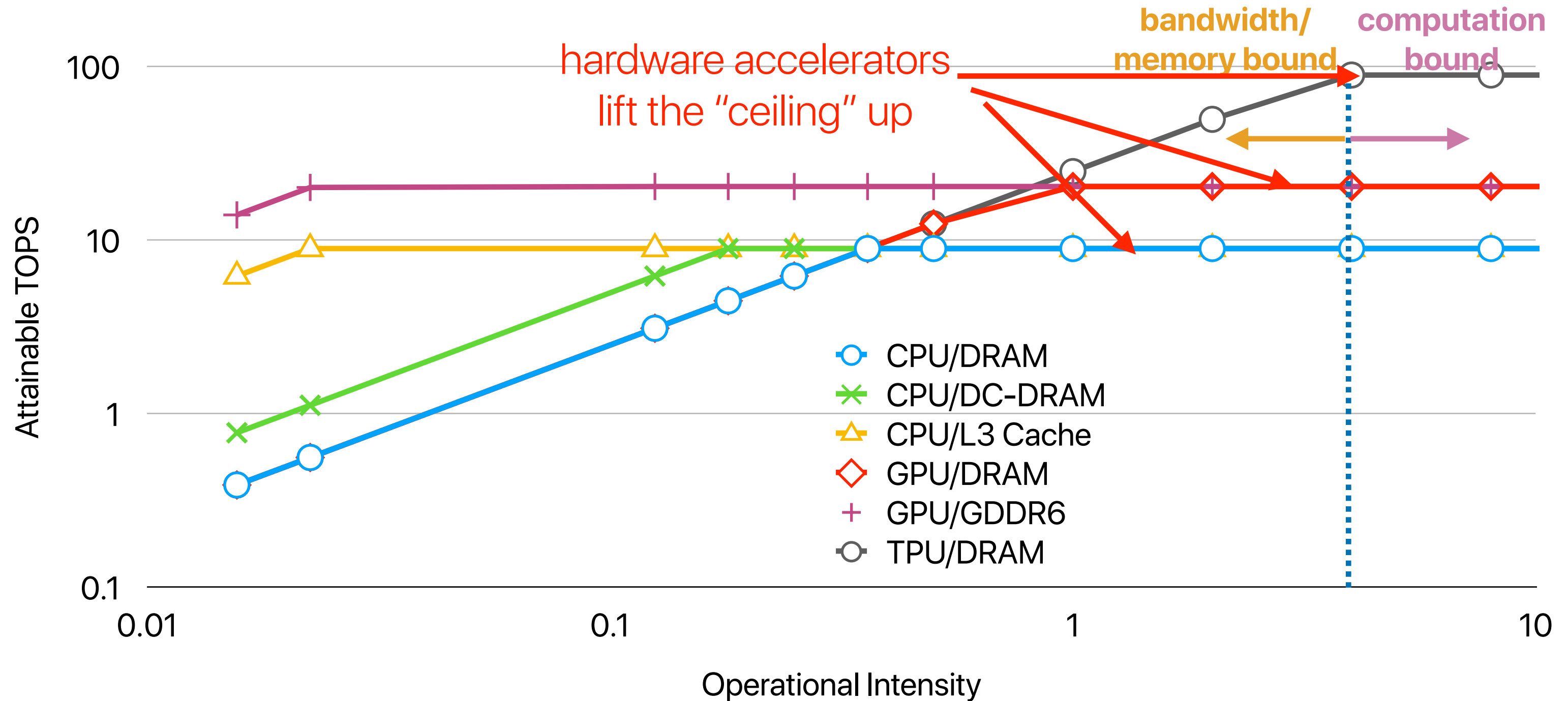# Modern Heterogeneous Computers: (4) Memory Components II

Hung-Wei Tseng

# Recap: the roofline after using hardware accelerators

# Memory technologies we have today

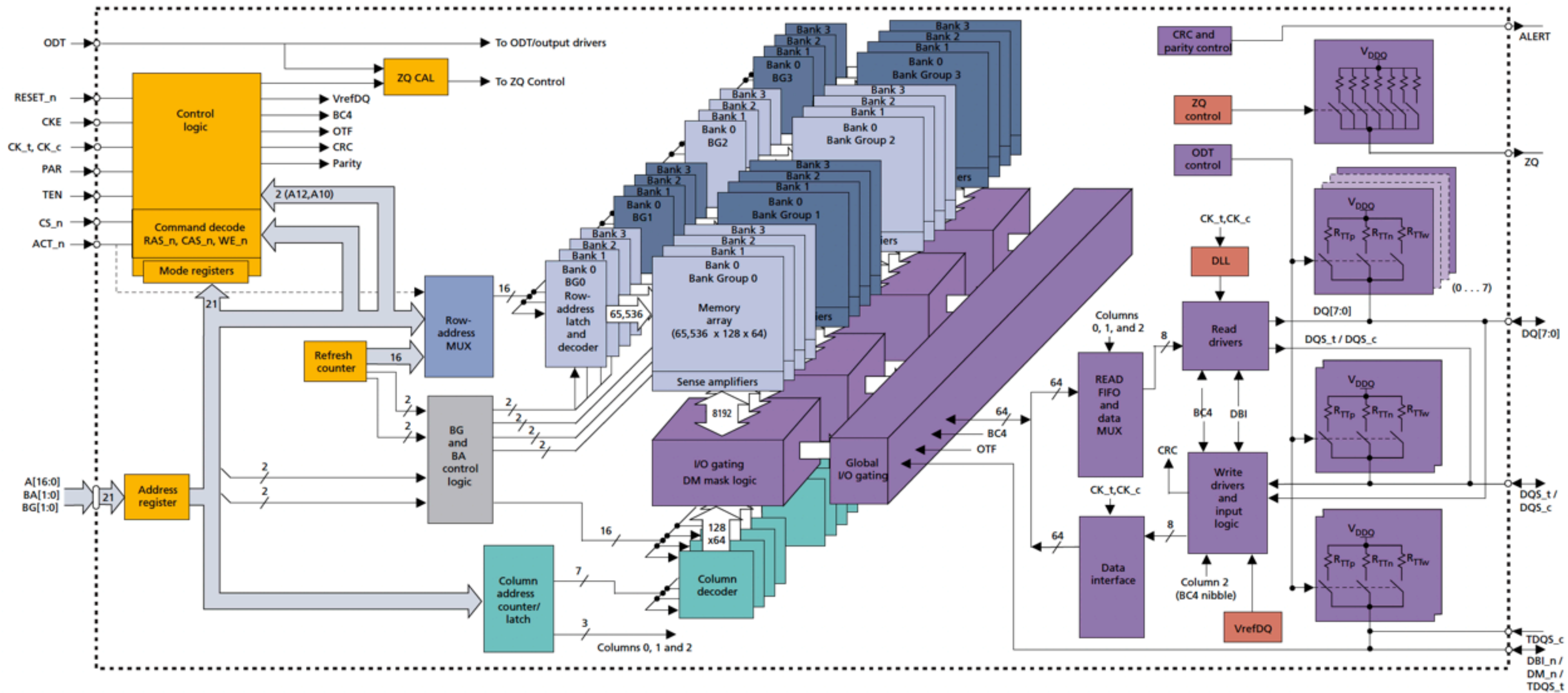**Volatile Memory**                    **Non-Volatile Memory**

**100ps**          **SRAM**

**ns**            **DRAM**                              **RRAM**

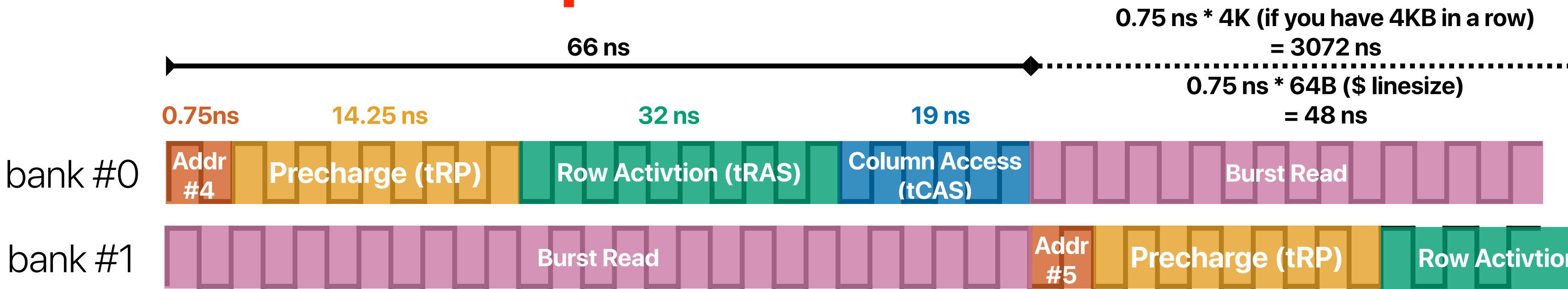                                              **PCM    3DXPoint**
                                        **Flash memory**

**us**

**ms**                                  **Hard Disk Drives**

4

# Recap: DRAM Accesses

0.75 ns * 4K (if you have 4KB in a row) = 3072 ns

66 ns

0.75 ns * 64B ($ linesize) = 48 ns

| 0.75ns | 14.25 ns | 32 ns | 19 ns | |

**bank #0**

| Addr #4 | Precharge (tRP) | Row Activtion (tRAS) | Column Access (tCAS) | Burst Read |

**bank #1**

| Burst Read | Addr #5 | Precharge (tRP) | Row Activtion |

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

## 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 | 14.06 | 14.16 | 14.25 | ns |
| tRP(min) | 13.75 | 13.92 | 14.06 | 14.16 | 14.25 | ns |
| tRAS(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)

The 8Gb DDR4 SDRAM C-die is organized as a 64Mbit x 16 I/Os x 8banks

# Recap: DRAM Performance

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

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

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

## 2. Key Features

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

| Speed | DDR4-1600 | DDR4-1866 | DDR4-2133 | DDR4-2400 | DDR4-2666 |
|---|---|---|---|---|---|
|  | 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 |
| CAS Latency | 11 | 13 | 15 | 17 | 19 |
| tRCD(min) | 13.75 | 13.92 | 14.06 | 14.16 | 14.25 |
| tRP(min) | 13.75 | 13.92 | 14.06 | 14.16 | 14.25 |
| tRAS(min) | 35 | 34 | 33 | 32 | 32 |
| tRC(min) | 48.75 | 47.92 | 47.06 | 46.16 | 46.25 |

- 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
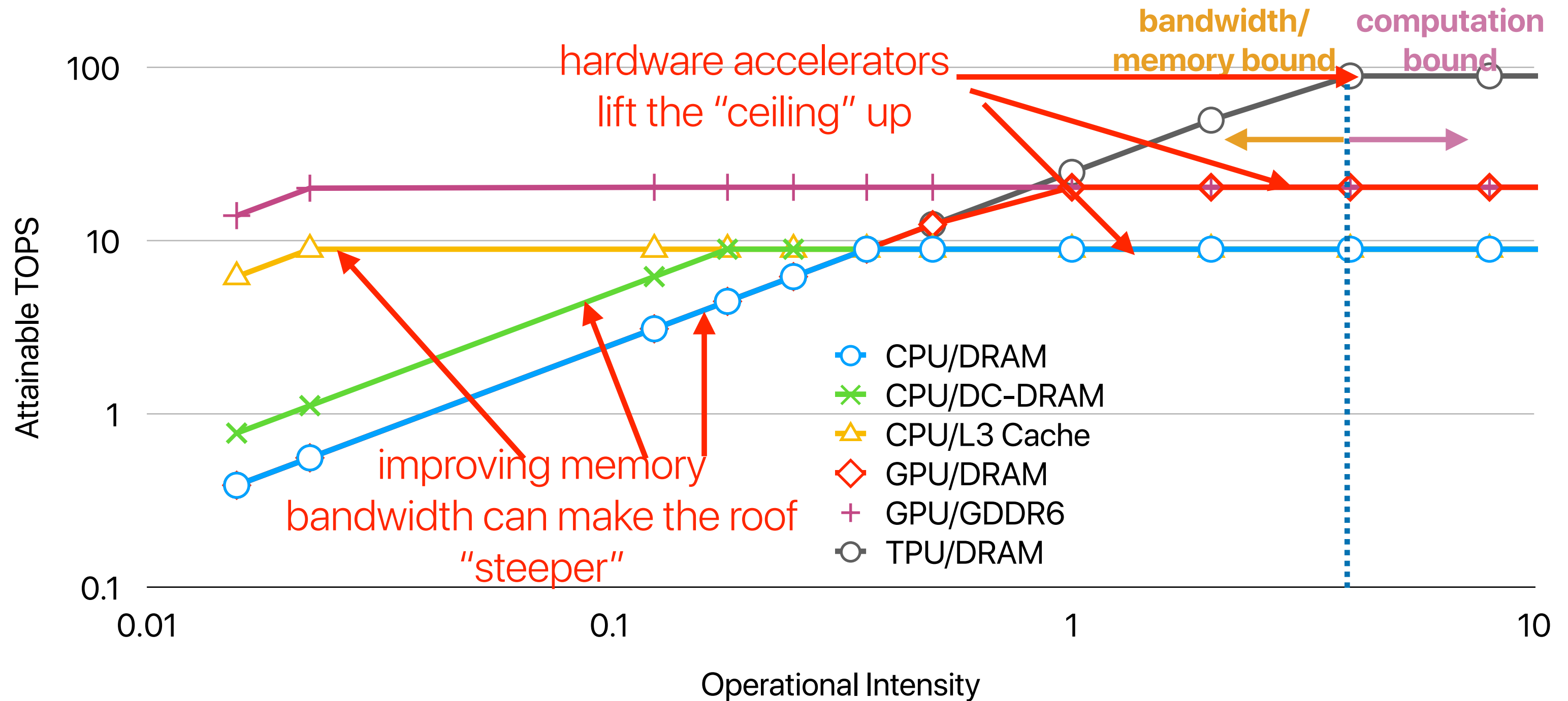- Connectivity Test Mode (TEN) is Supported

The 8Gb DDR4 SDRAM C-die is organized as a 64Mbit x device. This synchronous device achieves high speed transfer rates of up to 2666Mb/sec/pin (DDR4-2666) for tions.

The chip is designed to comply with the following key DD tures such as posted CAS, Programmable CWL, Internal On Die Termination using ODT pin and Asynchronous Re All of the control and address inputs are synchronized wi nally supplied differential clocks. Inputs are latched at the ferential clocks (CK rising and $\overline{CK}$ falling). All I/Os are syr pair of bidirectional strobes (DQS and $\overline{DQS}$) in a source s ion. The address bus is used to convey row, column, a information in a $\overline{RAS}/\overline{CAS}$ multiplexing style. The DDR4 with a single 1.2V (1.14V~1.26V) power supply, 1.2V(1.1 and 2.5V (2.375V~2.75V) $V_{PP}$.

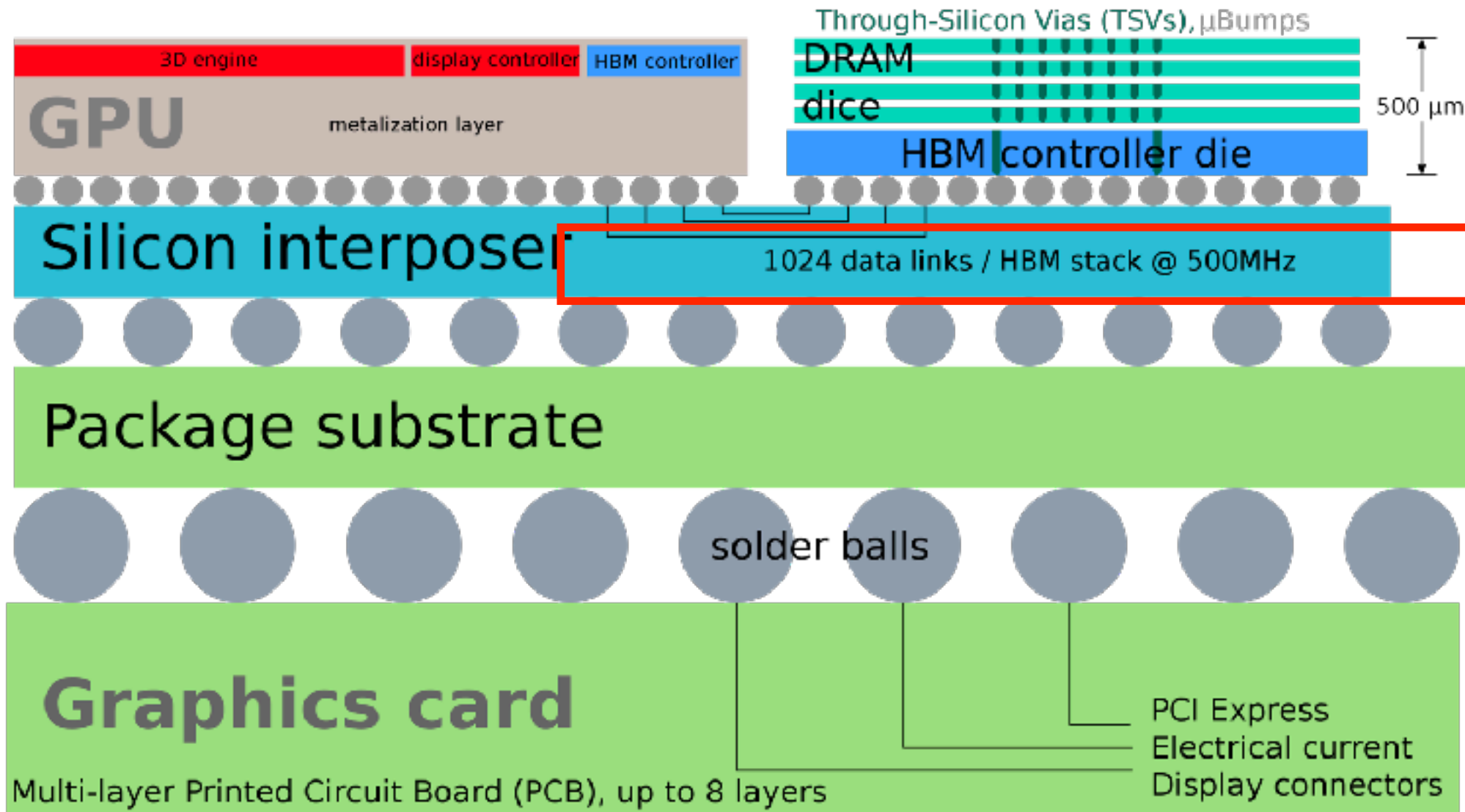The 8Gb DDR4 C-die device is available in 96ball FBGAs

# Recap: the roofline after using hardware accelerators

# **Ideas of increasing bandwidth**

- More parallel bits

- More banks

- More channels

- Widen the memory-processor bus

# GDDR

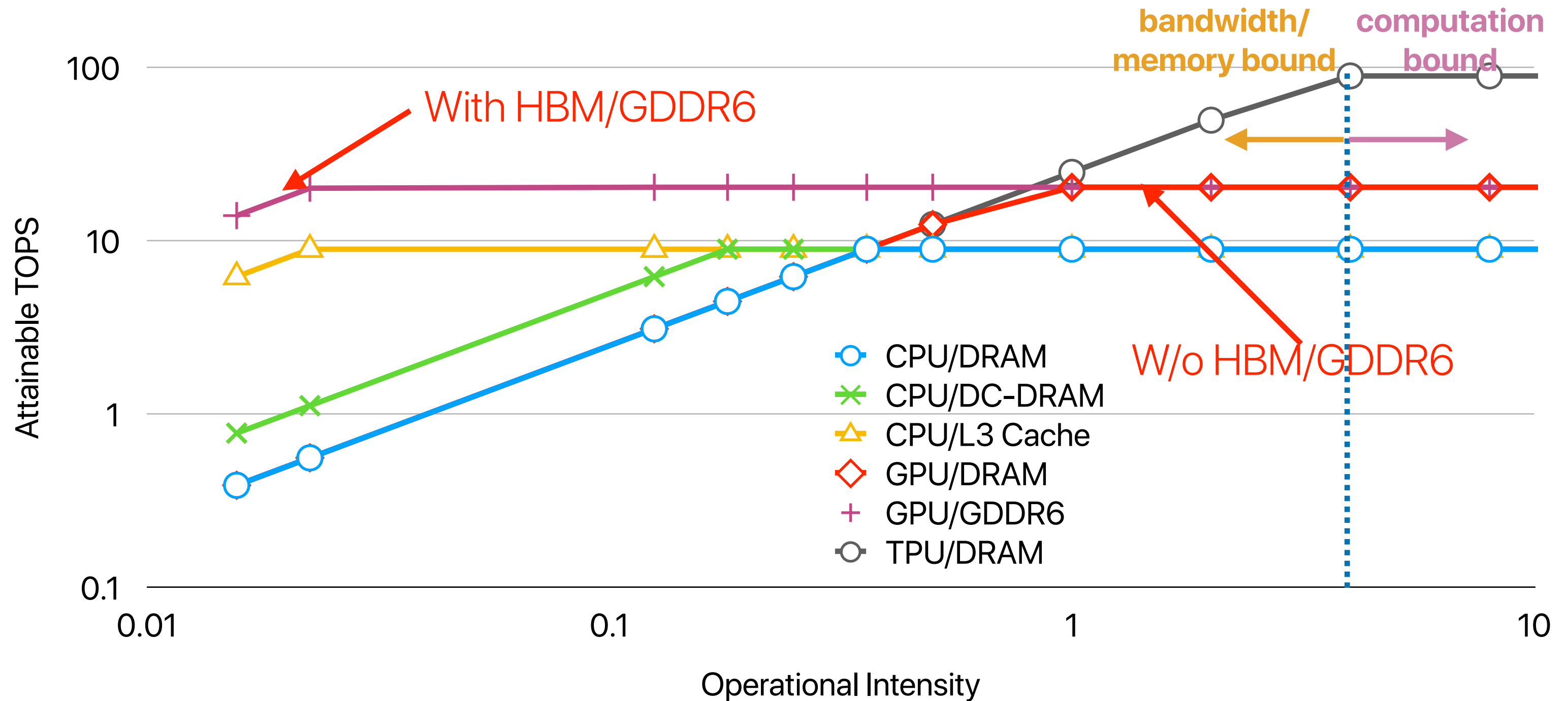**More parallel bits**
**32-bit—64-bit each column access**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| DRAM Array | DRAM Array | DRAM Array | DRAM Array | DRAM Array | DRAM Array | ...... | DRAM Array | DRAM Array | DRAM Array | DRAM Array |

**Module/Component**

Module/Component (top ×3)

Module/Component (left ×4)

**GPU**

Module/Component

Module/Component

Module/Component

Module/Component

**More modules and each as a channel**

If each DRAM array operates at 500 MHz, the effective bandwidth per component

$$32 \times 0.5G \times \frac{32bits}{8bits} = 64GB/sec$$

**Wider bus: 192-bit — 384-bit bus**

If the bus is 256-bit (32B) wide, the memory controller needs to be at

$$\frac{64GB/sec}{32B} = 2GHz$$

**If you have 12x modules — 12*64 = 768 GB/sec**

9

# HBM (High Bandwidth Memory)



Through-Silicon Vias (TSVs), µBumps

DRAM dice — 500 µm

HBM controller die

GPU — 3D engine, display controller, HBM controller, metalization layer

Silicon interposer — 1024 data links / HBM stack @ 500MHz

Package substrate

solder balls

Graphics card — PCI Express, Electrical current, Display connectors

Multi-layer Printed Circuit Board (PCB), up to 8 layers

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

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

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

# Recap: the roofline after using hardware accelerators

# 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

Which one do you prefer?
A. Single chip, heterogeneous processor (e.g., NVIDIA/intel) or
B. standalone accelerators? (e.g., TPU)

# Single-chip or standalone ones?

# Is system-wide heterogeneous processing a better idea?

- Easier for heat dissipation
  - Each processor can be more powerful
    - gamers/datacenters still prefers discrete GPUs
    - Cloud TPUs are standalone ones
- System size is larger, cost of ownership can be higher
- Data movement going through system interconnects
- Allow the "computation core" of the accelerator to enjoy larger bandwidth, if you do it right.

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

```c
#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);
```

```c
    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: %.0lf ns, GBps:
%lf\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): %.0lf ns, GBps
(10x average): %lf\n",total_time/10, (double)
((double)size*10*sizeof(size_t)/(total_time)));
    return 0;
}
```

https://github.com/hungweitseng/EE277/tree/main/demo/memory

# The program

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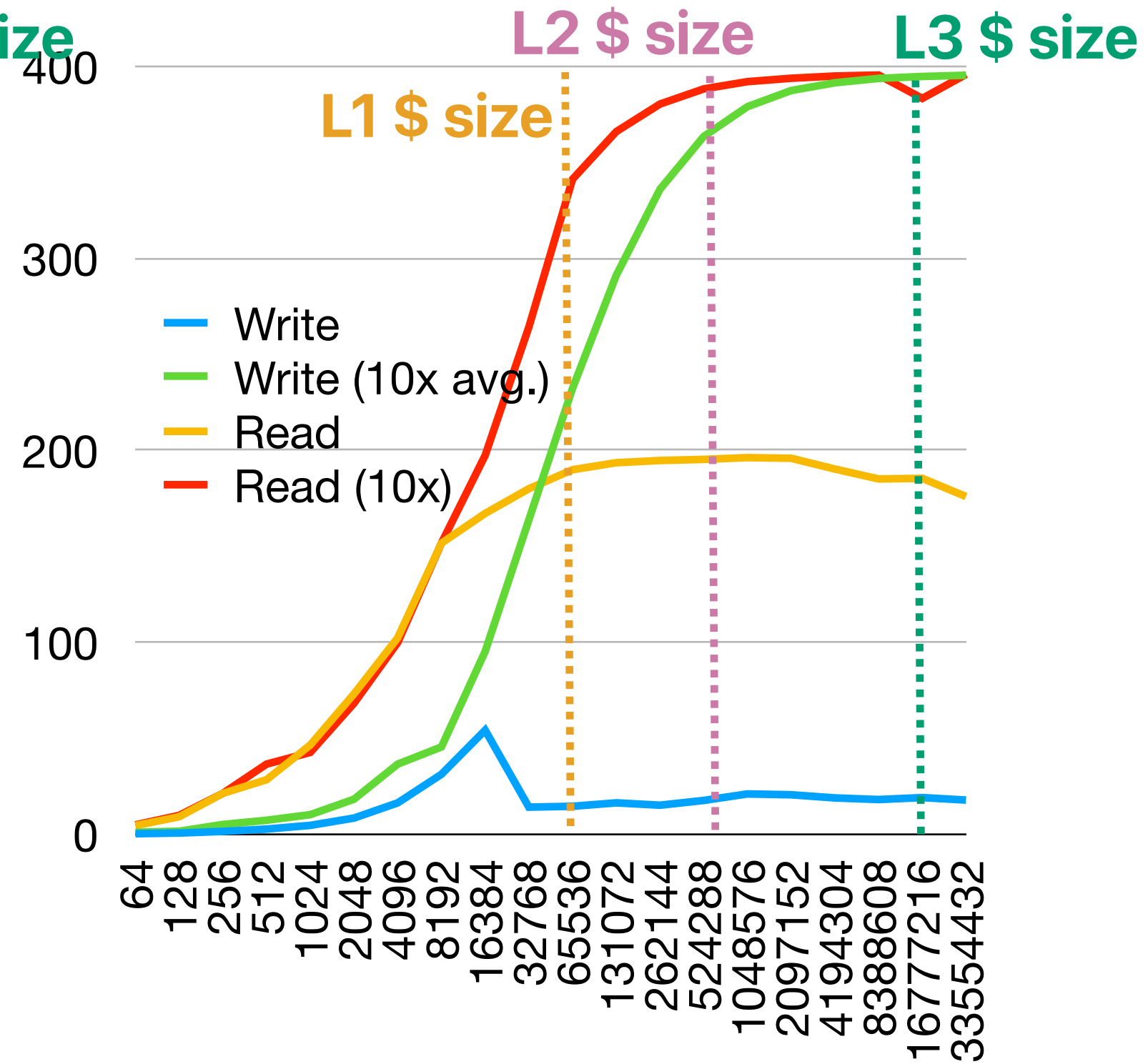
```c
    clock_gettime(CLOCK_MONOTONIC, &start); /* mark start
time */
    write_memory_avx(array, size);

    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: %.0lf ns, GBps:
%lf\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): %.0lf ns, GBps
(10x average): %lf\n",total_time/10, (double)
((double)size*10*sizeof(size_t)/(total_time)));
    return 0;
}
```

19

# Performance Chart (on AMD RyZen 5 2600)

# 3D-die stacking

- Providing huge bandwidth between memory chips and processors (e.g., HBM)
- The renaissance of near-memory processing!
  - Zhu, Qiuling, Berkin Akin, H. Ekin Sumbul, Fazle Sadi, James C. Hoe, Larry Pileggi, and Franz Franchetti. A 3D-stacked logic-in-memory accelerator for application-specific data intensive computing.In 3DIC. 2013.
  - Dongping Zhang, Nuwan Jayasena, Alexander Lyashevsky, Joseph L. Greathouse, Lifan Xu, and Michael Ignatowski. TOP-PIM: throughput-oriented programmable processing in memory. HPDC '14. 2014
  - Farmahini-Farahani, Amin, Jung Ho Ahn, Katherine Morrow, and Nam Sung Kim. NDA: Near-DRAM acceleration architecture leveraging commodity DRAM devices and standard memory modules. In HPCA. 2015.
  - Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi. A scalable processing-in-memory accelerator for parallel graph processing. ISCA '15. 2015.
  - Youngeun Kwon, Yunjae Lee, and Minsoo Rhu. TensorDIMM: A Practical Near-Memory Processing Architecture for Embeddings and Tensor Operations in Deep Learning. In MICRO '52. 2019

# Electrical
## Computer Science
## Engineering

つづく