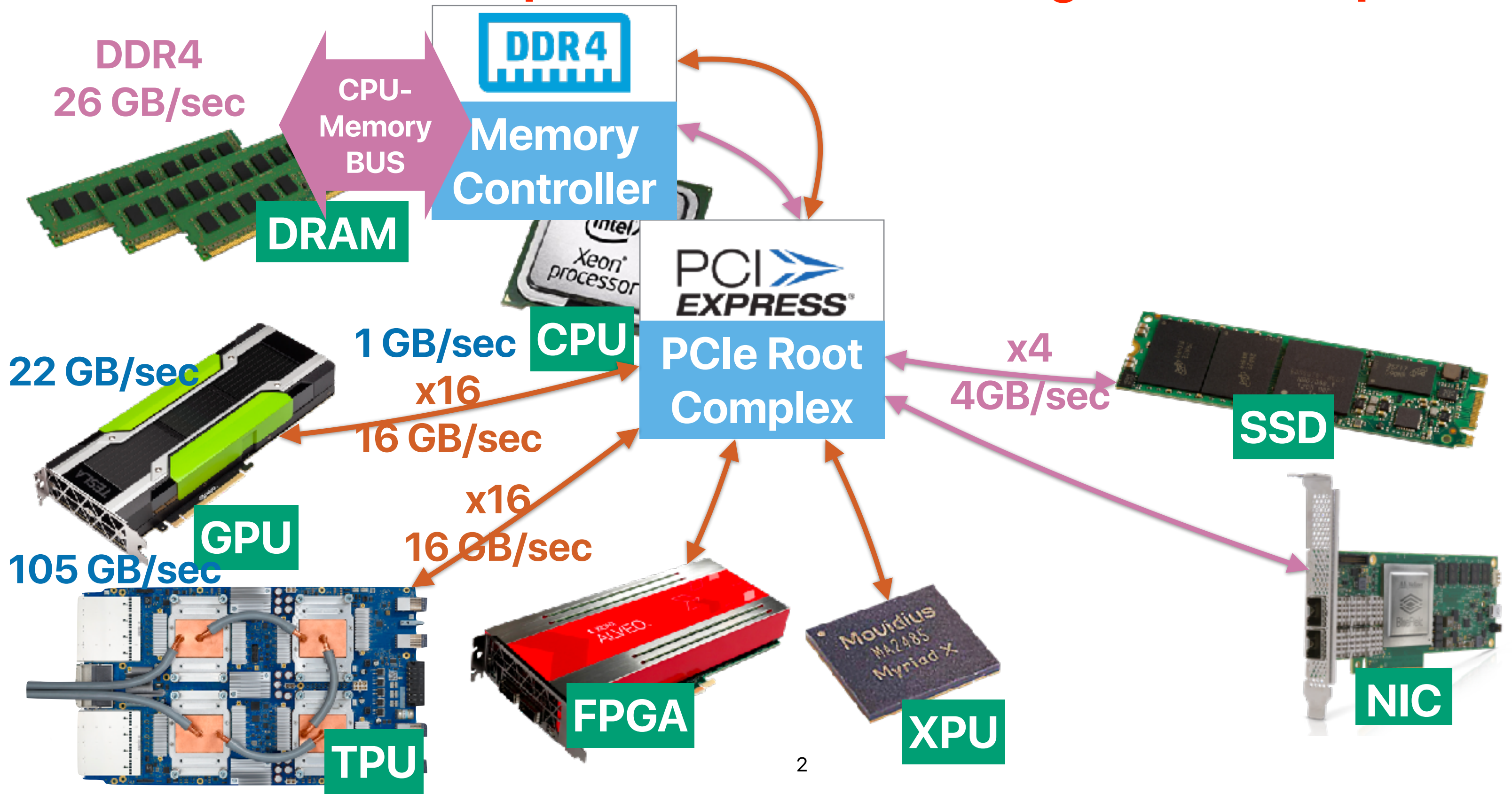


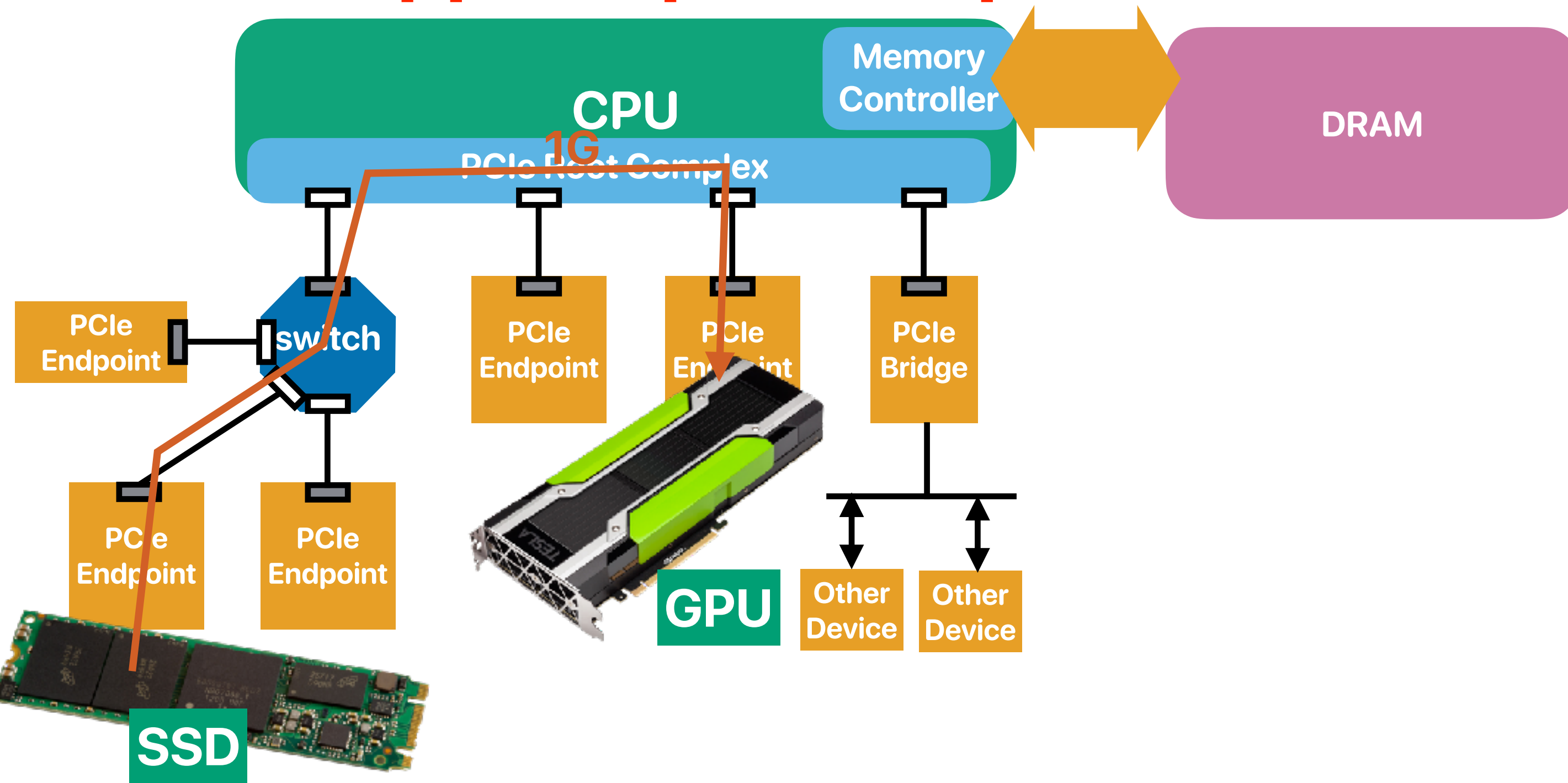
The Concept of Near Data Processing & In-Storage Processing

Hung-Wei Tseng

Review: The "data path" in modern heterogeneous computers



PCIe supports peer-to-peer communication!



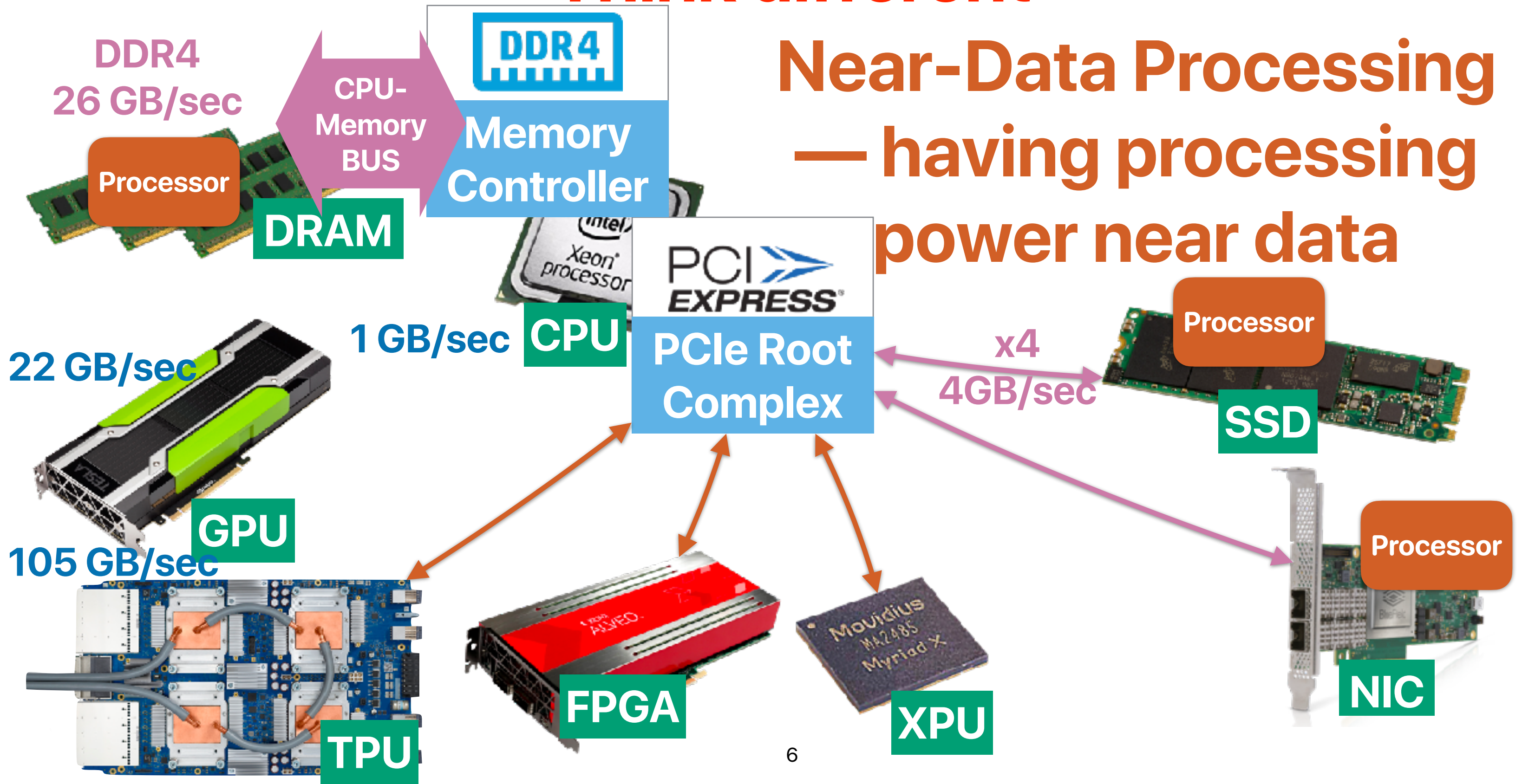
Recap: simply P2P is not enough

- If your data resides in SSD — the total out-going bandwidth from SSD is still way lower than the throughput of your computing resources
- If your dataset is large, you will need to partition your computation — DRAM is a better buffer than an SSD

**Think different — do we really need
to separate data and computation?**

Think different

Near-Data Processing — having processing power near data



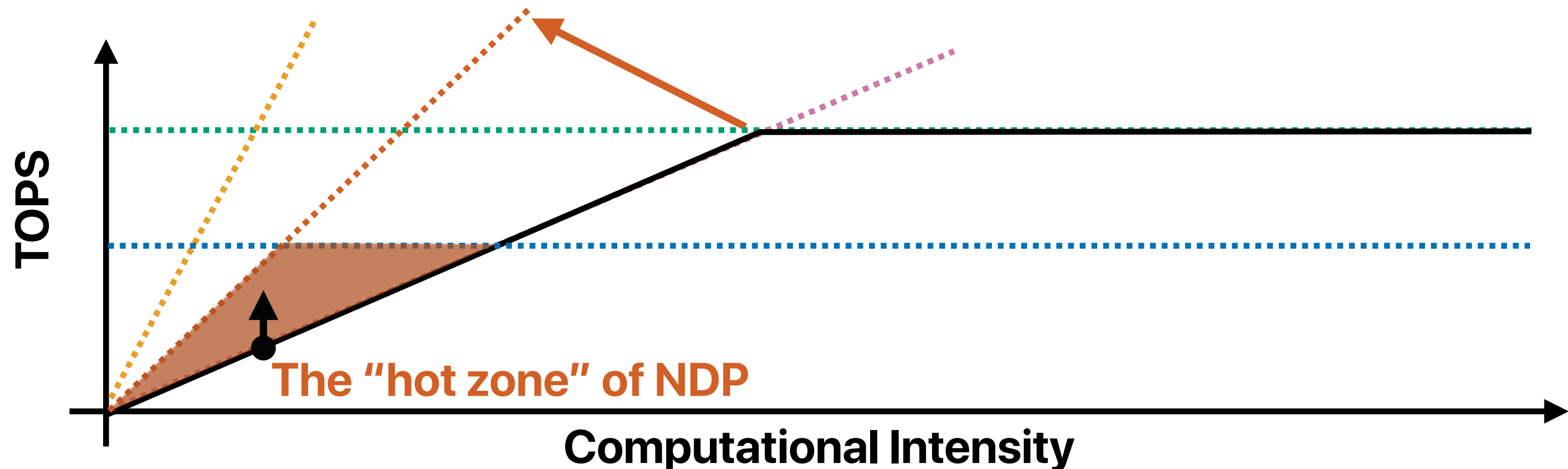
3:00

**Under what criteria do you think
NDP can improve performance?**

Criteria where NDP can help

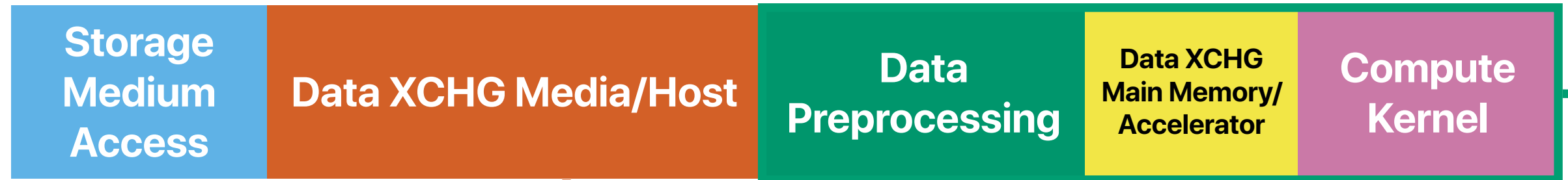
Reviewing the roofline model

$\min(\text{Peak OPS of target computing resource})$
 $\min(\text{Peak memory bandwidth} \times \text{computational intensity} \div \text{reduction of data volume})$
 $\min(\text{Peak OPS of target NDP device})$
 $\text{Peak device internal bandwidth} \times \text{computational intensity of NDP program}$



The “Winning Formula” of Near-Data Processing

Conventional Model



NDP Model



$$\frac{DataVolume_{raw}}{Bandwidth_{device}} + \frac{DataVolume_{raw} \times ComputationIntensity_{NDP}}{ComputationThroughput_{device}} + \frac{DataVolume_{afterNDP}}{Bandwidth_{device2host}} + \frac{DataVolume_{afterNDP} \times ComputationIntensity_{afterNDP}}{ComputationThroughput}$$

$$< = \frac{DataVolume_{raw}}{Bandwidth_{device2host}} + \frac{DataVolume_{raw} \times ComputationIntensity}{ComputationThroughput}$$

The “Winning Formula” of Near-Data Processing

Conventional Model



NDP Model



- The NDP computation can help offload computation
- The NDP computation can help reduce data volume
- The NDP computation can facilitate data preprocessing
- The NDP computation itself cannot be too slow

Not the cases of Near-Data Processing

Conventional Model



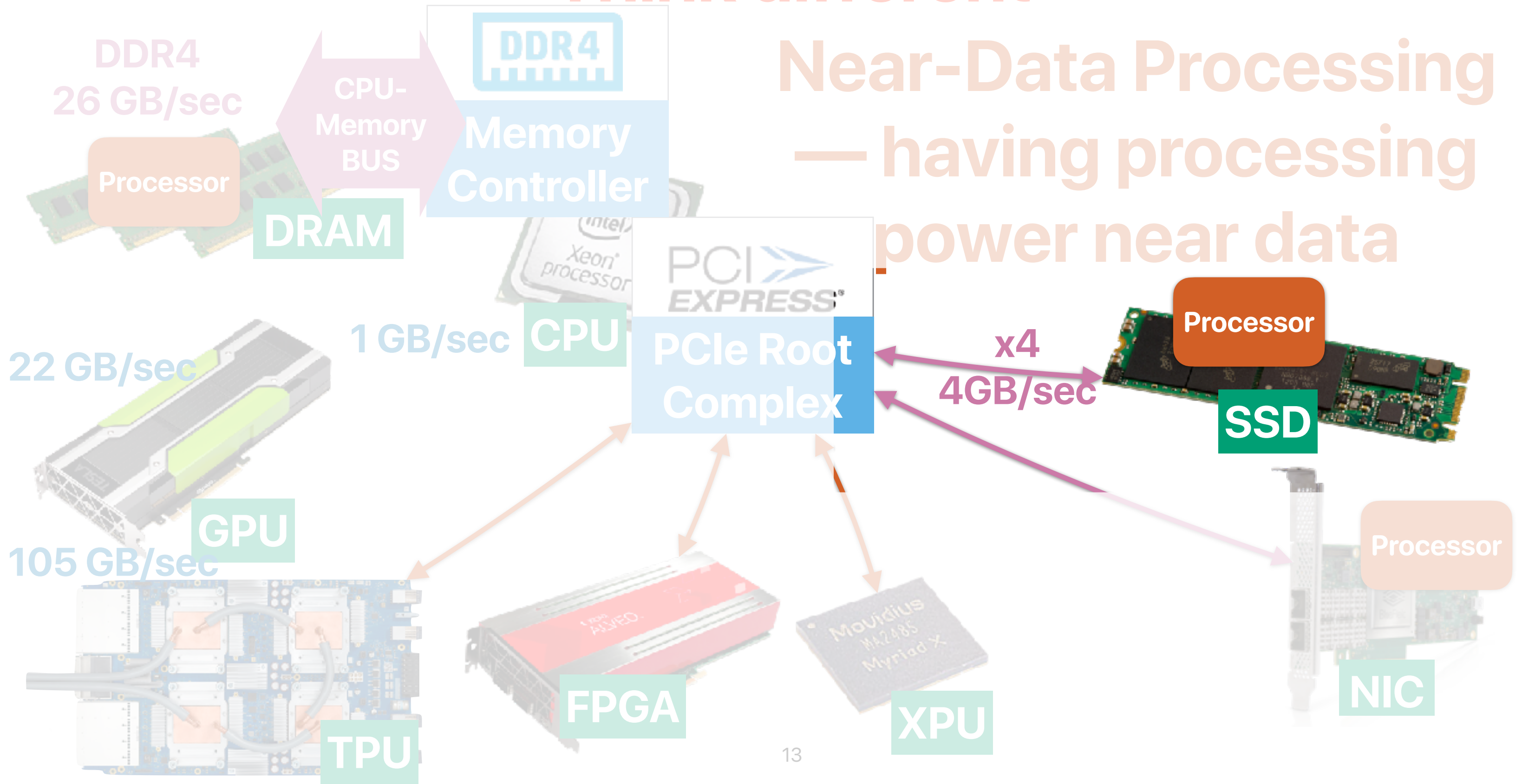
NDP Model



- The NDP processing exceeds NDP computing resources' capabilities
- The NDP processing does not help reduce data volume
- The NDP device itself does not offer rich internal bandwidth to the computing resource near-by

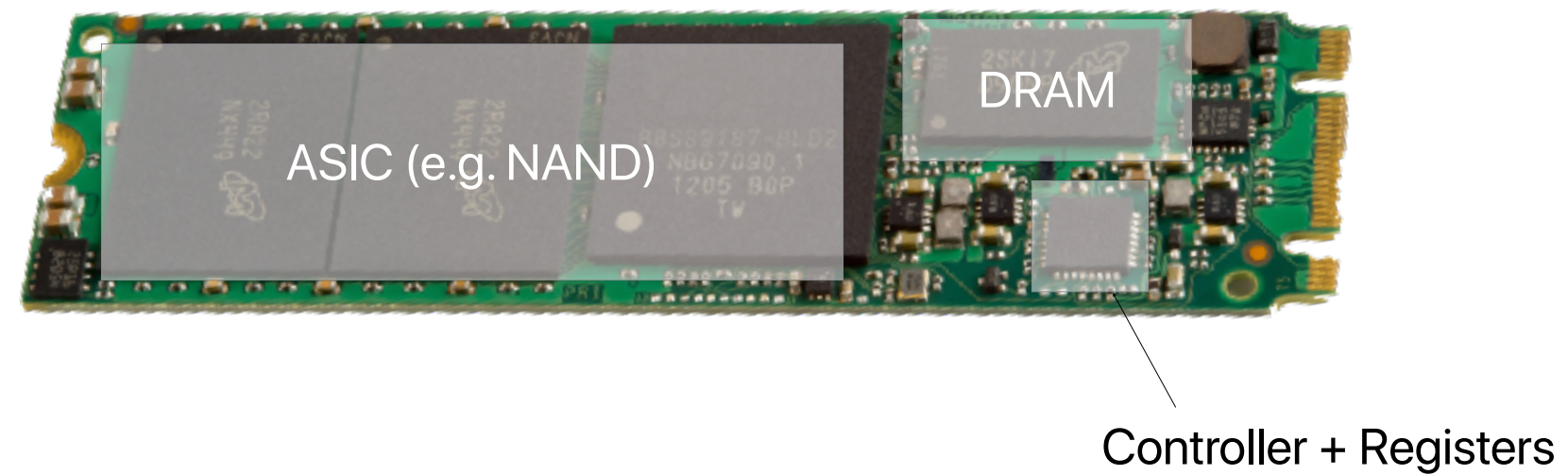
Think different

Near-Data Processing — having processing power near data



What's inside an SSD and why do we need them?

What an SSD looks like

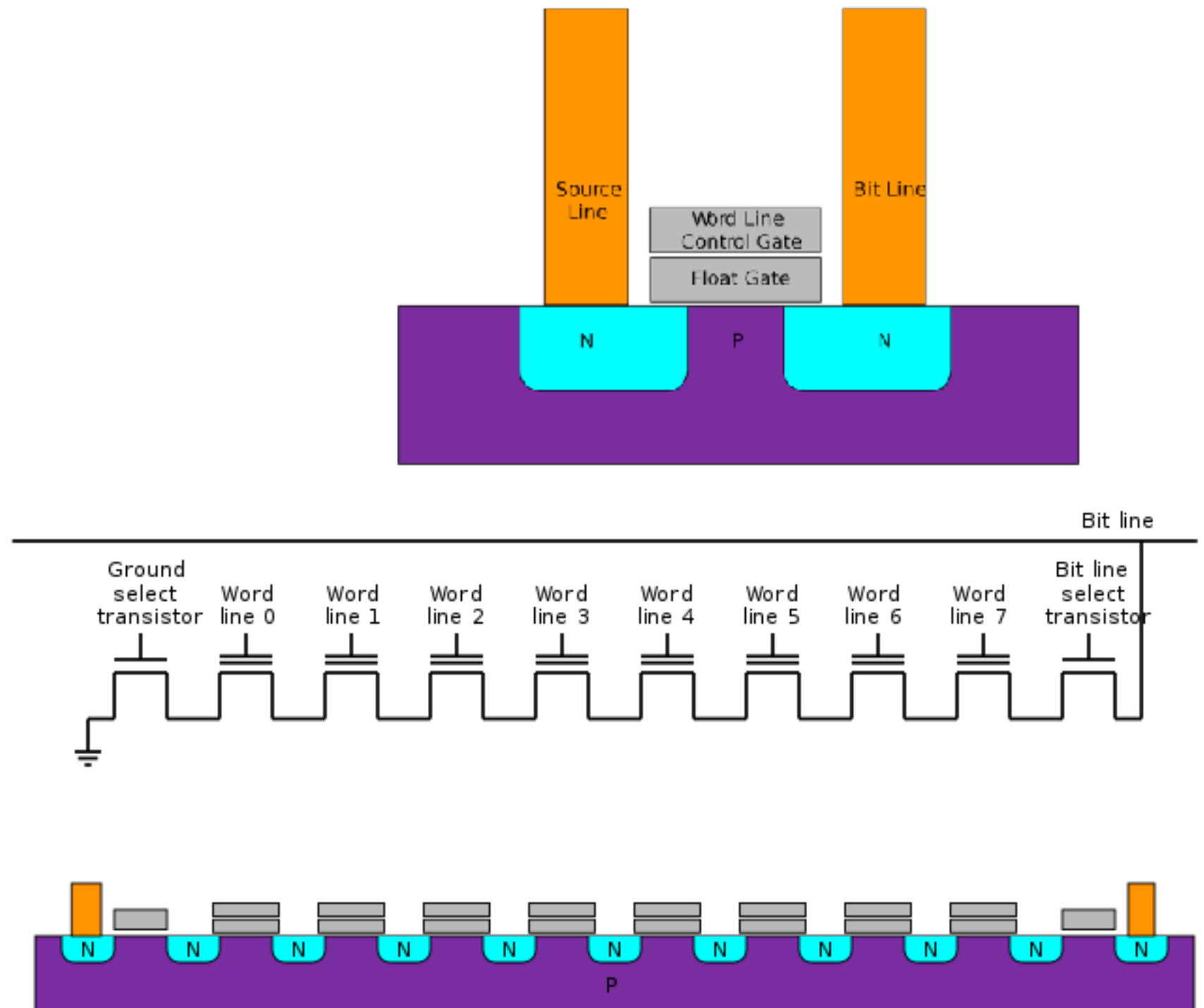


Why do we need controllers in SSDs?

- Interfacing with the interconnect
- Maintaining the block device abstraction
- Dealing with the “weird” device characteristics

Flash memory

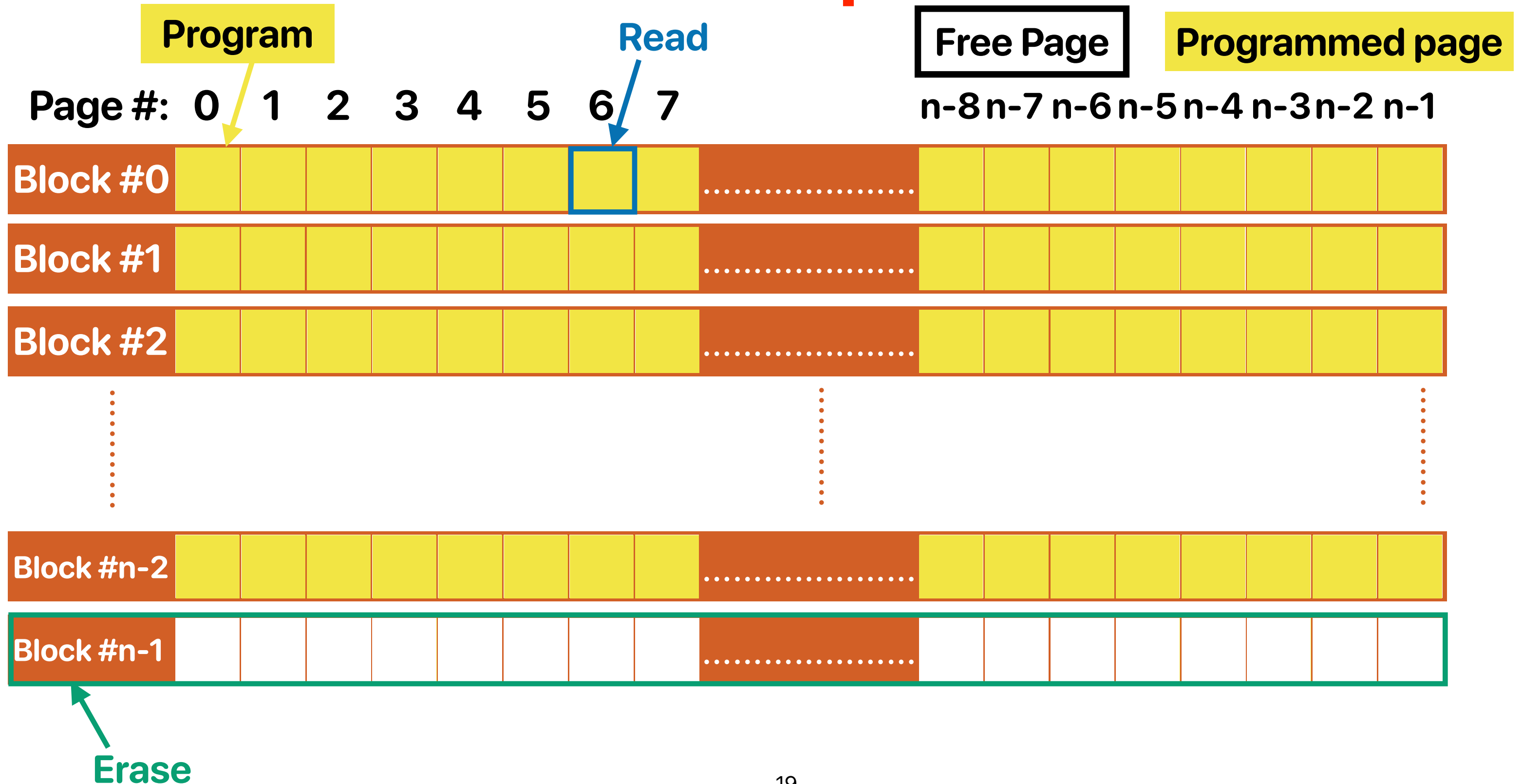
- Floating gate made by polycrystalline silicon trap electrons
- The voltage level within the floating gate determines the value of the cell
- The floating gates will wear out eventually



NAND flash is just odd!

- Modern SSDs are based on NAND flash memory
- Different operation granularities
 - Read/Program in pages
 - Erase in blocks (64-384 pages)
- Performance of operation varies
 - Read — tens of us
 - Program — hundreds of us
 - Erase — ms
- Limited erase cycles
 - Only 1000 times for "QLC"

Basic flash operations



Why we need a processor in SSDs?

- Modern SSDs are based on NAND flash memory
- Different operation granularities
 - Read/Program in pages
 - Erase in blocks (64-384 pages)
- Performance of operation varies
 - Read — tens of us
 - Program — hundreds of us
 - Erase — ms
- Limited erase cycles
 - Only 1000 times for "QLC"

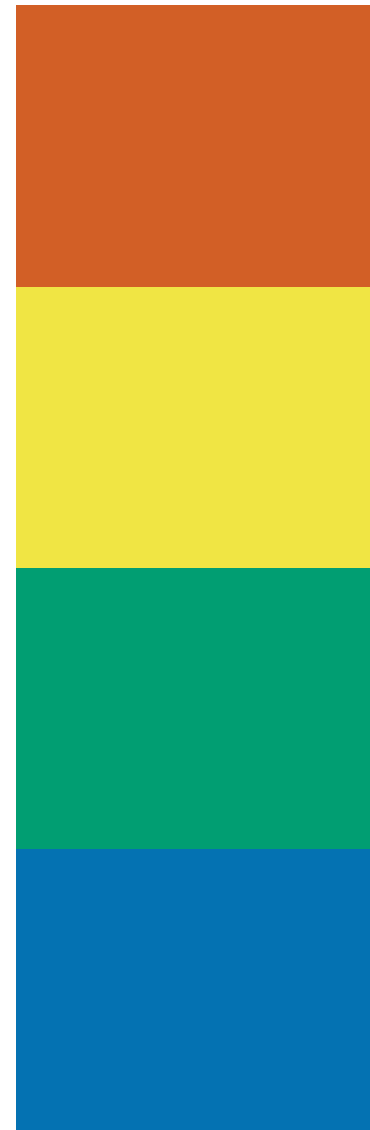
Types of Flash Chips

2 voltage levels,
1-bit



**Single-Level Cell
(SLC)**

4 voltage levels,
2-bit



**Multi-Level Cell
(MLC)**

8 voltage levels,
3-bit



**Triple-Level Cell
(TLC)**

16 voltage levels,
4-bit



**Quad-Level Cell
(QLC)**

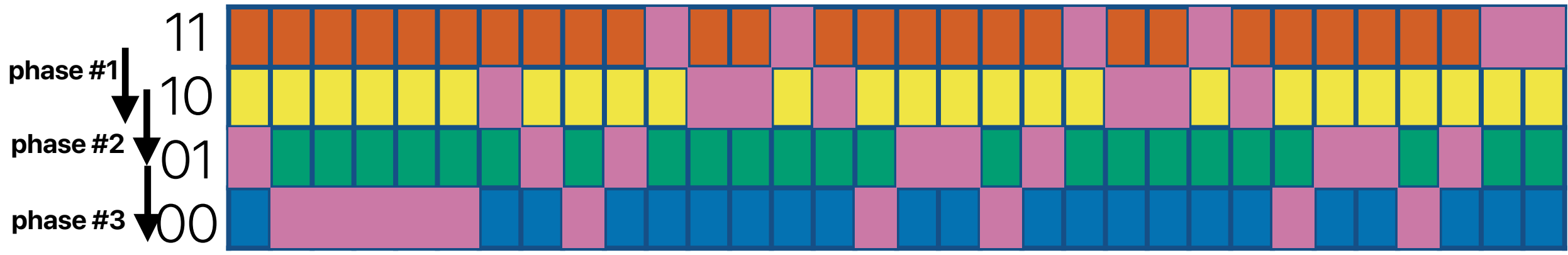
Programming in MLC

4 voltage levels,
2-bit



Multi-Level Cell
(MLC)

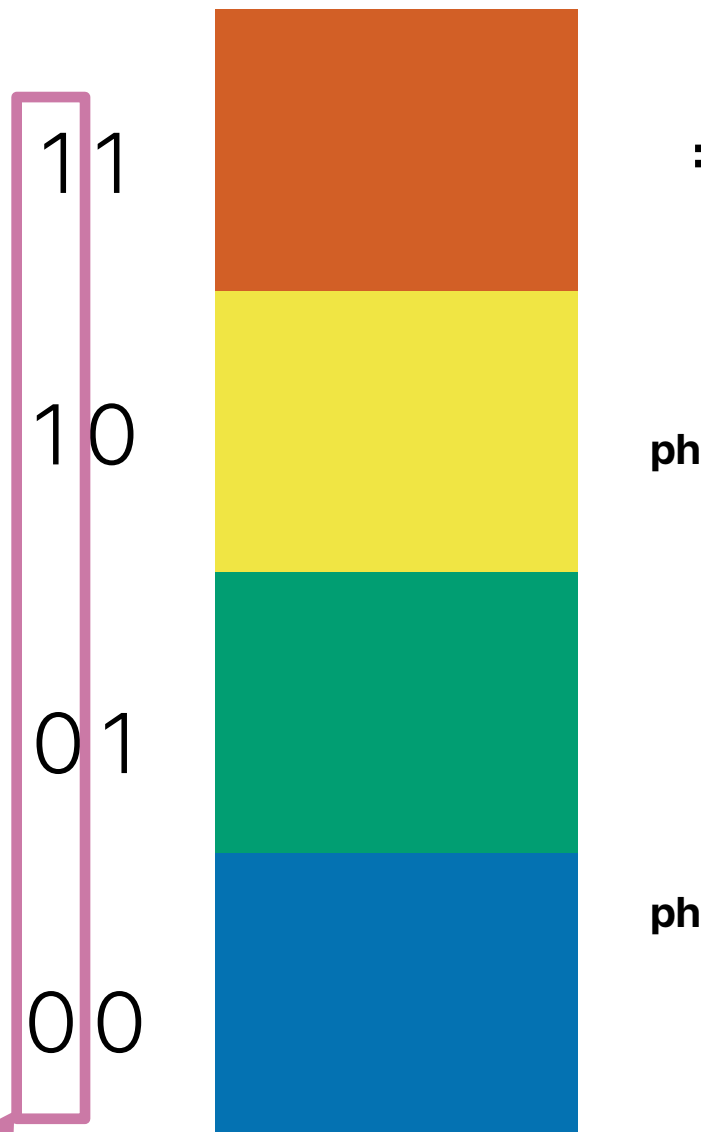
3.1400000000000000001243449787580
= 0x40091EB851EB851F
= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111



3 Cycles/Phases to finish programming

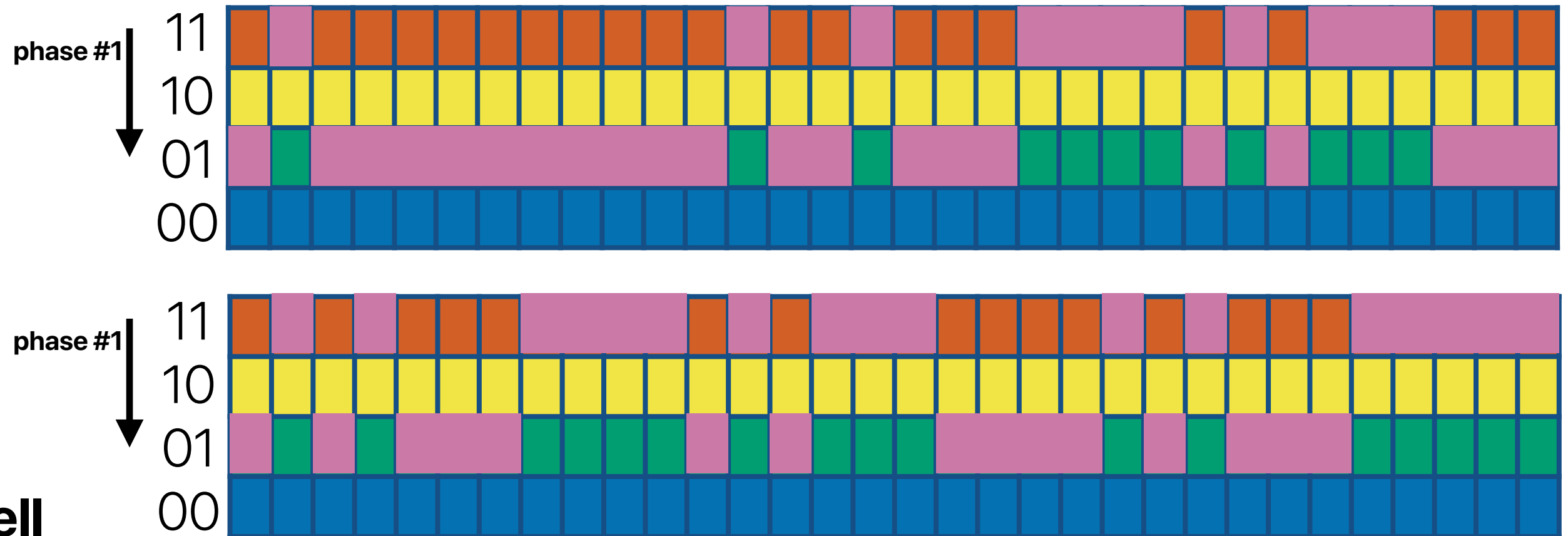
Programming in MLC

4 voltage levels,
2-bit



**Multi-Level Cell
(MLC)**

3.140000000000000000001243449787580
= 0x40091EB851EB851F
= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111



1 Phase to finish programming the first page!

Programming the 2nd page in MLC

4 voltage levels,
2-bit

2nd page

3.140000000000000000001243449787580

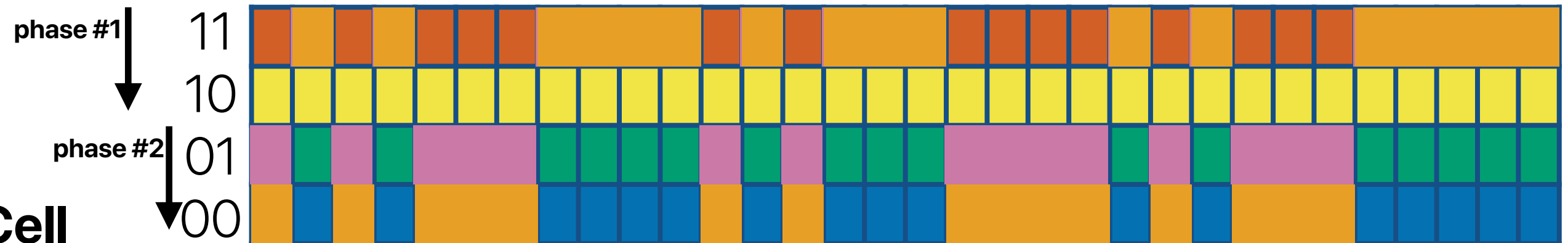
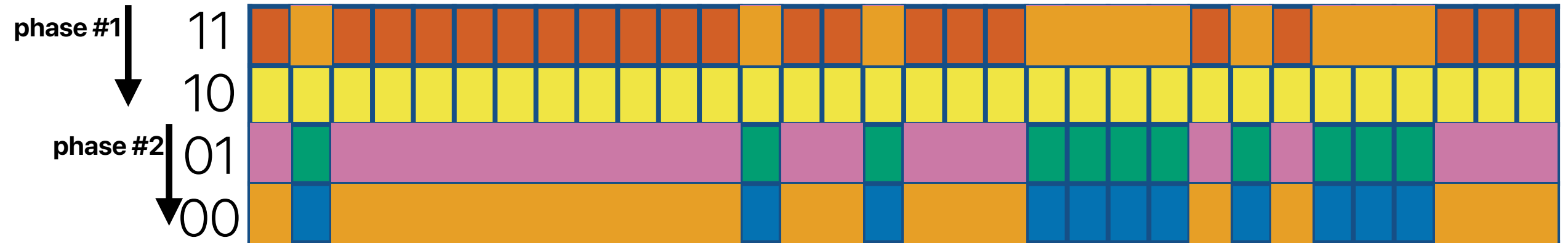
= 0x40091EB851EB851F

= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111

= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111



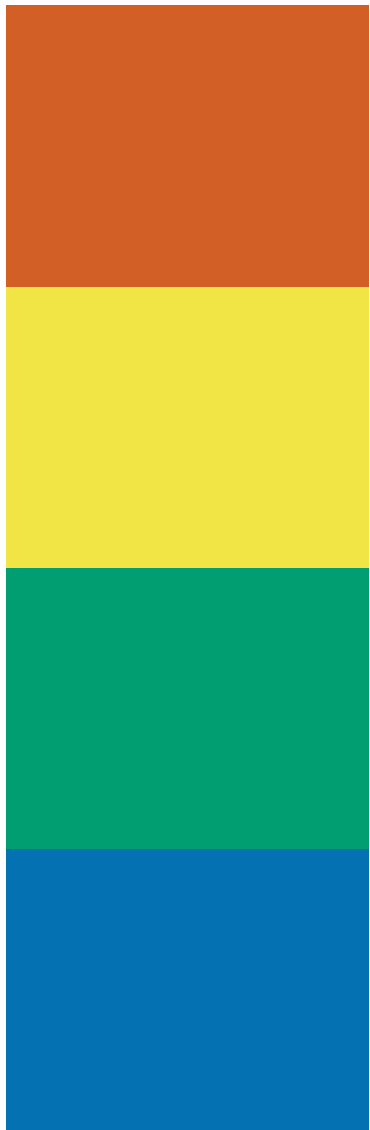
Multi-Level Cell
(MLC)



2 Phase to finish programming the second page!

Optimizing 1st Page Programming in MLC

4 voltage levels,
2-bit

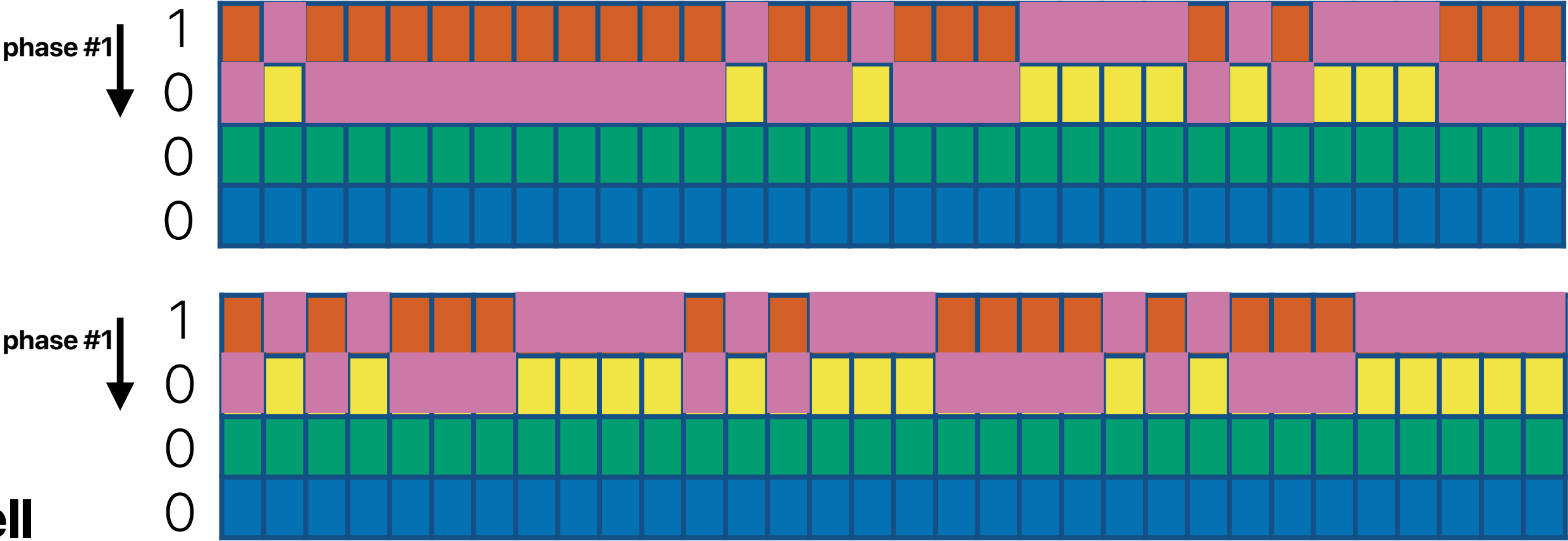


Multi-Level Cell
1st page (MLC)

3.1400000000000000001243449787580

= 0x40091EB851EB851F

= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111



1 Phase to finish programming the first page!

25 — the phase is shorter now

2nd Page Programming in MLC

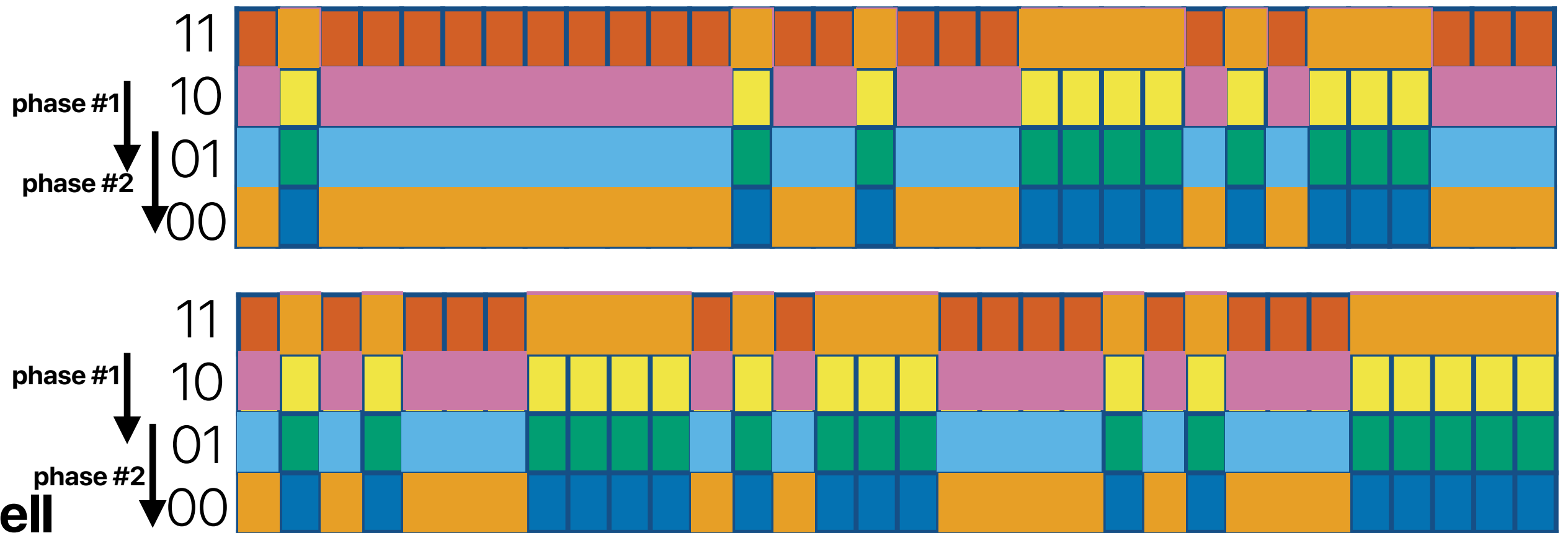
2nd page 4 voltage levels,
2-bit

[illegible]

= 0x40091EB851EB851F

= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111

= 01000000 00001001 00011110 10111000 01010001 11101011 10000101 00011111



2 Phase to finish programming the second page!

QLC = More Density Per NAND Cell



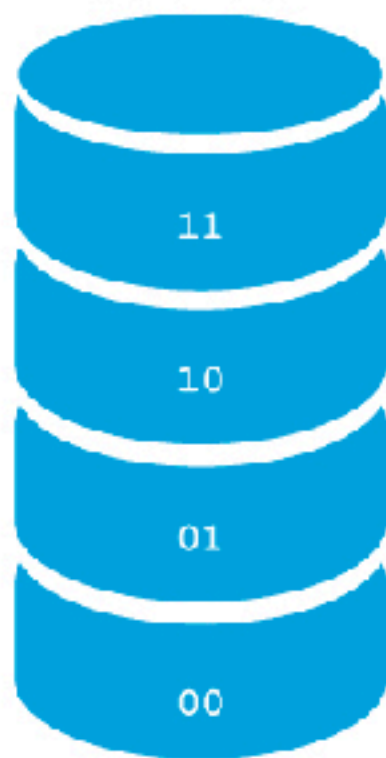
SLC



1 Bit Per Cell

First SSD NAND technology

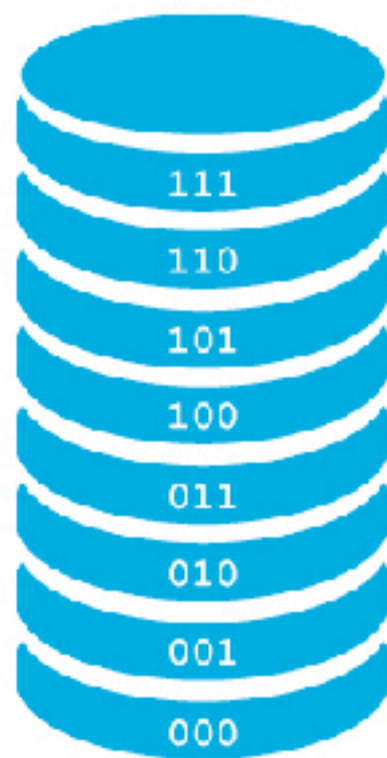
MLC



2 Bits Per Cell

100% increase

TLC



3 Bits Per Cell

50% increase

QLC



4 Bits Per Cell

33% increase



100K P/E Cycles
(at technology introduction)

10K P/E Cycles

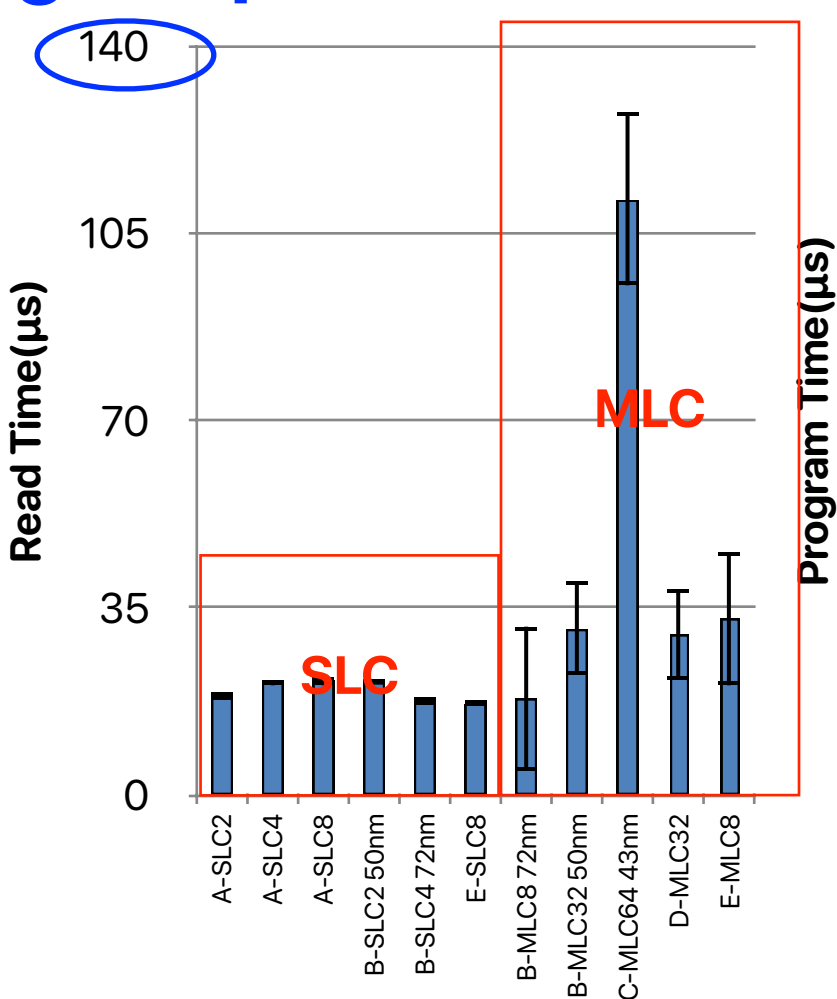
3K P/E Cycles

1K P/E Cycles

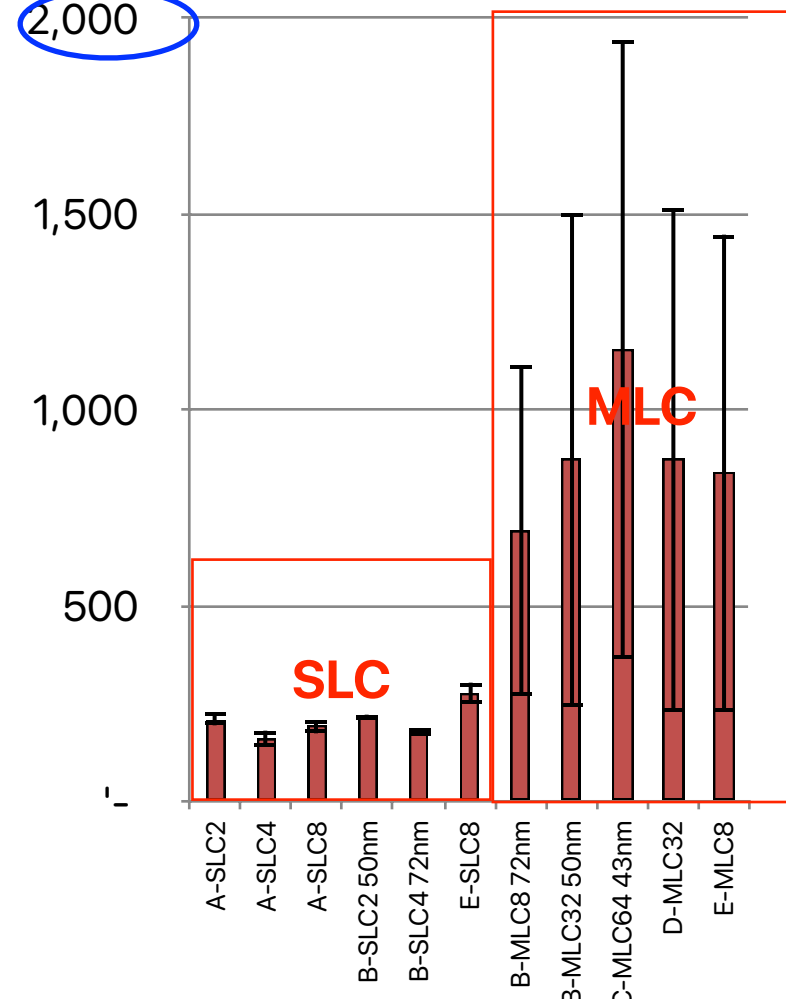
Fewer writes per cell

Not a good practice

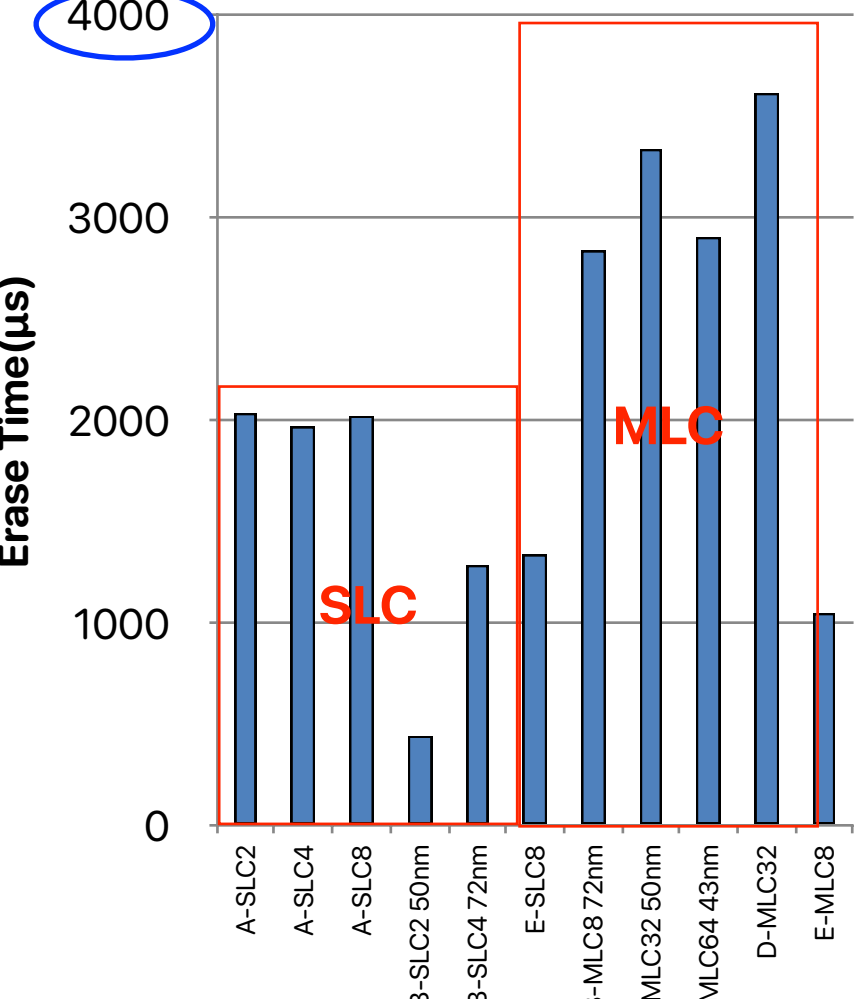
Flash performance



Reads:
less than 150us



Program/write:
less than 2ms

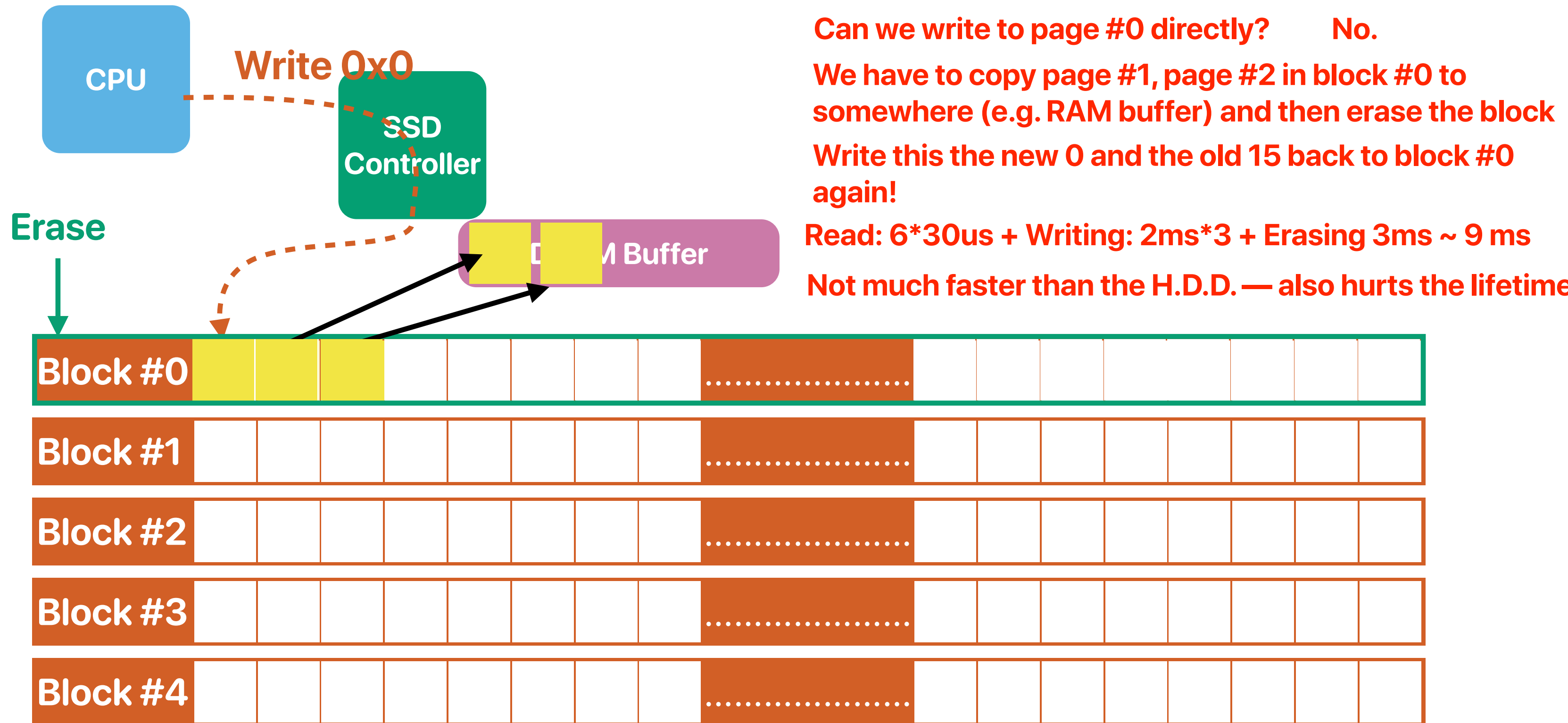


Erase:
less than 3.6ms

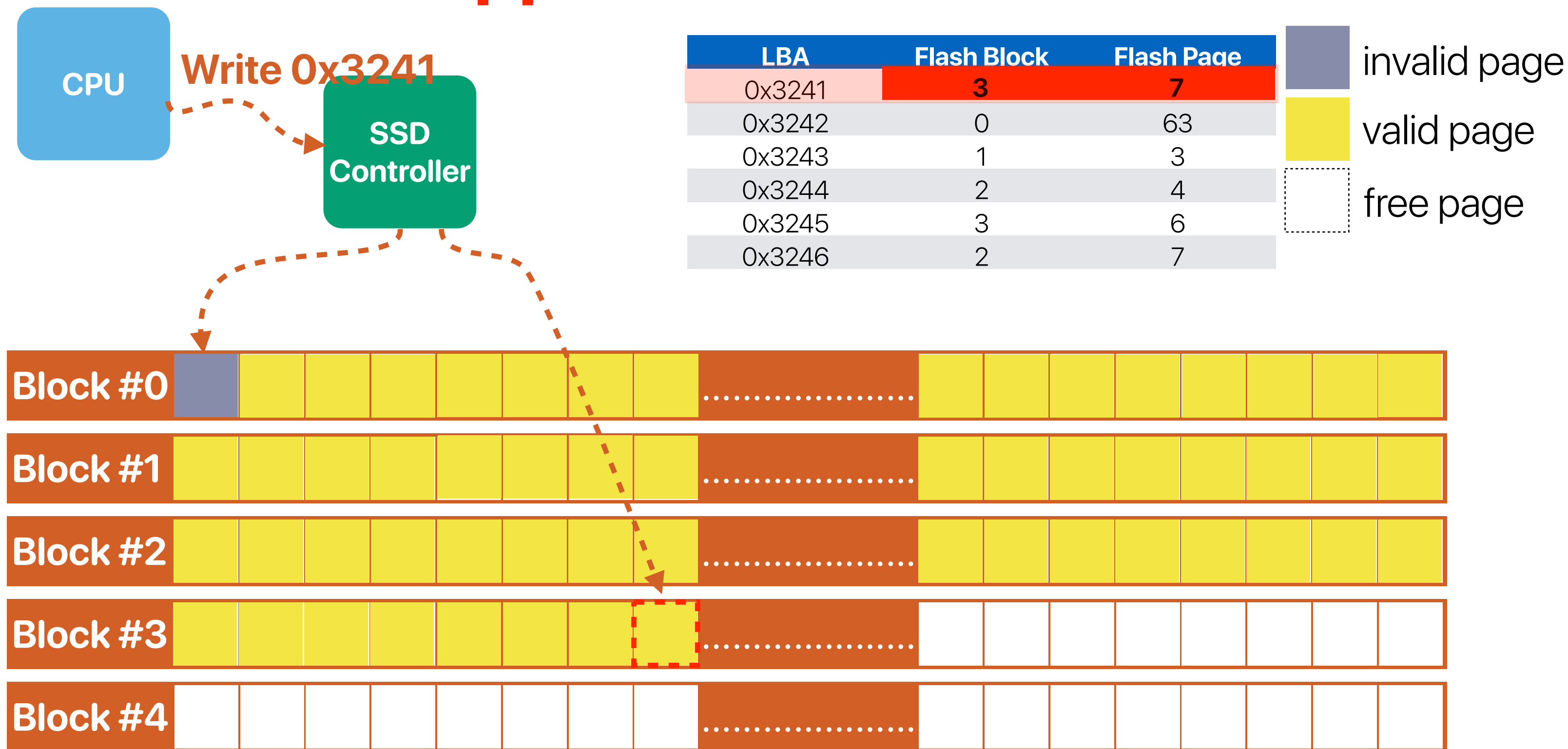
Similar relative performance for reads, writes and erases

Laura M. Grupp, Adrian M. Caulfield, Joel Coburn, Steven Swanson, Eitan Yaakobi, Paul H. Siegel, and Jack K. Wolf.
Characterizing flash memory: anomalies, observations, and applications. In MICRO 2009.

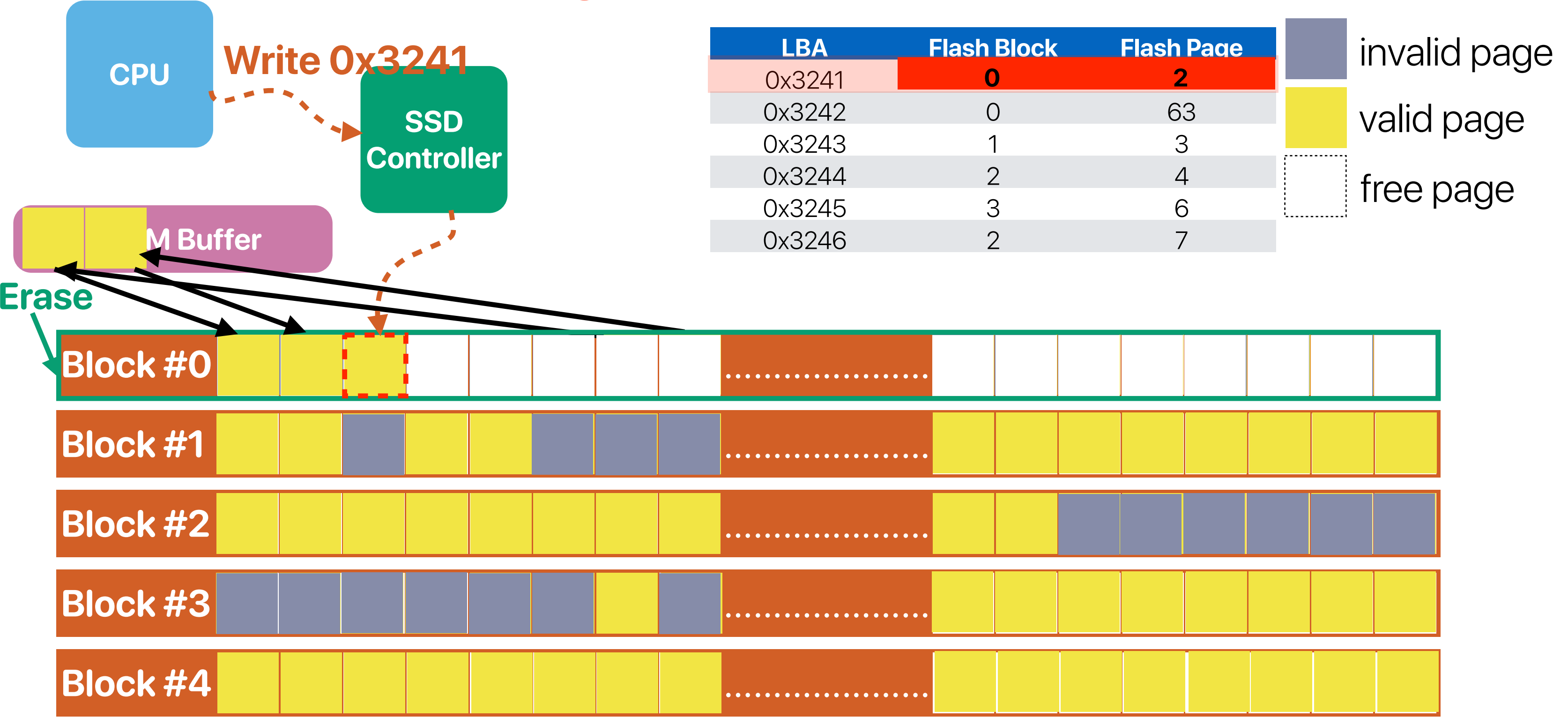
What happens on a write if we use the same abstractions as H.D.D.



What happens on a write with FTL



Garbage Collection in FTL



Multiple channels to improve bandwidth

Flashtec™ NVMe2032 and NVMe2016 Controllers

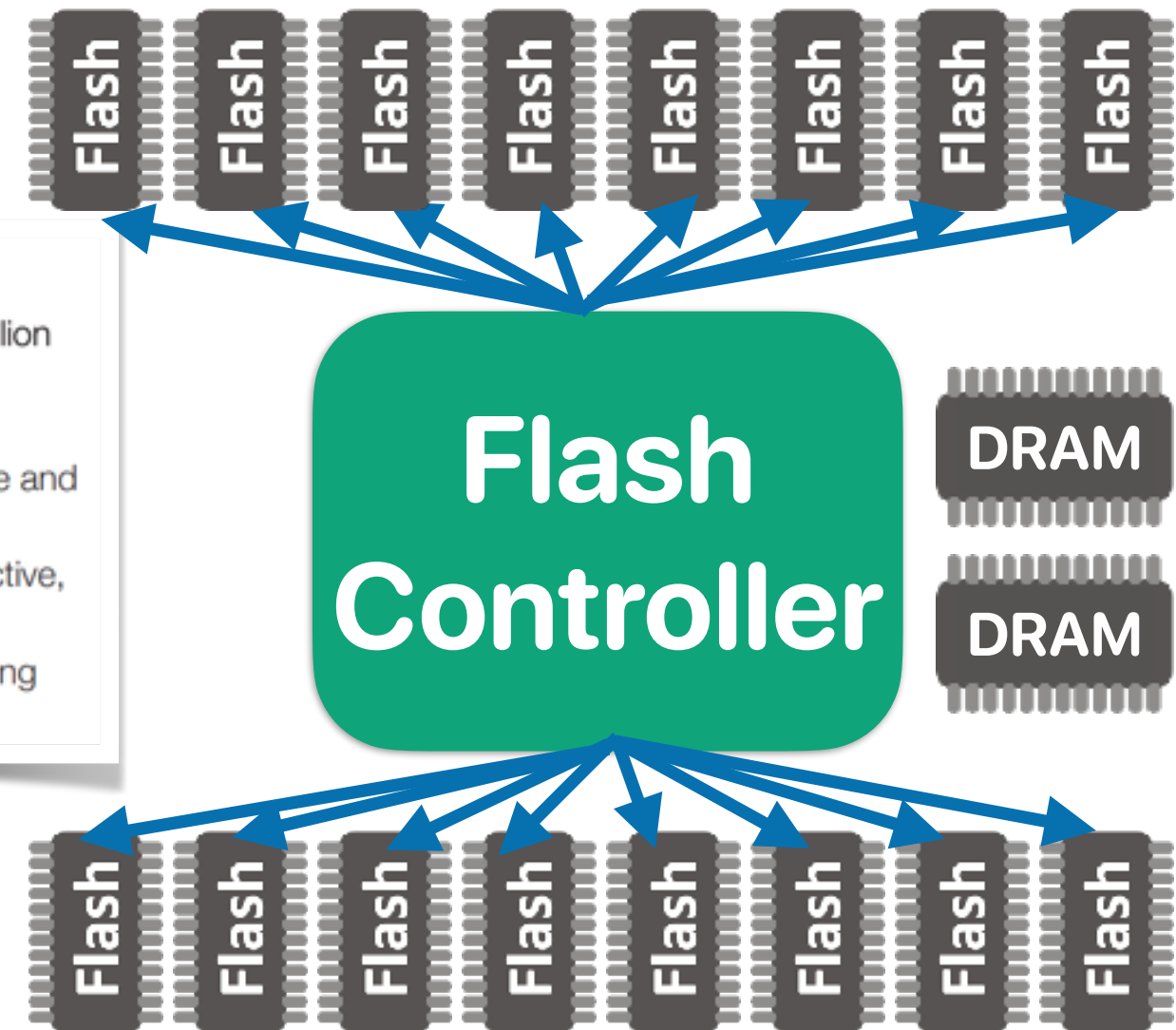
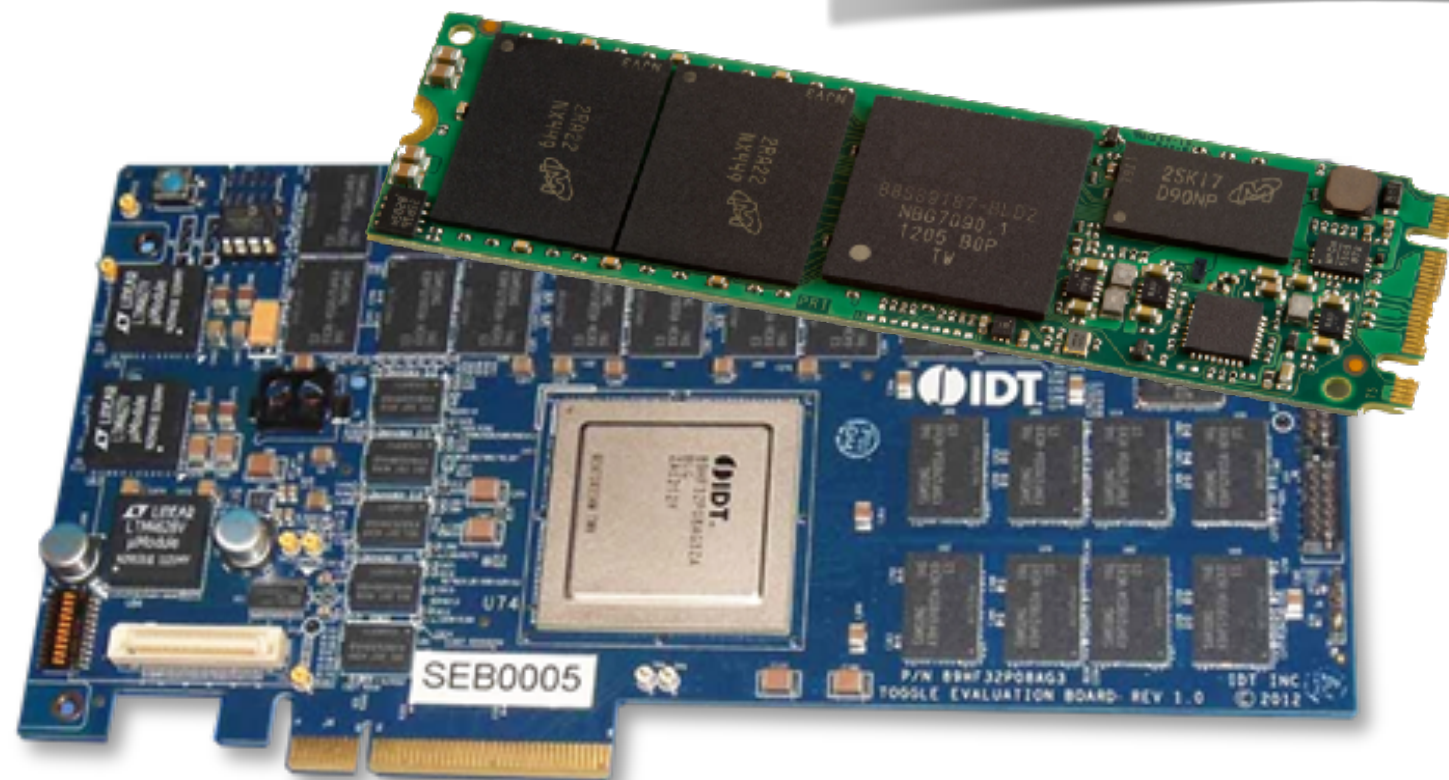
32- and 16-Channel PCIe Flash Controller Products

Summary

The Flashtec™ 2nd generation NVMe Controller Family enables the server and data centers to realize the highest performance SSDs utilizing new technologies. Combining world-class capacity and flexibility, the Flashtec is the reliable choice. The Flashtec NVMe2032 and NVMe2016 controllers support the Express (NVMe) host interface and are optimized for high-performance operations, performing all flash management operations on-chip and processing and memory resources.

Features

- Flashtec NVMe2032 controller can achieve up to 1 million random read IOPS on 4 KB operations
- Up to 20 TB Flash capacity using 256 GB Flash
- SLC, MLC, Enterprise MLC, and TLC Flash with toggle and ONFI interface
- PCIe Gen 3 x8 or dual independent PCIe Gen 3 x4 (active, active/standby) host interface
- 16 and 32 independent Flash channels, each supporting up to 8 CE



Electrical Computer Science Engineering

277

つくづく

