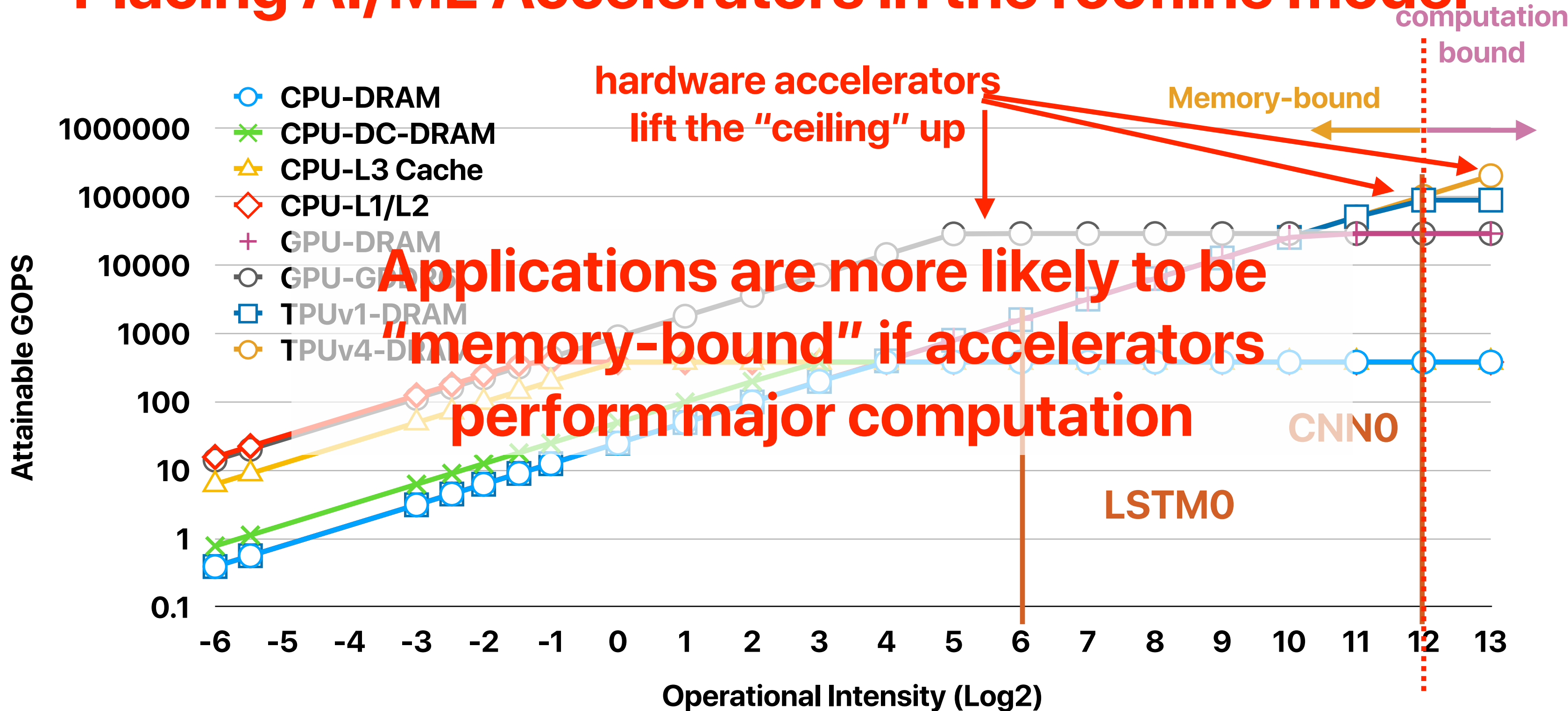
Memory subsystem (2) & In/Near Memory Processing

Hung-Wei Tseng

Recap: The landscape of modern computers



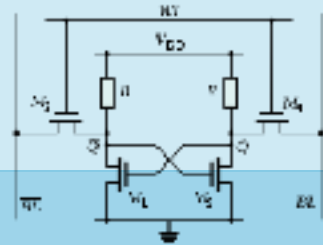
Placing AI/ML Accelerators in the roofline model



Recap: Memory technologies we have today

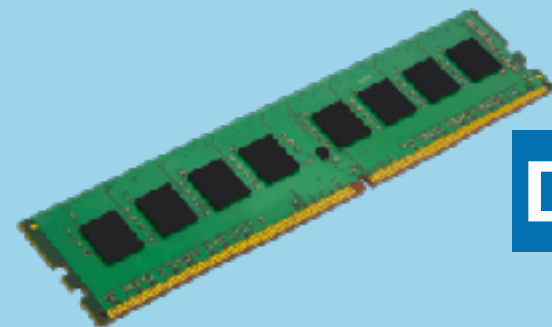
Volatile Memory

100ps



SRAM

ns

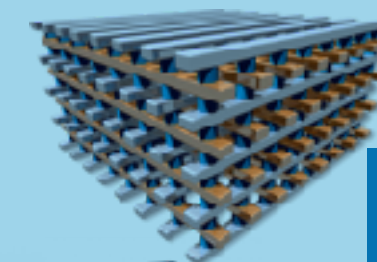


DRAM

us

ms

Non-Volatile Memory



RRAM



PCM

3D XPoint



Flash memory



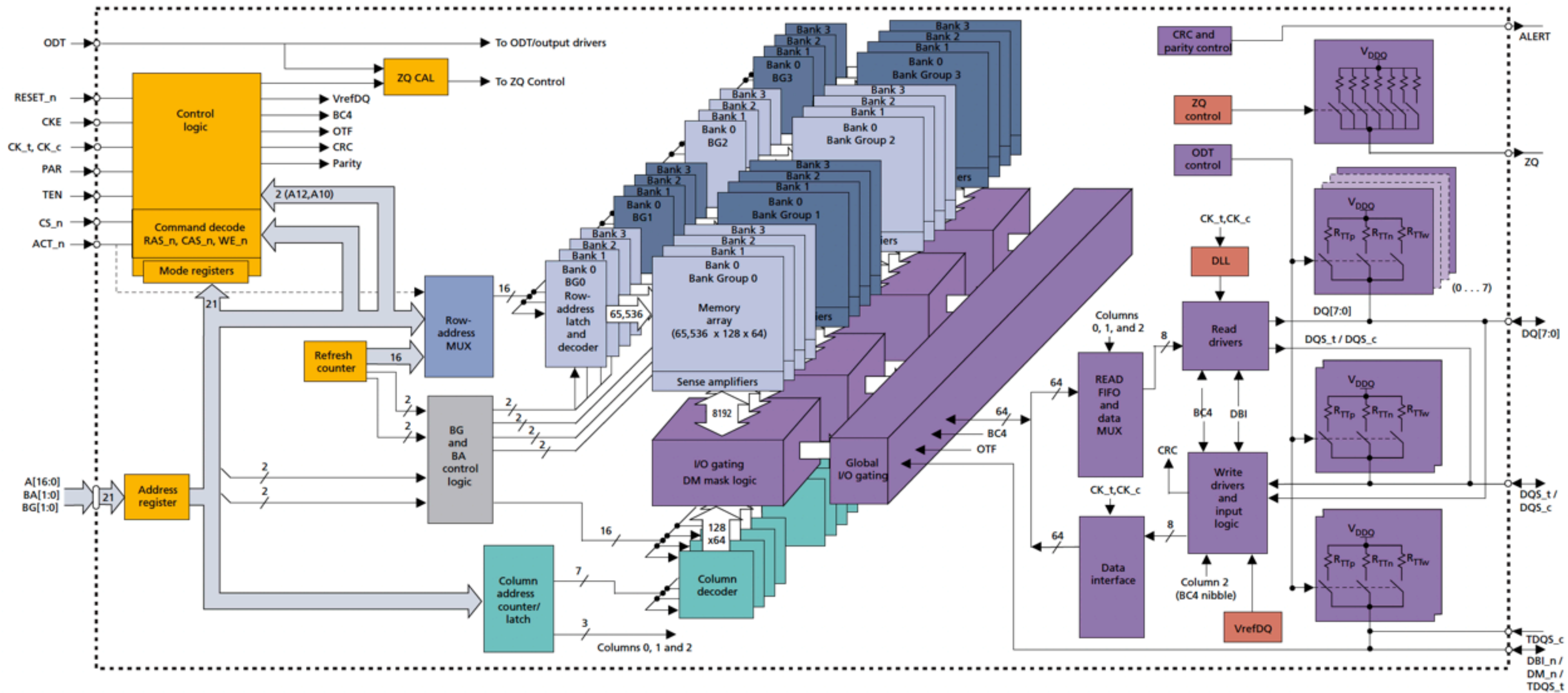
Hard Disk Drives

Recap: Improving memory-bounded applications

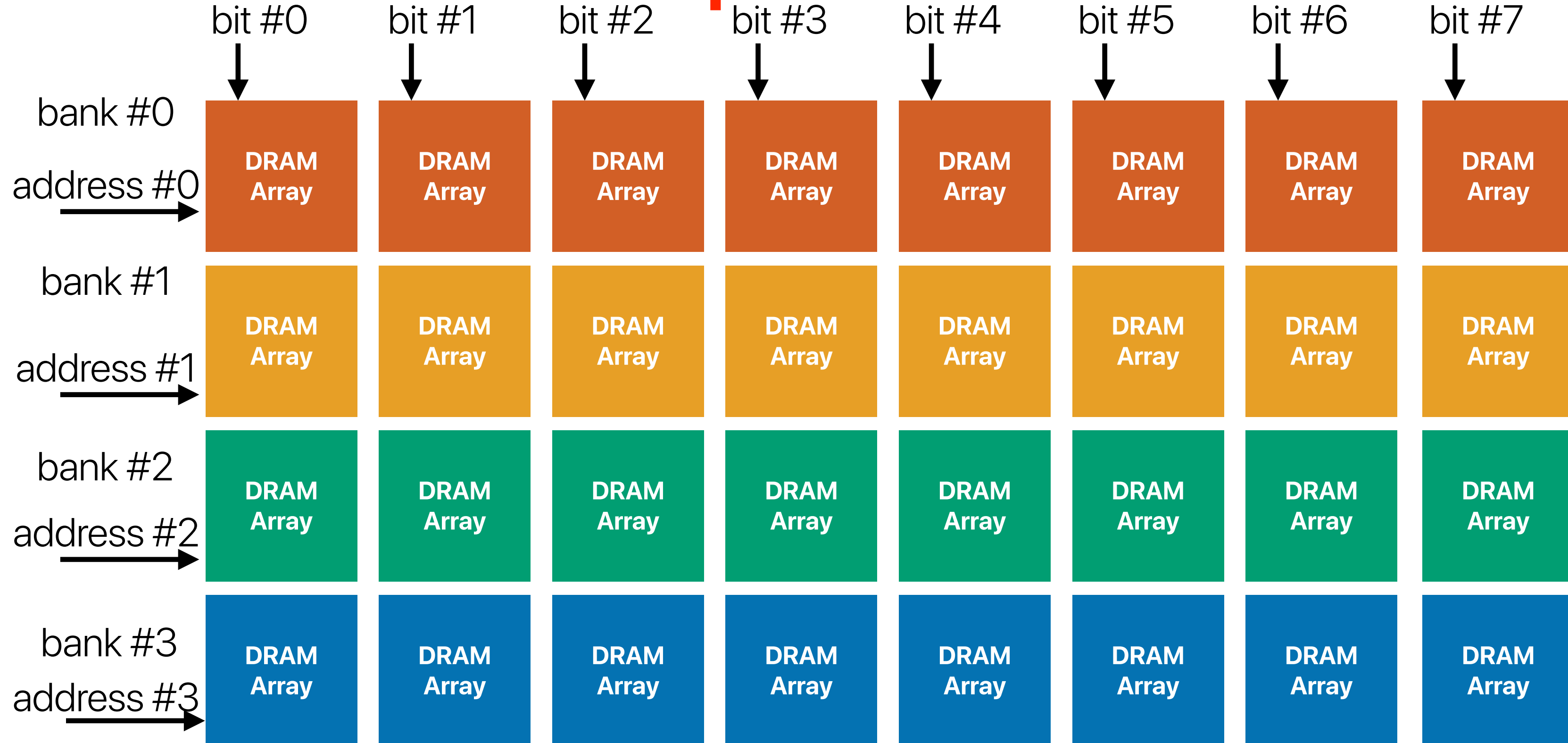
- Faster memory technologies
- Higher memory bandwidth
- Lower data volume

Outline

- The Main Memory subsystem (cont.)
- Performance of main memory subsystem
- Near/In-memory processing

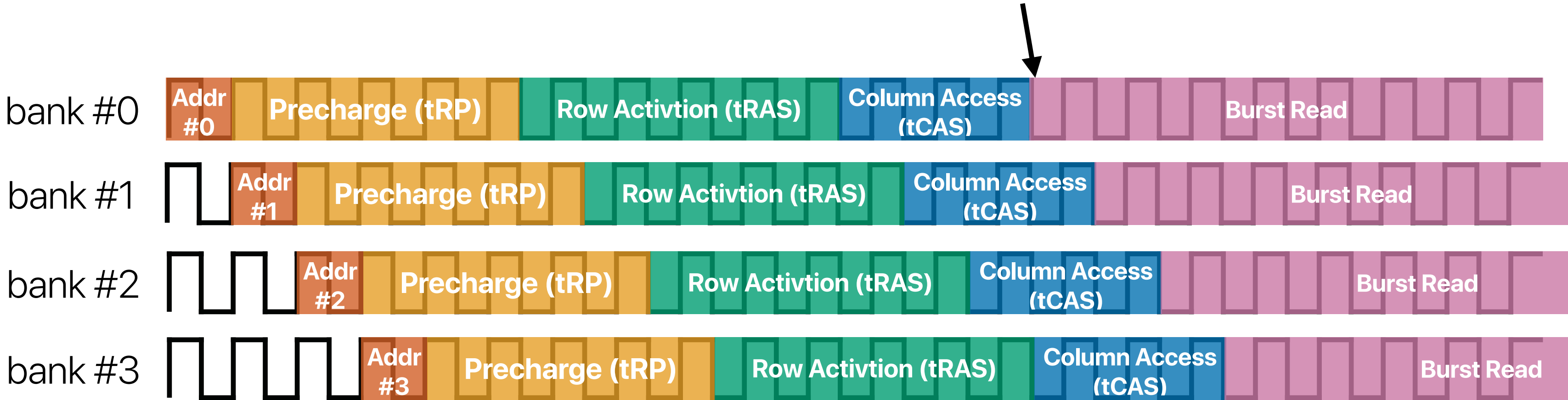


Multiple Banks



Multi-bank access

we can start output a "byte" from every 8 chips each cycle after this

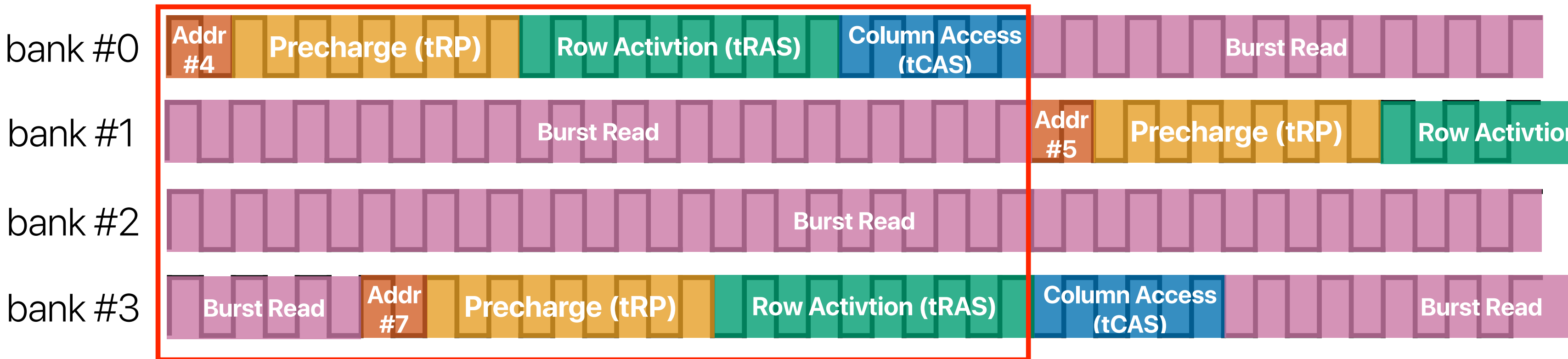


only one bank can accept request each cycle

the memory bandwidth
can be fully utilized after
this

Multi-bank access

The latency of pre-charge, row/column accesses is fully covered!



DRAM Performance

- Latency per "8-bit" — 0.75 ns (if it's row-buffered)

- Bandwidth per die = $\frac{1}{0.75ns} = 1.33GB/sec$
- 16 chips = $16 \times \frac{1}{0.75ns} = 21.33GB/sec$

2. Key Features

[Table 2] 8Gb DDR4 C-die Speed bins

Speed	DDR4-1600	DDR4-1866	DDR4-2133	DDR4-2400	DDR4-2666	Unit
	11-11-11	13-13-13	15-15-15	17-17-17	19-19-19	
tCK(min)	1.25	1.071	0.937	0.833	0.75	ns
CAS Latency	11	13	15	17	19	nCK
tRCD(min)	13.75	13.92	14.08	14.18	14.25	ns
tRP(min)	13.75	13.92	14.08	14.18	14.25	ns
IRAS(min)	35	34	33	32	32	ns
tRC(min)	48.75	47.92	47.06	46.16	46.25	ns

- JEDEC standard 1.2V (1.14V~1.26V)
- V_{DDQ} = 1.2V (1.14V~1.26V)
- V_{PP} = 2.5V (2.375V~2.75V)
- 800 MHz f_{CK} for 1600Mb/sec/pin, 933 MHz f_{CK} for 1866Mb/sec/pin, 1067MHz f_{CK} for 2133Mb/sec/pin, 1200MHz f_{CK} for 2400Mb/sec/pin, 1333MHz f_{CK} for 2666Mb/sec/pin
- 8 Banks (2 Bank Groups)
- Programmable CAS Latency (posted CAS): 10,11,12,13,14,15,16,17,18,19,20
- Programmable CAS Write Latency (CWL) = 9,11 (DDR4-1600), 10,12 (DDR4-1866), 11,14 (DDR4-2133), 12,16 (DDR4-2400) and 14,18 (DDR4-2666)
- 8-bit pre-fetch
- Burst Length: 8, 4 with tCCD = 4 which does not allow seamless read or write [either On the fly using A12 or MRS]
- Bi-directional Differential Data-Strobe
- Internal (self) calibration: Internal self calibration through ZQ pin (RZQ: 240 ohm ± 1%)
- On Die Termination using ODT pin
- Average Refresh Period 7.8us at lower than T_{CASE} 85°C, 3.9us at 85°C < T_{CASE} ≤ 95 °C
- Connectivity Test Mode (TEN) is Supported

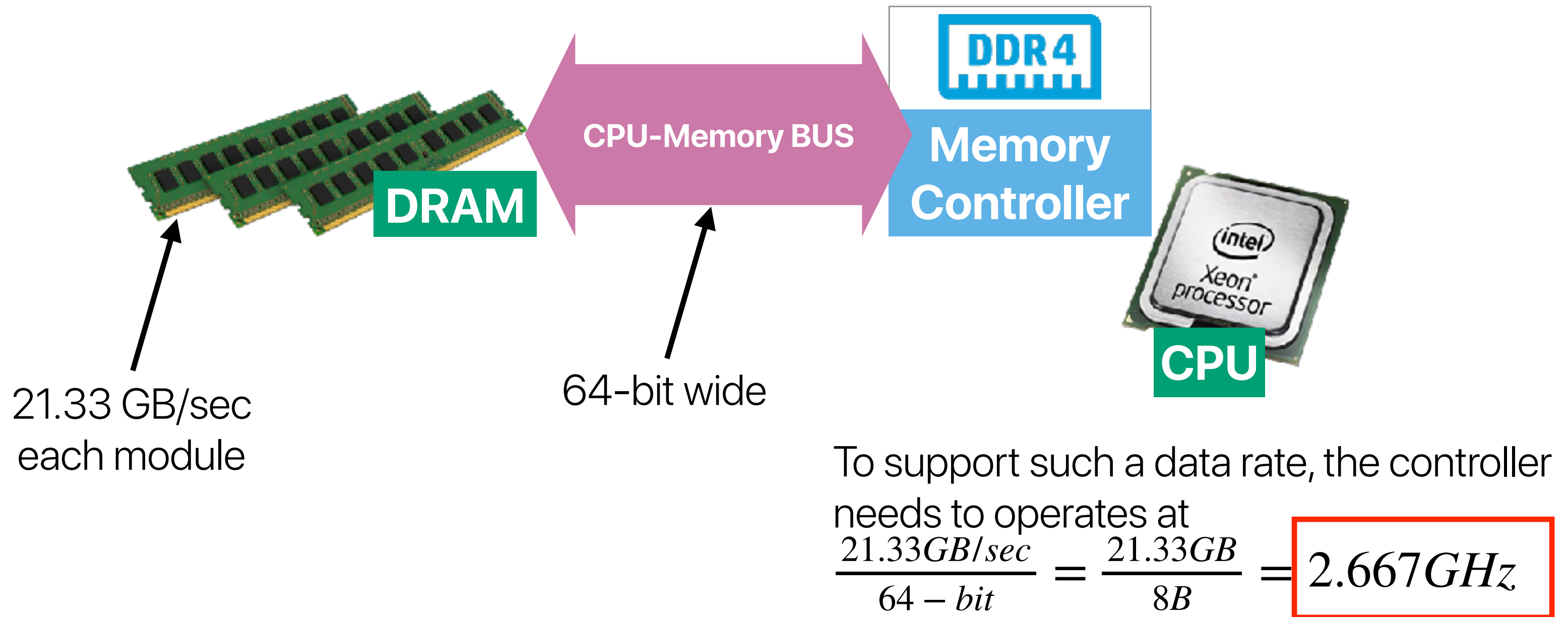
The 8Gb DDR4 SDRAM C-die is organized as a 84Mbit x 16 I/Os x 8banks device. This synchronous device achieves high speed double-data-rate transfer rates of up to 2666Mb/sec/pin (DDR4-2666) for general applications.

The chip is designed to comply with the following key DDR4 SDRAM features such as posted CAS, Programmable CWL, Internal (Self) Calibration, On Die Termination using ODT pin and Asynchronous Reset.

All of the control and address inputs are synchronized with a pair of externally supplied differential clocks. Inputs are latched at the crosspoint of differential clocks (CK rising and CK falling). All I/Os are synchronized with a pair of bidirectional strobes (DQS and DQS) in a source synchronous fashion. The address bus is used to convey row, column, and bank address information in a RAS/CAS multiplexing style. The DDR4 device operates with a single 1.2V (1.14V~1.26V) power supply, 1.2V(1.14V~1.26V) V_{DDQ} and 2.5V (2.375V~2.75V) V_{PP}.

The 8Gb DDR4 C-die device is available in 96ball FBGAs(x16).

How fast is the memory controller frequency?



**What's the bandwidth of
DDR5-5600?**

DDR5-5600

To support such a data rate, the controller needs to operate at

$$\frac{x \text{ GB/sec}}{64 - \text{bit}} = \frac{x \text{ GB/sec}}{8B} = 5.6\text{GHz}$$

$$x = 8 \times 5.6\text{GHz} = 44.8\text{GB/sec}$$

Memory bandwidth

**How can you increase the DRAM
bandwidth?**

Ideas of increasing bandwidth

Ideas of increasing bandwidth

- More parallel bits
- More banks
- More channels
- Widen the memory-processor bus

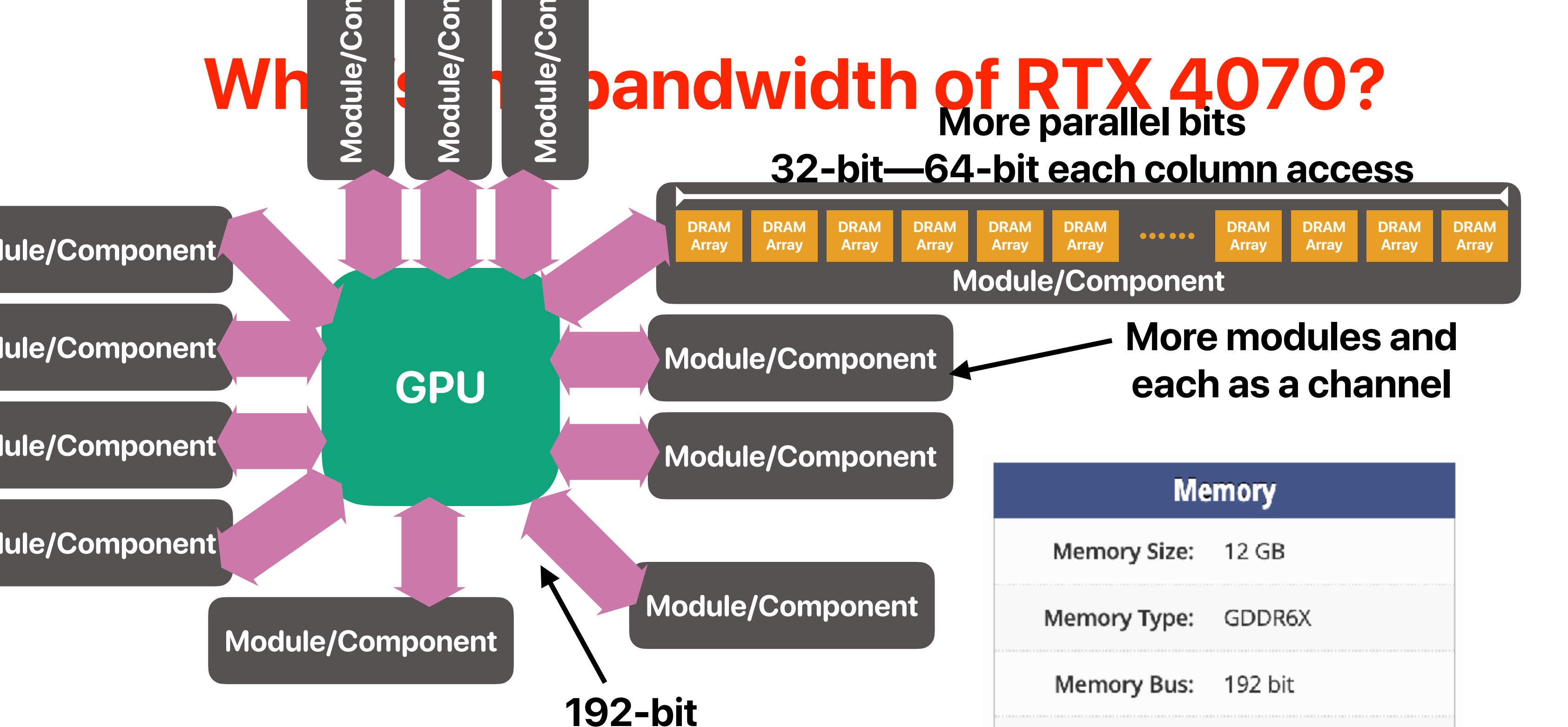
What's the bandwidth of RTX 4070?

	GeForce RTX 4070 Ti SUPER	GeForce RTX 4070 Ti	GeForce RTX 4070 SUPER	GeForce RTX 4070
GPU Engine Specs:				
NVIDIA CUDA® Cores	8448	7680	7168	5888
Boost Clock (GHz)	2.61	2.61	2.48	2.48
Base Clock (GHz)	2.34	2.31	1.98	1.92
Memory Specs:				
Standard Memory Config	16 GB GDDR6X	12 GB GDDR6X	12 GB GDDR6X	12 GB GDDR6X
Memory Interface Width	256-bit	192-bit	192-bit	192-bit

Why is the bandwidth of RTX 4070?

More parallel bits

32-bit—64-bit each column access



More modules and
each as a channel

192-bit

Memory	
Memory Size:	12 GB
Memory Type:	GDDR6X
Memory Bus:	192 bit
Bandwidth:	504.2 GB/s

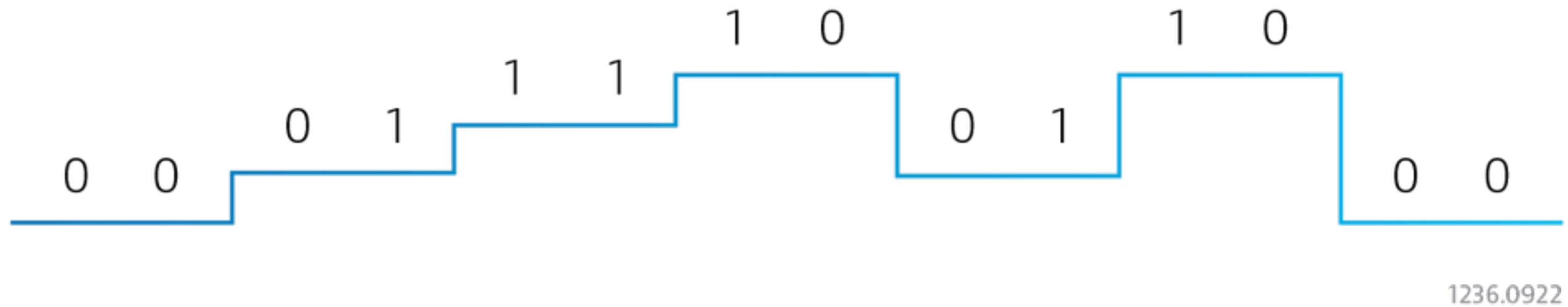
If the bus is **192**-bit (24B) wide, bandwidth is

$$\frac{192}{8} \times 10501 \times 10^6 = 252024 MB/sec = 252 GB/sec$$

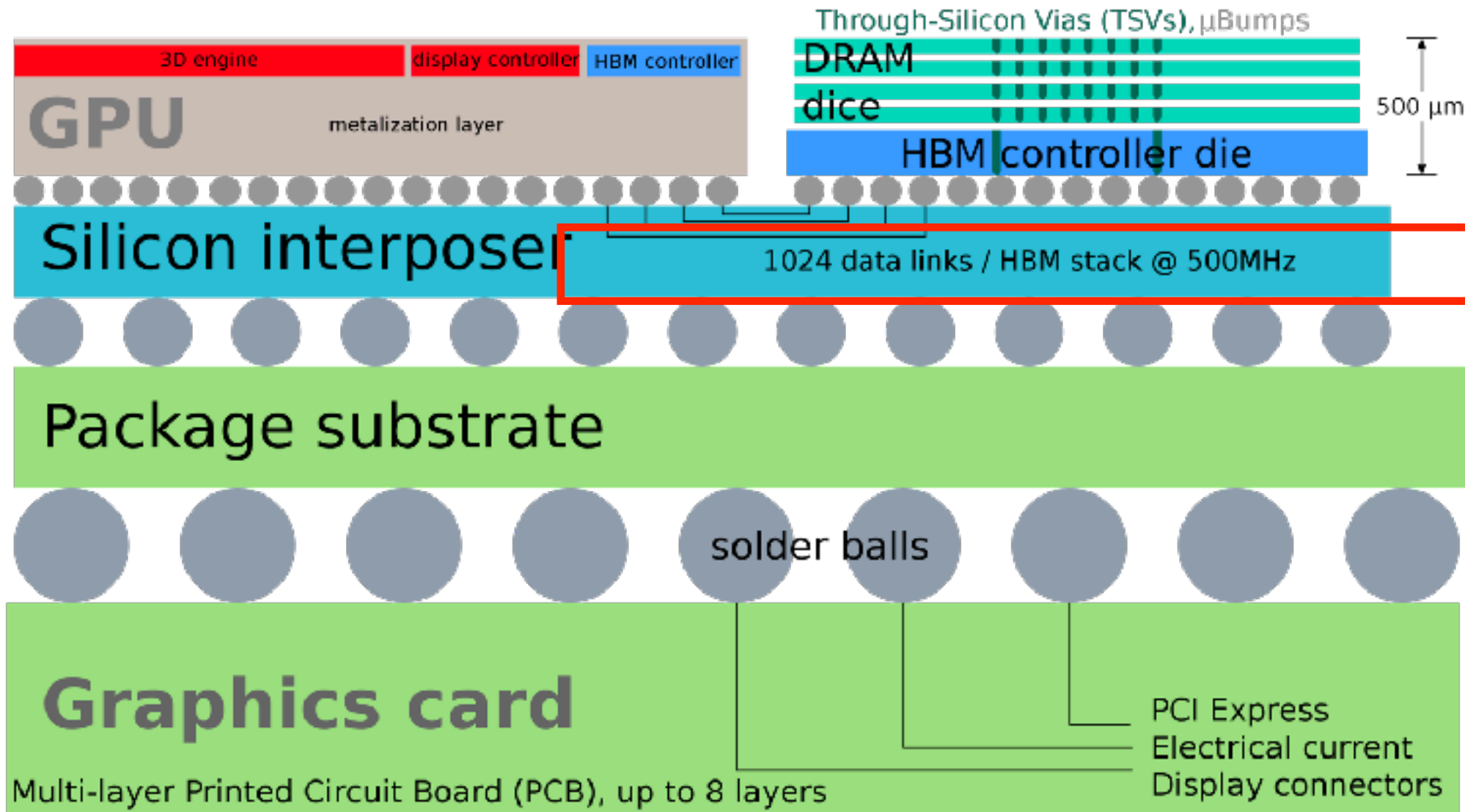
<https://www.techpowerup.com/gpu-specs/geforce-rtx-4070.c3924>

Signaling in DDR6X

PAM4 Signaling Technology



HBM (High Bandwidth Memory)



$$0.5G \times \frac{1024bits}{8bits} = 64GB/sec$$

If you have 4x modules —
4*64 = 256 GB/sec

**HBM2 increase the clock
rate to 2GHz — 1TB/sec**

What's the pros & cons of GDDR v.s. HBM

GDDR v.s. HBM

GDDR v.s. HBM

- HBM's bandwidth is wider
- HBM is more expensive — pin counts are more expensive
- HBM is limited by the per-chip heat dissipation — less powerful

The program

What're the differences between the 1st and the rest runs?

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/time.h>
#include <time.h> /* for clock_gettime */
#include <malloc.h>
#include <immintrin.h>
```

```
void write_memory_avx(void* array, size_t size) {
    __m256i* varray = (__m256i*) array;

    __m256i vals = _mm256_set1_epi32(1);
    size_t i;
    for (i = 0; i < size / sizeof(__m256i); i++) {
        _mm256_store_si256(&varray[i], vals);
        // This will generate the vmovaps instruction.
    }
}
```

What's the purpose of using AVX/Vector Instructions here?

```
int main(int argc, char **argv)
{
    size_t *array;
    size_t size;
    double total_time;
    struct timespec start, end;
    int i;
    size = atoi(argv[1])/sizeof(size_t);
```

```
    array = (size_t *)memalign(32, sizeof(size_t)*size);

    clock_gettime(CLOCK_MONOTONIC, &start);

    write_memory_avx(array, size);

    clock_gettime(CLOCK_MONOTONIC, &end);

    total_time = ((end.tv_sec * 1000000000.0 +
end.tv_nsec) - (start.tv_sec * 1000000000.0 +
start.tv_nsec));
    fprintf(stderr, "Latency: %.01f ns, GBps:
%1f\n", total_time, ((double)((double)size*sizeof(size_t)/
(total_time)));
    clock_gettime(CLOCK_MONOTONIC, &start);

    for(i = 0 ; i< 10; i++)
        write_memory_avx(array, size);

    clock_gettime(CLOCK_MONOTONIC, &end);
    total_time = ((end.tv_sec * 1000000000.0 +
end.tv_nsec) - (start.tv_sec * 1000000000.0 +
start.tv_nsec));
    fprintf(stderr, "Latency (10x average): %.01f ns, GBps
(10x average): %1f\n", total_time/10, ((double)
((double)size*10*sizeof(size_t)/(total_time)));
    return 0;
}
```

The program

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/time.h>
#include <time.h> /* for clock_gettime */
#include <malloc.h>
#include <immintrin.h>

void write_memory_avx(void* array, size_t size) {
    __m256i* varray = (__m256i*) array;

    __m256i vals = _mm256_set1_epi32(1);
    size_t i;
    for (i = 0; i < size / sizeof(__m256i); i++) {
        _mm256_store_si256(&varray[i], vals); // This will
generate the vmovaps instruction.
    }
}

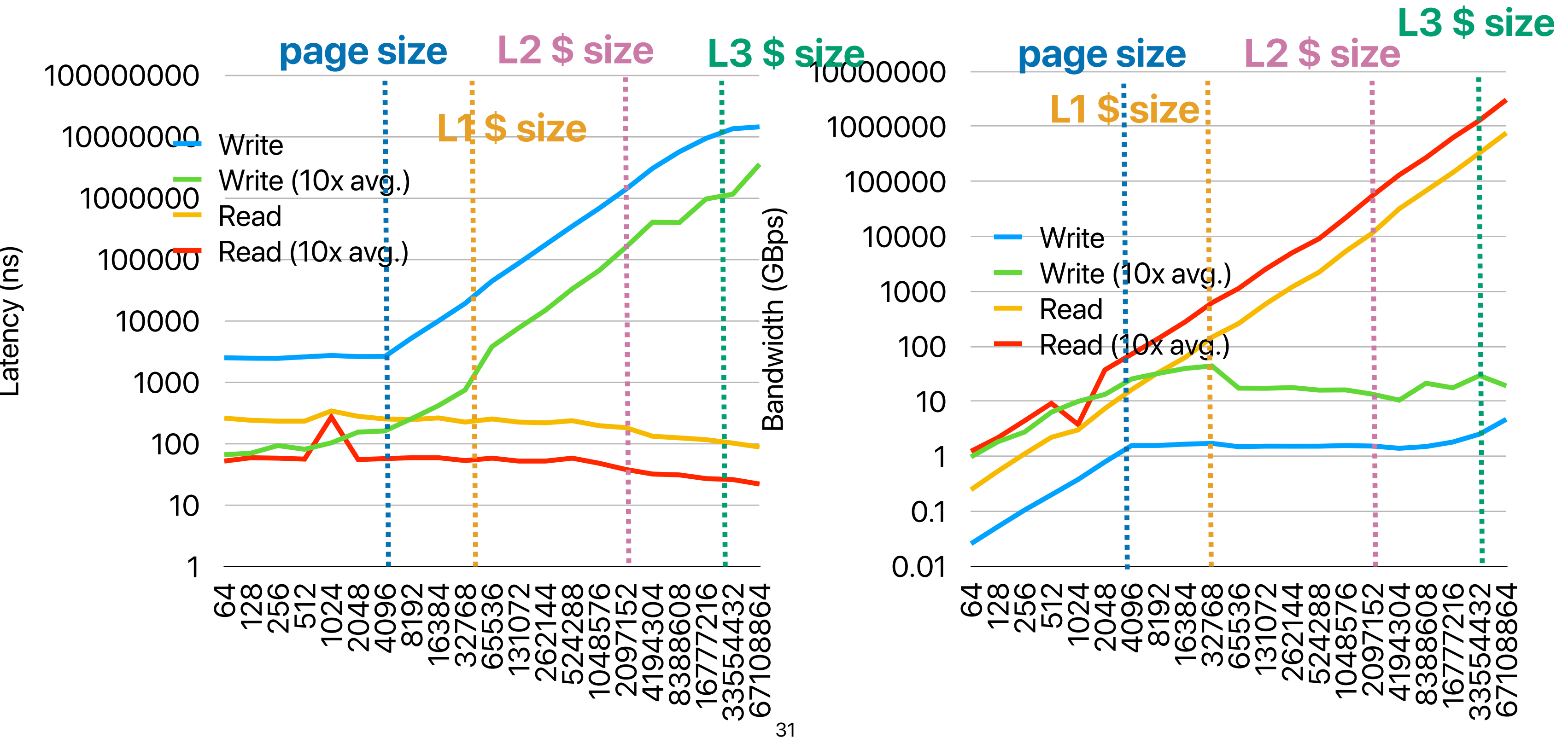
int main(int argc, char **argv)
{
    size_t *array;
    size_t size;
    double total_time;
    struct timespec start, end;
    int i;
    size = atoi(argv[1])/sizeof(size_t);
    array = (size_t *)memalign(32, sizeof(size_t)*size);
```

```
    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    write_memory_avx(array, size);

    clock_gettime(CLOCK_MONOTONIC, &end); /* mark the end
time */
    total_time = ((end.tv_sec * 1000000000.0 +
end.tv_nsec) - (start.tv_sec * 1000000000.0 +
start.tv_nsec));
    fprintf(stderr, "Latency: %.01f ns, GBps:
%.01f\n", total_time, (double)((double)size*sizeof(size_t)/
(total_time)));
    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    for(i = 0 ; i< 10; i++)
        write_memory_avx(array, size);

    clock_gettime(CLOCK_MONOTONIC, &end); /* mark the end
time */
    total_time = ((end.tv_sec * 1000000000.0 +
end.tv_nsec) - (start.tv_sec * 1000000000.0 +
start.tv_nsec));
    fprintf(stderr, "Latency (10x average): %.01f ns, GBps
(10x average): %.01f\n", total_time/10, (double)
((double)size*10*sizeof(size_t)/(total_time)));
    return 0;
}
```

Performance Chart (on Intel Core i7 13700)



The program

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <sys/time.h>
#include <time.h> /* for clock_gettime */
#include <malloc.h>

void write_memory_loop(void* array, size_t size) {
    size_t* carray = (size_t*) array;
    size_t i;
    for (i = 0; i < size / sizeof(size_t); i++) {
        carray[i] = 1;
    }
}

int main(int argc, char **argv)
{
    size_t *array;
    size_t *dest;
    size_t size;
    double total_time;
    struct timespec start, end;
    struct timeval time_start, time_end;
    size = atoi(argv[1])/sizeof(size_t);
    array = (size_t *)malloc(sizeof(size_t)*size);
```

```
    dest = (size_t *)malloc(sizeof(size_t)*size);
    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    #ifdef SIMD
    write_memory_avx(array, size);
    #else
    write_memory_loop(array, size);
    #endif
    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    memcpy(dest, array, size*sizeof(size_t));
    clock_gettime(CLOCK_MONOTONIC, &end); /* mark start
time */
    total_time = ((end.tv_sec * 1000000000.0 +
end.tv_nsec) - (start.tv_sec * 1000000000.0 +
start.tv_nsec));
    fprintf(stderr, "Latency: %.01f ns, GBps: %lf,
%llu\n", total_time, (double)((double)size*sizeof(size_t)/
(total_time)), dest[rand()%size]);
    return 0;
}
```

memmove & memcpy: 5% cycles in Google's datacenter

Svilen Kanev, Juan Pablo Darago, Kim Hazelwood,
Parthasarathy Ranganathan, Tipp Moseley, Gu-Yeon Wei, and
David Brooks. ISCA '15

Why is memcpy so often?

Network protocol stack

How do applications(e.g., server/client) interpret data: HTTPS, FTP, SSH, RTSP ...



How do applications connect to each other (end-to-end reliability): TCP, UDP, RTP



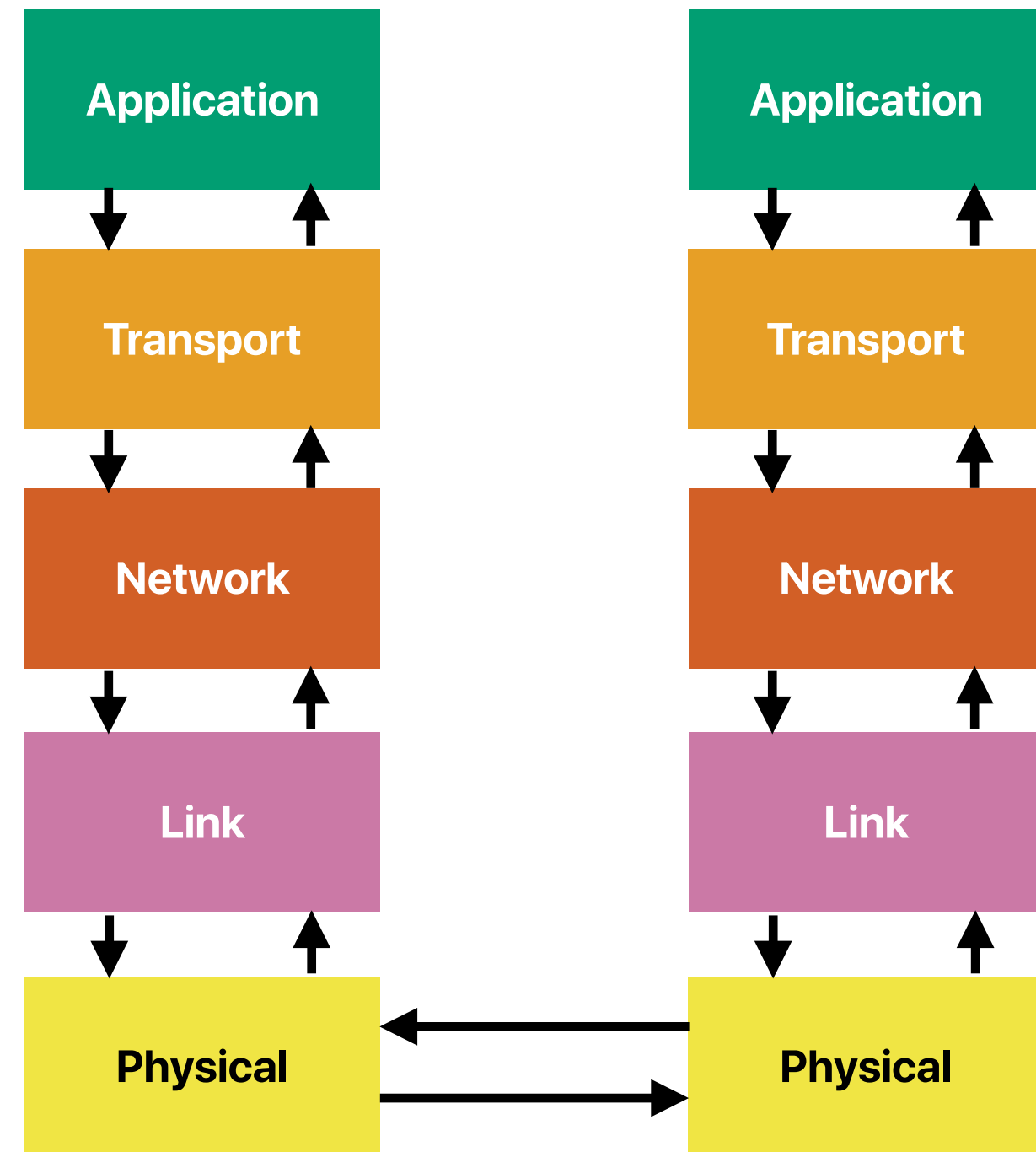
How do computers find each other: IP, ARP, ICMP



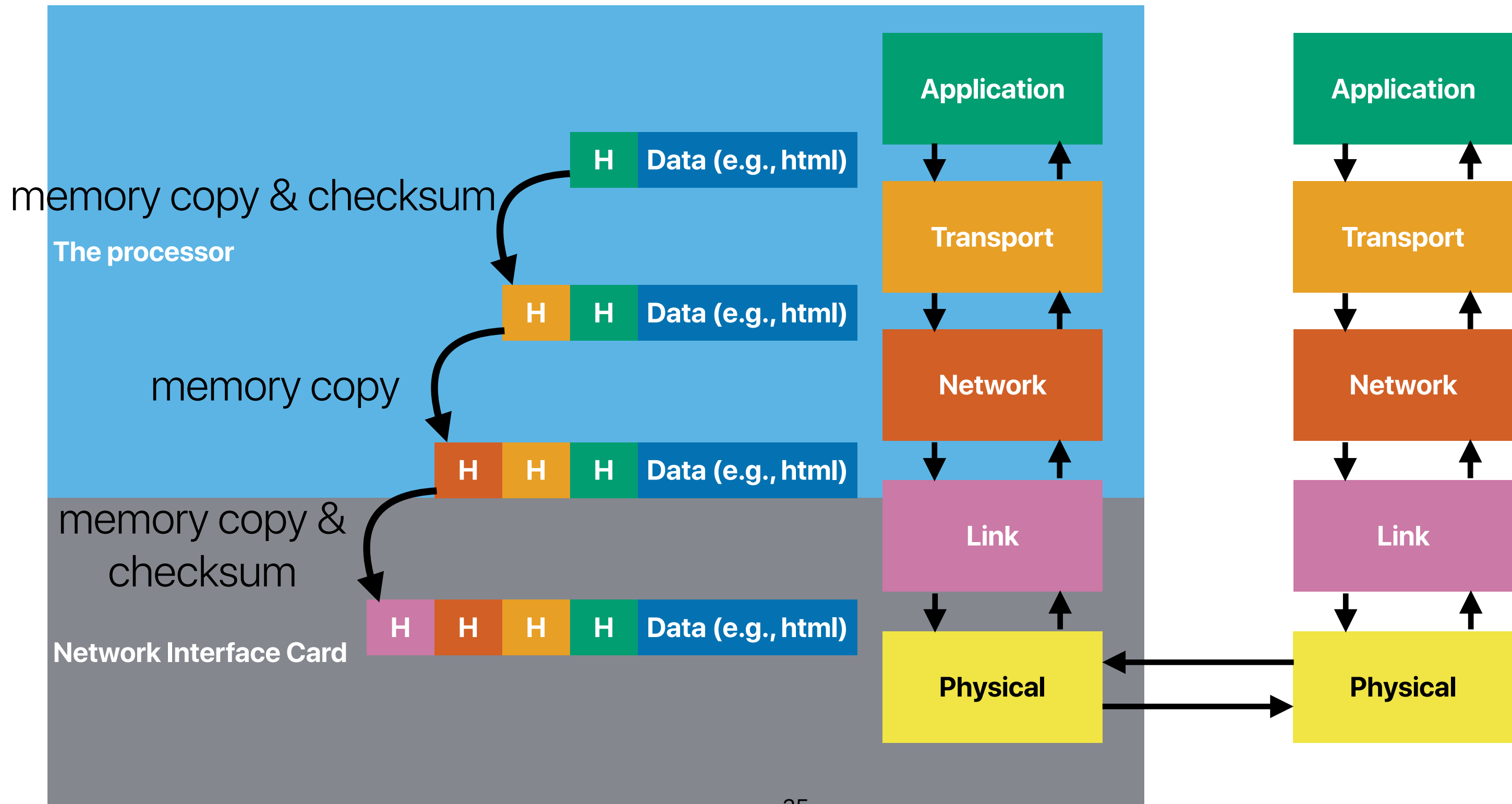
How do a computer use the media:
IEEE 802.3 (ethernet), IEEE 802.11 (WiFi)



How do the computer encode data on the media



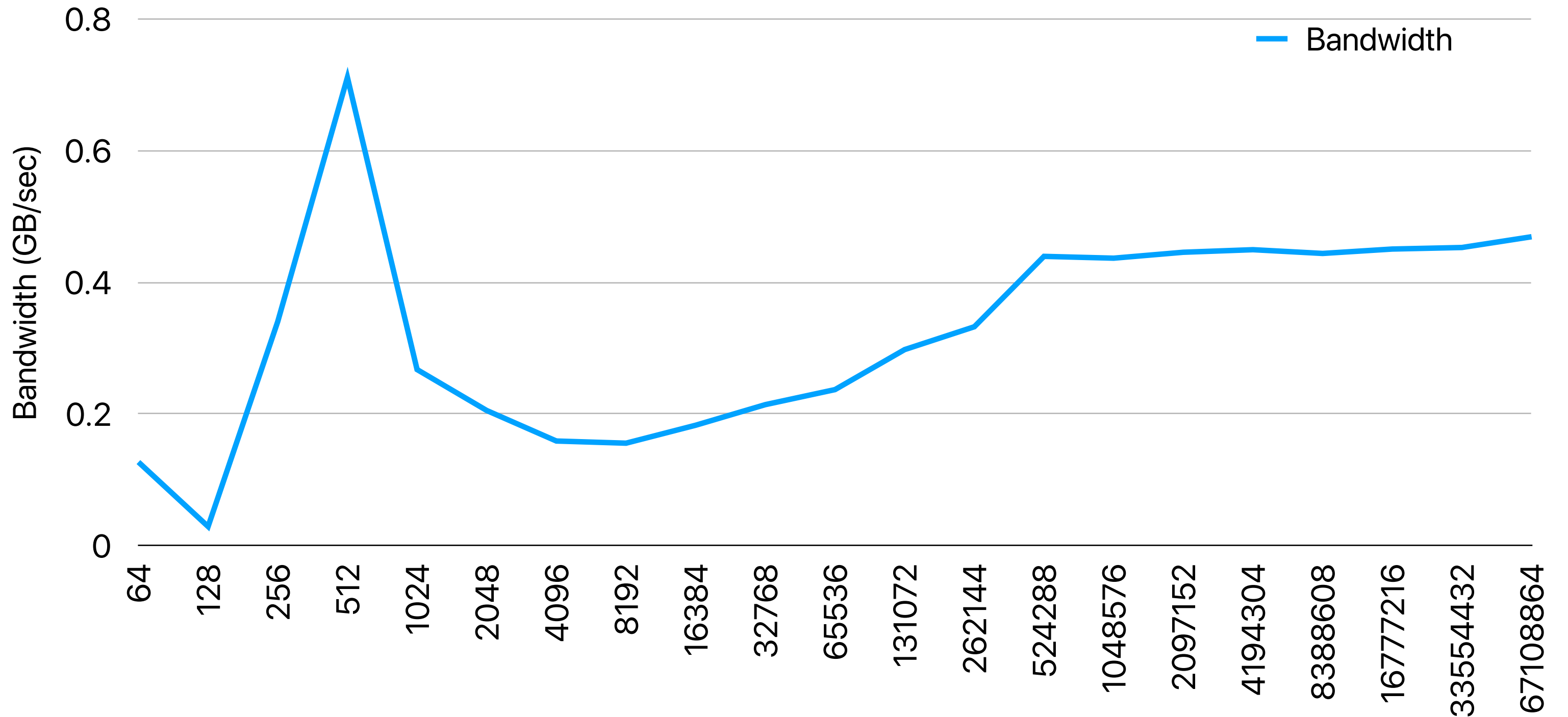
Network protocol stack



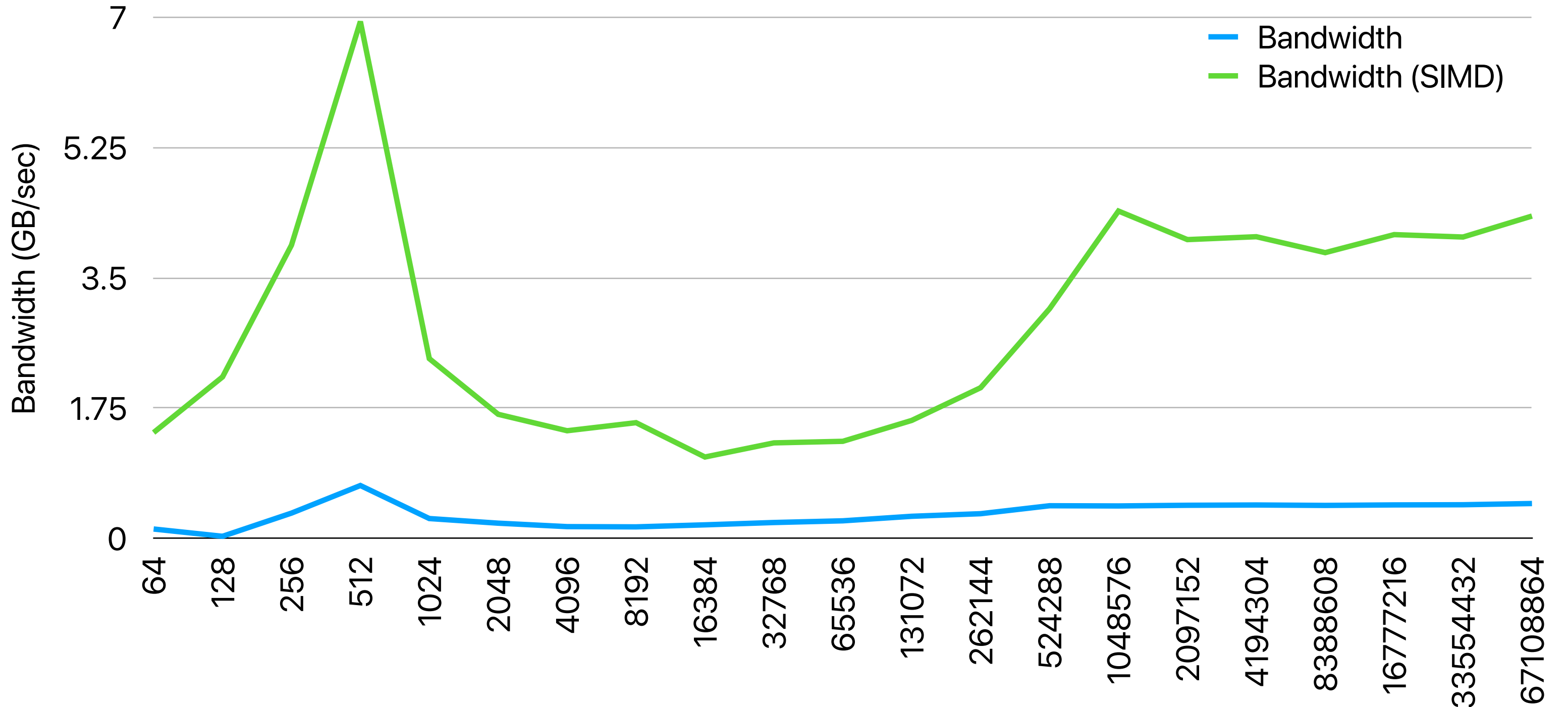
Use cases of memcpy

- Network packets (each layer needs to adjust the header and payload locations)
- Copy-on-write (spawning new process)

Memory copy bandwidth



Memory copy bandwidth using SIMD



Google Project Presentations

- Make an appointment on 2/6 and 2/8 through the Google Calendar **ASAP**
- 12 minute presentation with 3 minute Q & A
- Why & what & how!!! — considering you're giving a presentation at Apple's keynote
 - 6-minute why — why should everyone care about this problem? Why is this still a problem?
 - 4-minute what — what are you proposing in this project to address the problem?
 - 2-minute how — expected platforms/engineering efforts, milestones and workload distribution among members
 - Please reference this article to make a good presentation <https://cseweb.ucsd.edu/~swanson/GivingTalks.html>

Electrical Computer Science Engineering

277

つく

