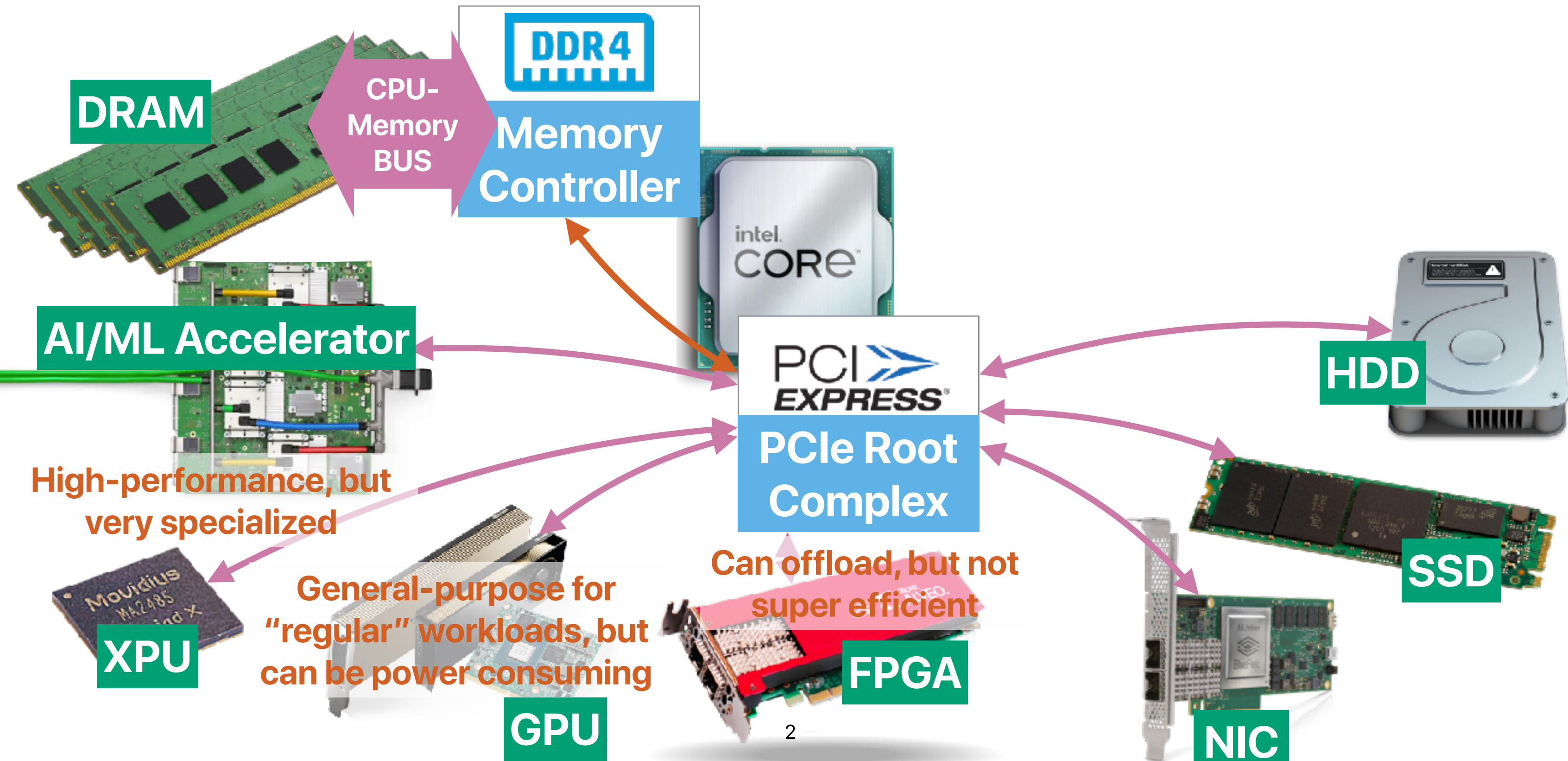


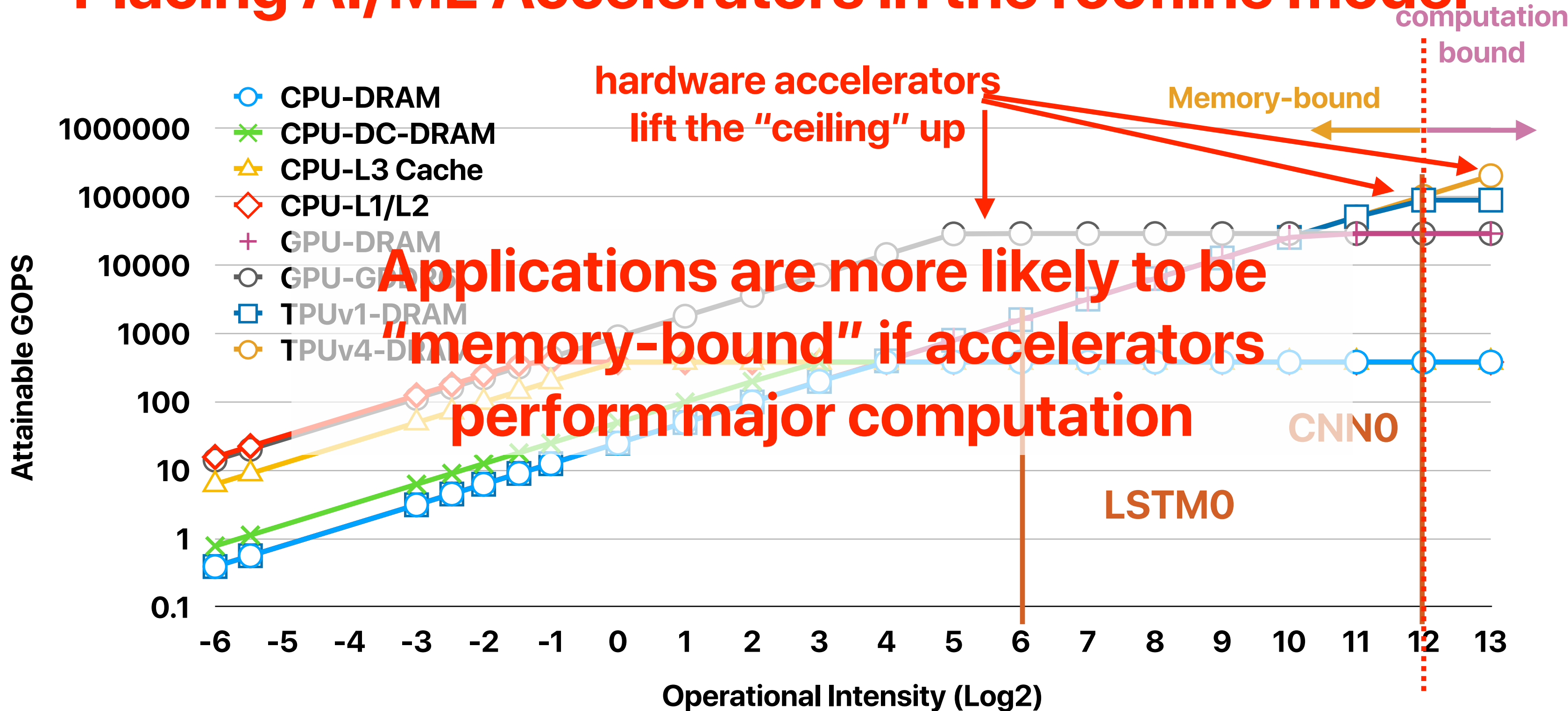
Memory subsystem (3)

Hung-Wei Tseng

Recap: The landscape of modern computers



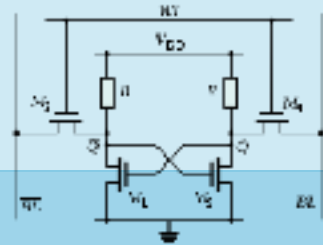
Placing AI/ML Accelerators in the roofline model



Recap: Memory technologies we have today

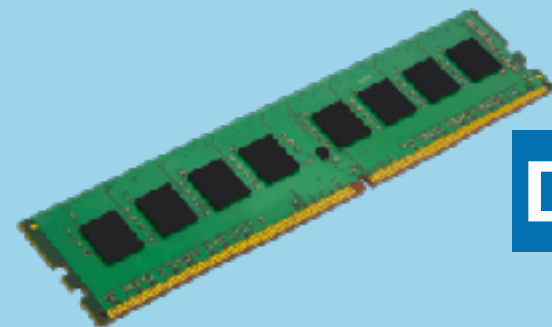
Volatile Memory

100ps



SRAM

ns

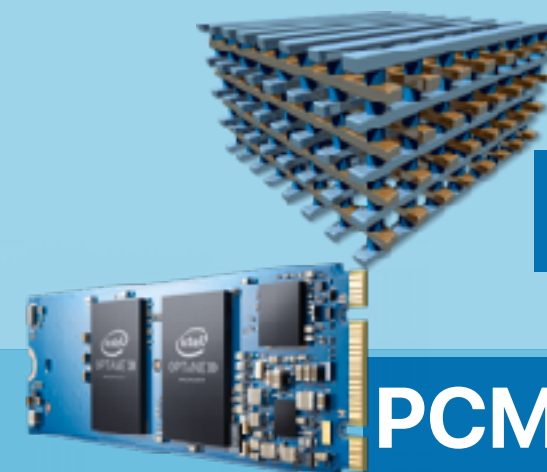


DRAM

us

ms

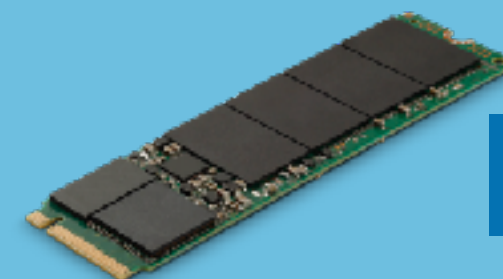
Non-Volatile Memory



RRAM

PCM

3D XPoint



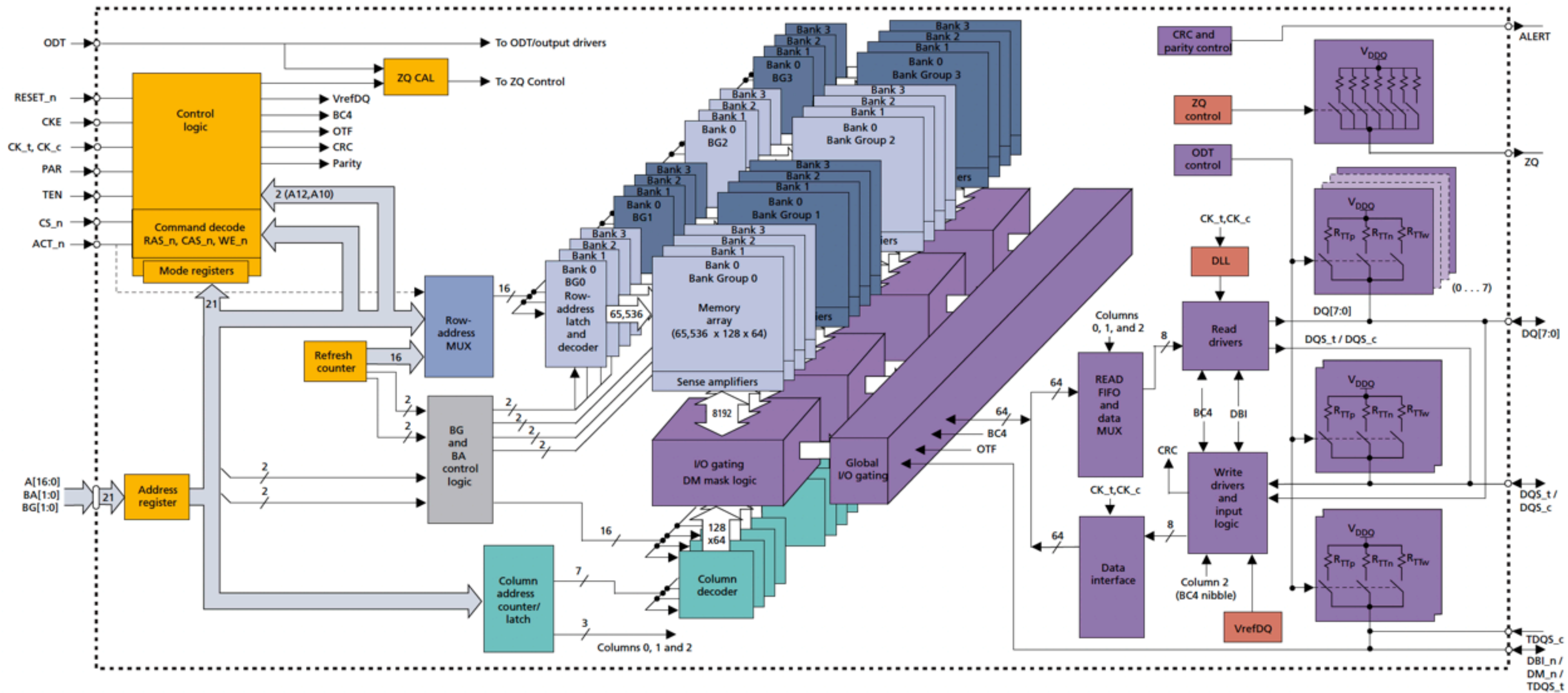
Flash memory



Hard Disk Drives

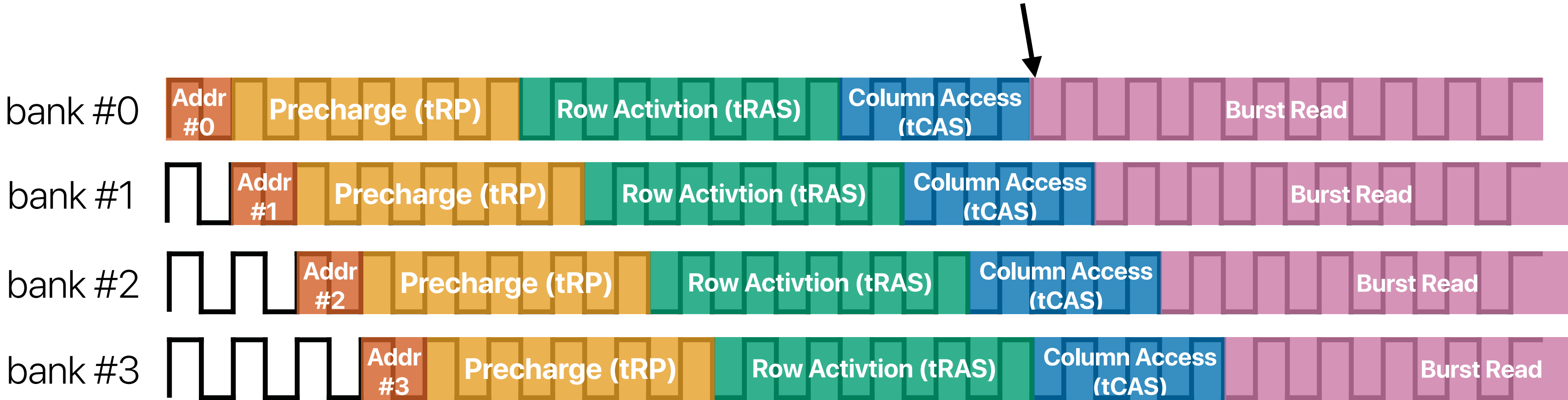
Recap: Improving memory-bounded applications

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



Recap: Multi-bank access

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

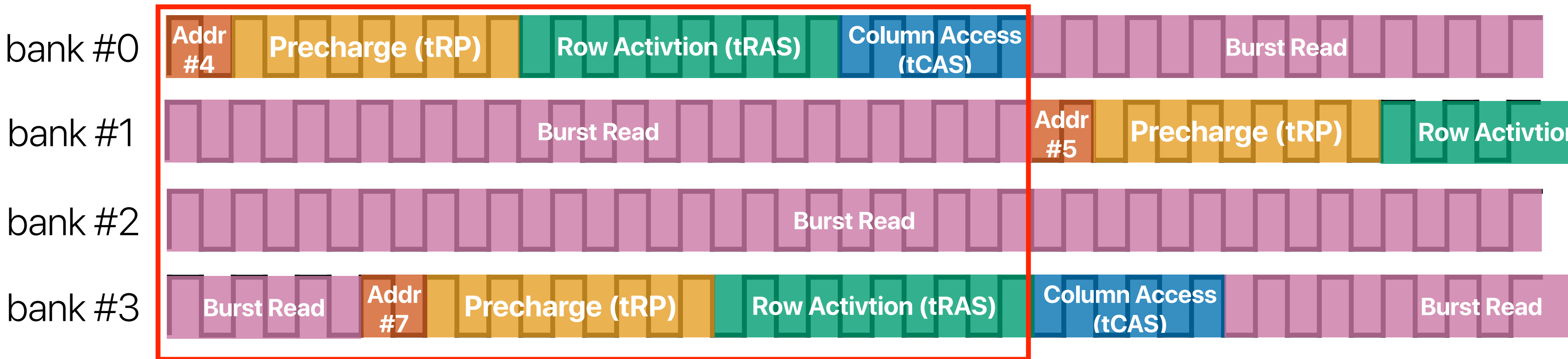


only one bank can accept request each cycle

the memory bandwidth
can be fully utilized after
this

Multi-bank access

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



Outline

- Improve the performance of main memory subsystem
- Near/In-memory processing

Memory bandwidth

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

Ideas of increasing bandwidth

Ideas of increasing bandwidth

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

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

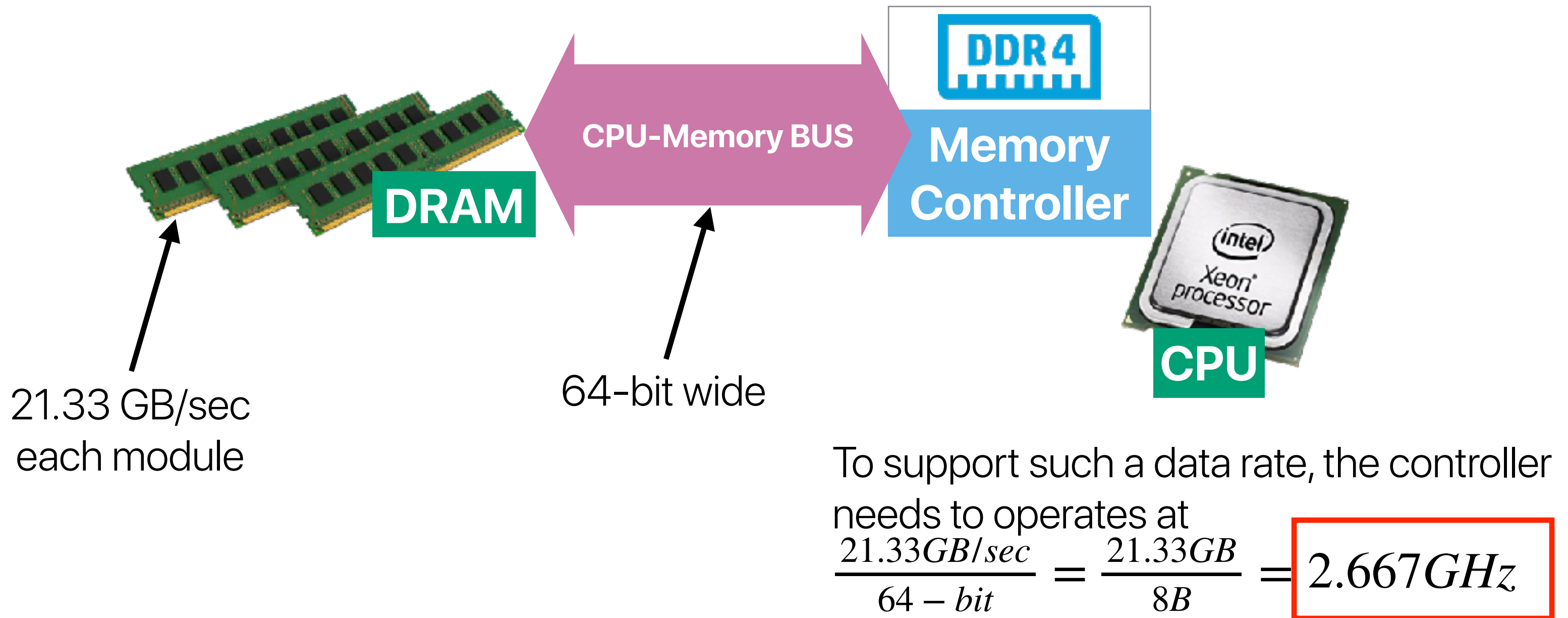
The 8Gb DDR4 SDRAM C-die is organized as a 84Mbit x 16 I/Os x 8banks device. This synchronous device achieves high speed double-data-rate transfer rates of up to 2666Mb/sec/pin (DDR4-2666) for general applications.

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?



DDR5-5600

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

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

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

Case: what's the goal of the following program and what do you expect the performance would change when size varies?

```
#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 */

    write_memory_loop(array, size);

    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    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: %.0lf ns, GBps: %lf,
%llu\n", total_time, (double)((double)size*sizeof(size_t)/
(total_time)), dest[rand()%size]);
    return 0;
}
```

Electrical Computer Science Engineering

277

つくづく

