
All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE |

All Tomorrow's Memories

Bruce Jacob

**Keystone Professor
University of Maryland**



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 2

Stealth Revolution



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 2

Stealth Revolution



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 2

Stealth Revolution



Fusion-ish SSD



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 2

Stealth Revolution



Fusion-ish SSD

All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 3

Stealth Revolution



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 3

Stealth Revolution



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 3

Stealth Revolution

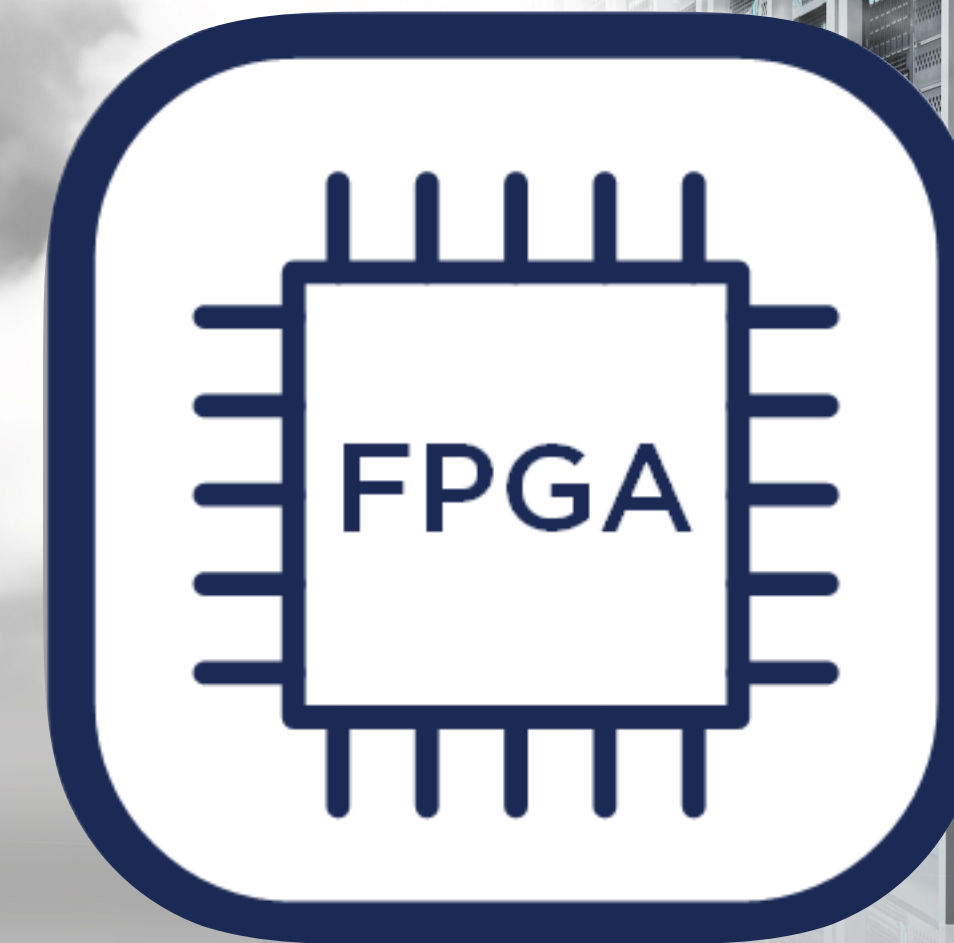


All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

Stealth Revolution



All Tomorrow's
Memories

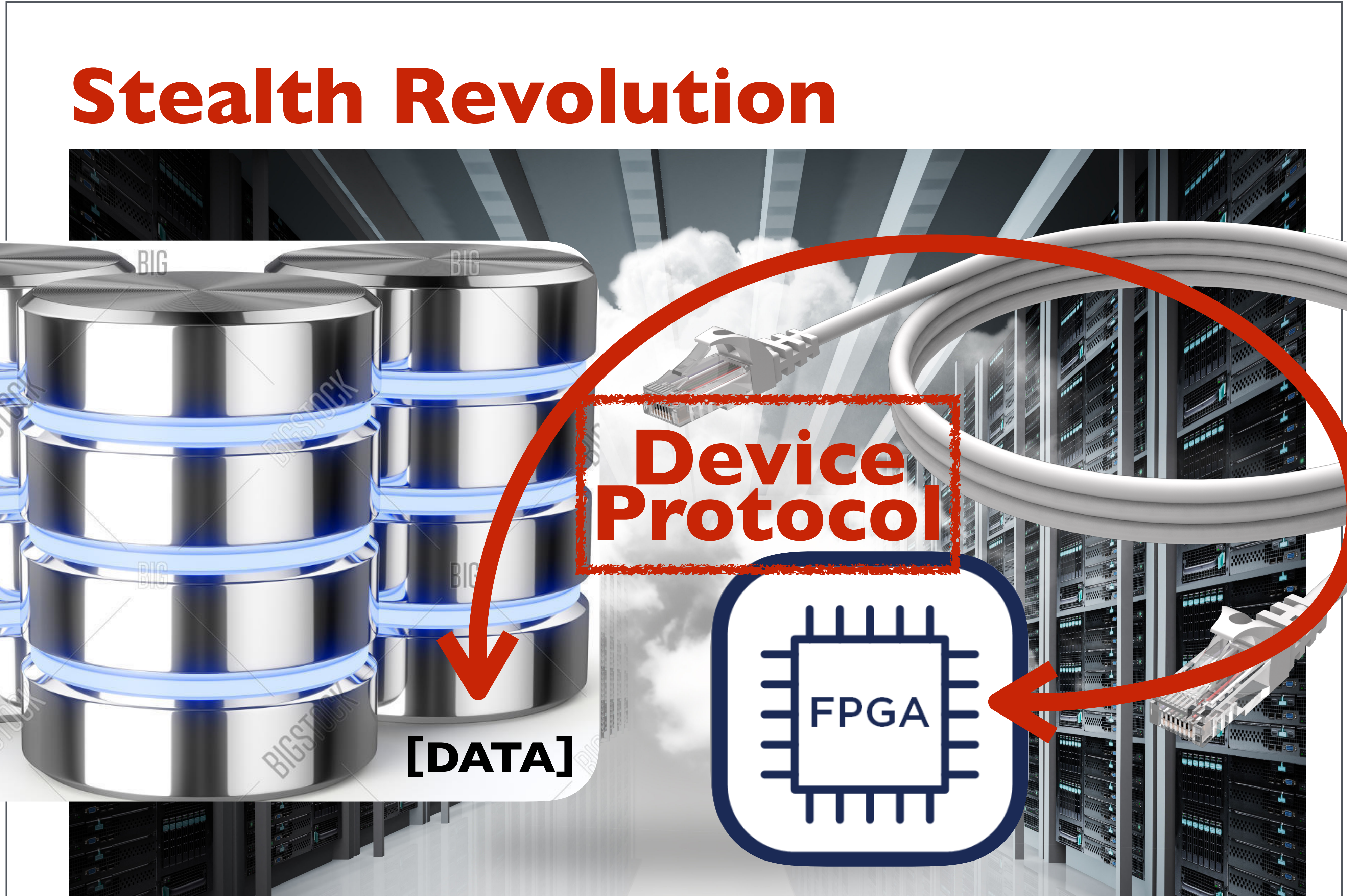
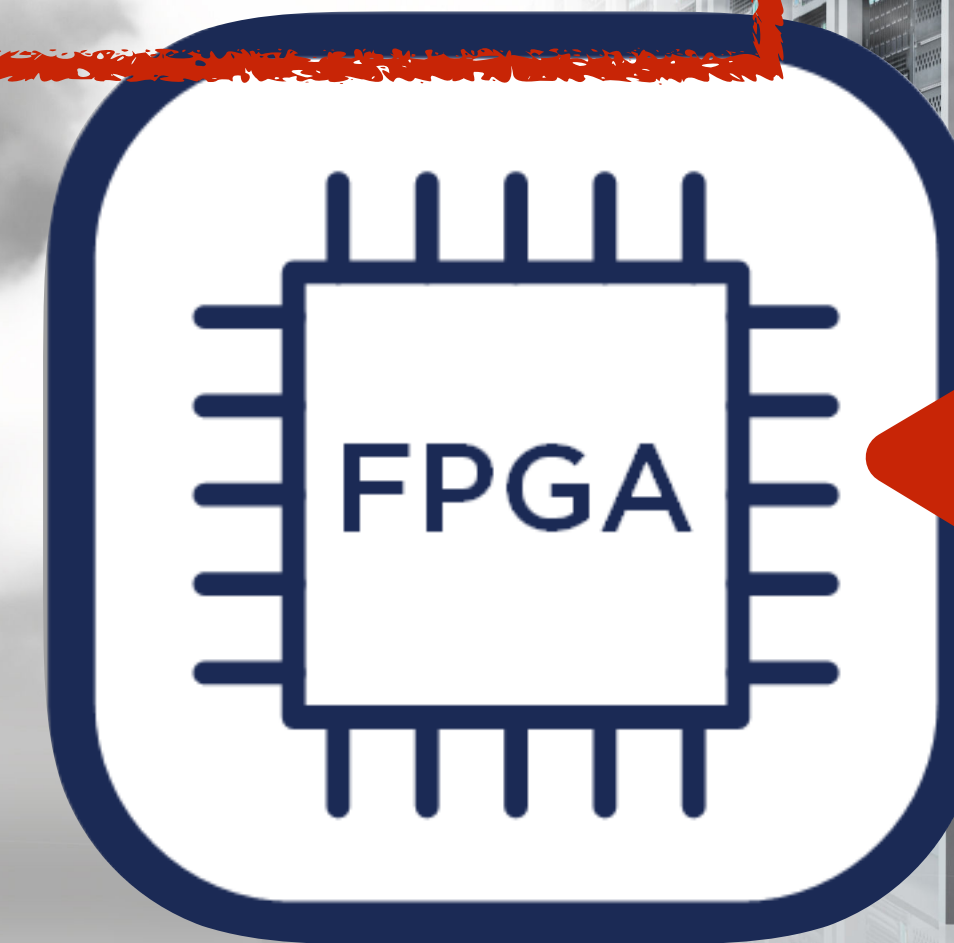
Bruce Jacob

University of
Maryland

Stealth Revolution



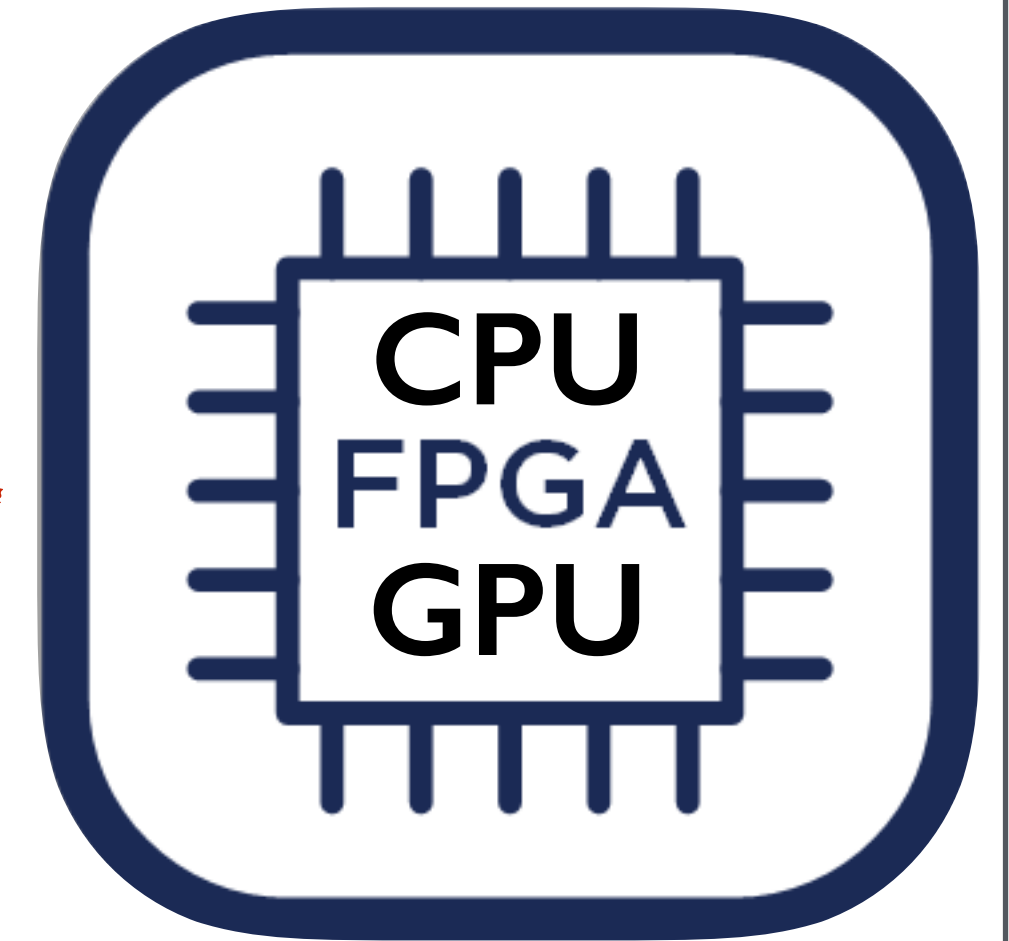
Device
Protocol



All Tomorrow's
Memories

Bruce Jacob

Stealth Revolution



Question: where is data?

Answer: <data owner>

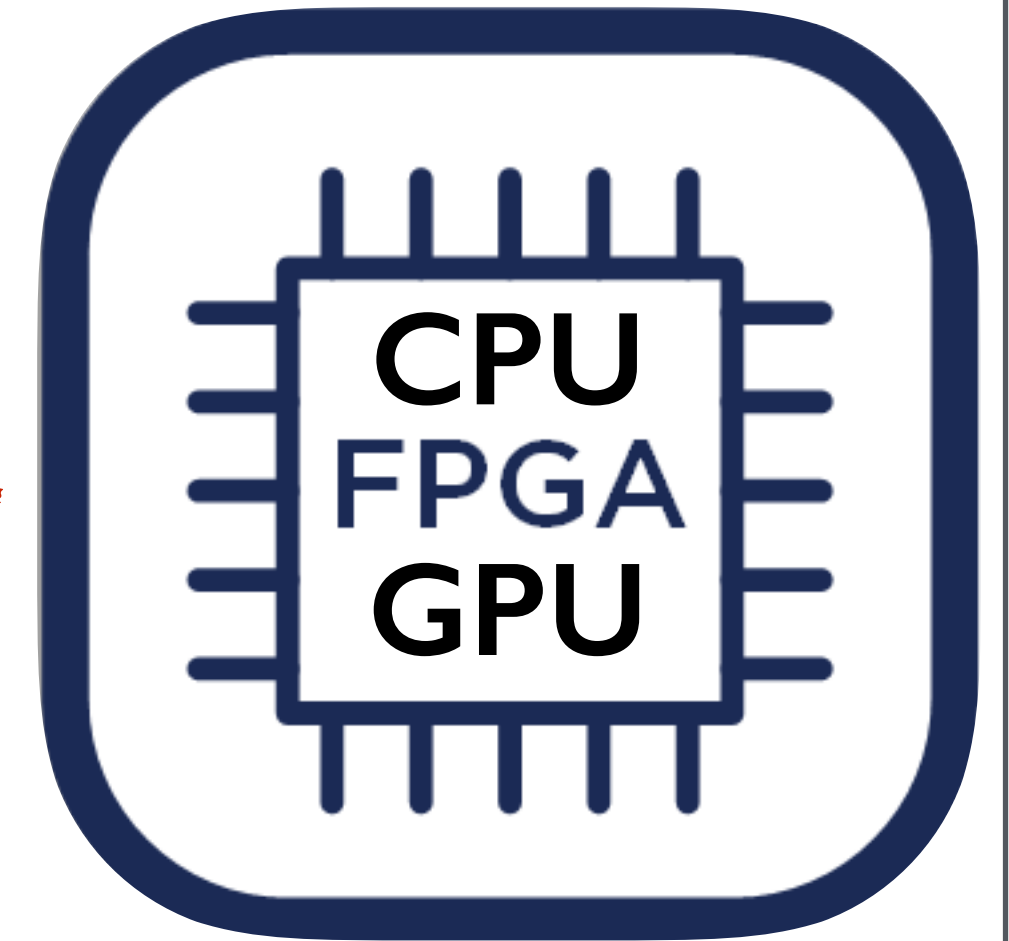
To data owner: read X

To requester: mem[X]

All Tomorrow's
Memories

Bruce Jacob

Stealth Revolution



Question: where is data?

Answer: <data owner>

To data owner: read X

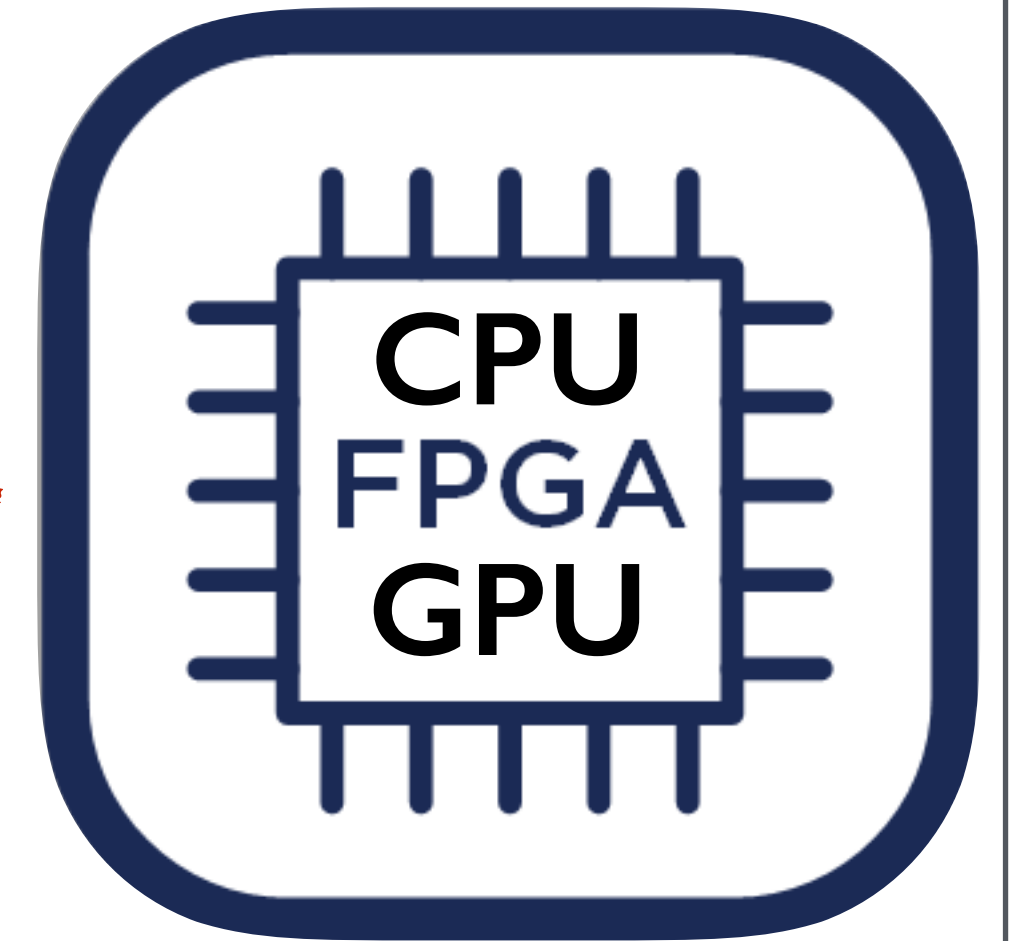
To requester: mem[X]

**Major implications
for OS and applications**

All Tomorrow's
Memories

Bruce Jacob

Stealth Revolution



Question: where is data?

Answer: <data owner>

To data owner: read X

To requester: mem[X]

*(esp. considering
Fusion-io like
capabilities)*

**Major implications
for OS and applications**

Background: Wish List

Fine-Grained Access

Bandwidth

Capacity

Low Power

Nonvolatility

DRAM -

HBM/HMC*

Flash, 3DXP, RRAM,
PCM, etc - **NVMM***

HBNV*

* Things we did and/or are doing now (I'll cover in talk)

Background: Wish List

Fine-Grained Access

Bandwidth

Capacity

Low Power

Nonvolatility

DRAM -

HBM/HMC*

Flash, 3DXP, RRAM,
PCM, etc - **NVMM***

HBNV*

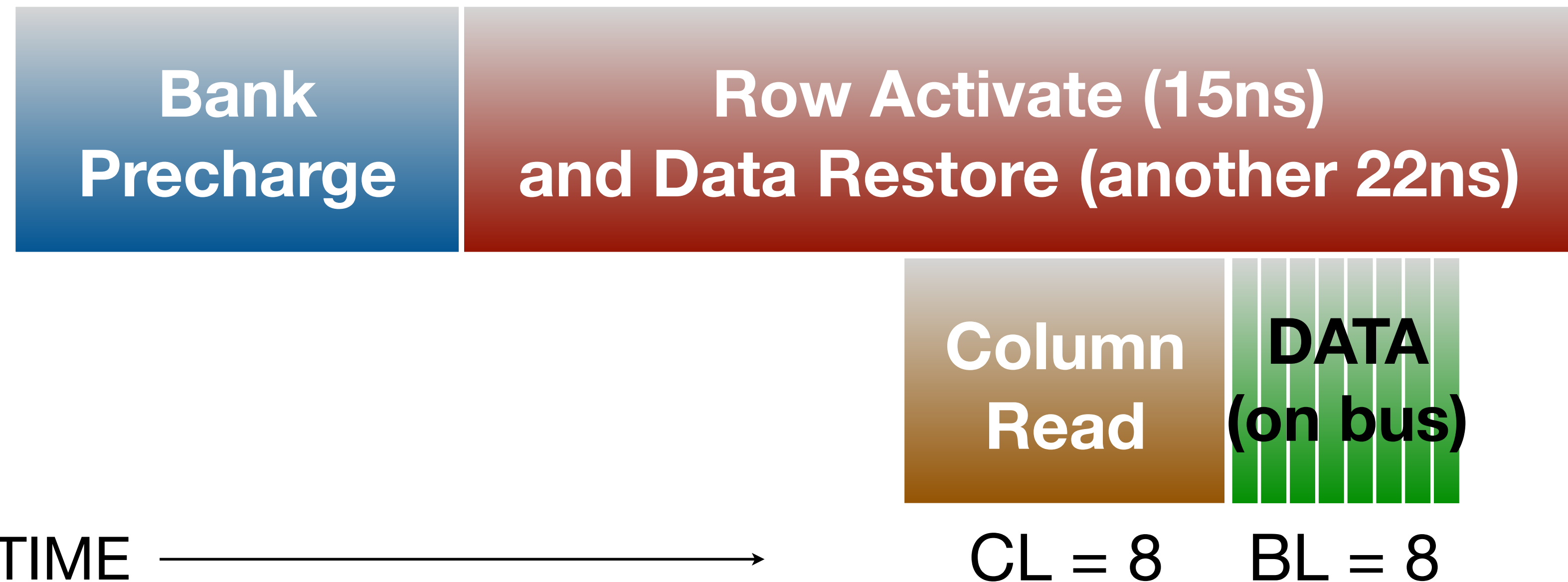
**Major implications
for OS* and applications**

* Things we did and/or are doing now (I'll cover in talk)

Background: Memory Latency

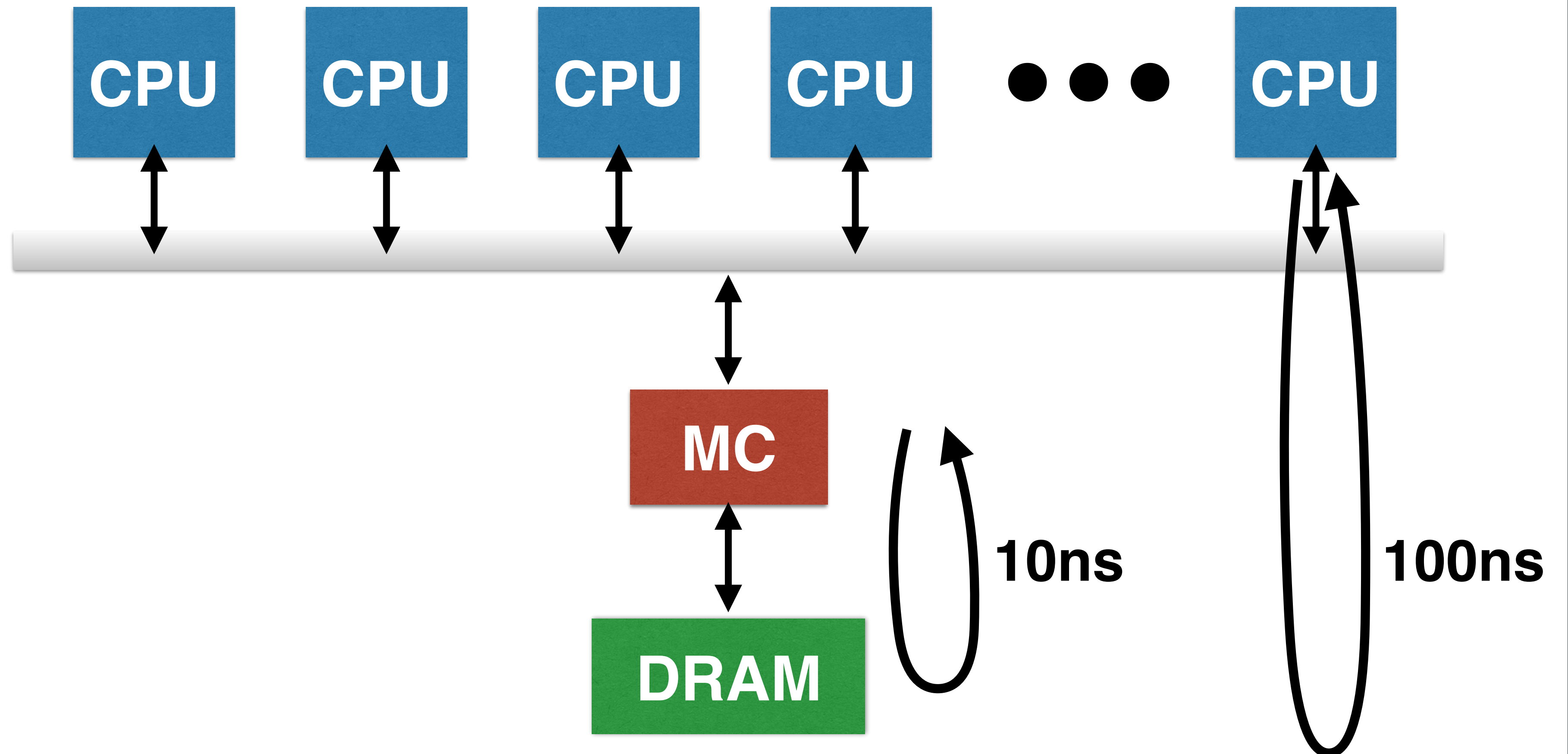
$t_{RP} = 15\text{ns}$

$t_{RCD} = 15\text{ns}$, $t_{RAS} = 37.5\text{ns}$



Cost of access is high; requires **significant effort** to amortize this over the (increasingly short) payoff.

Background: Memory Latency

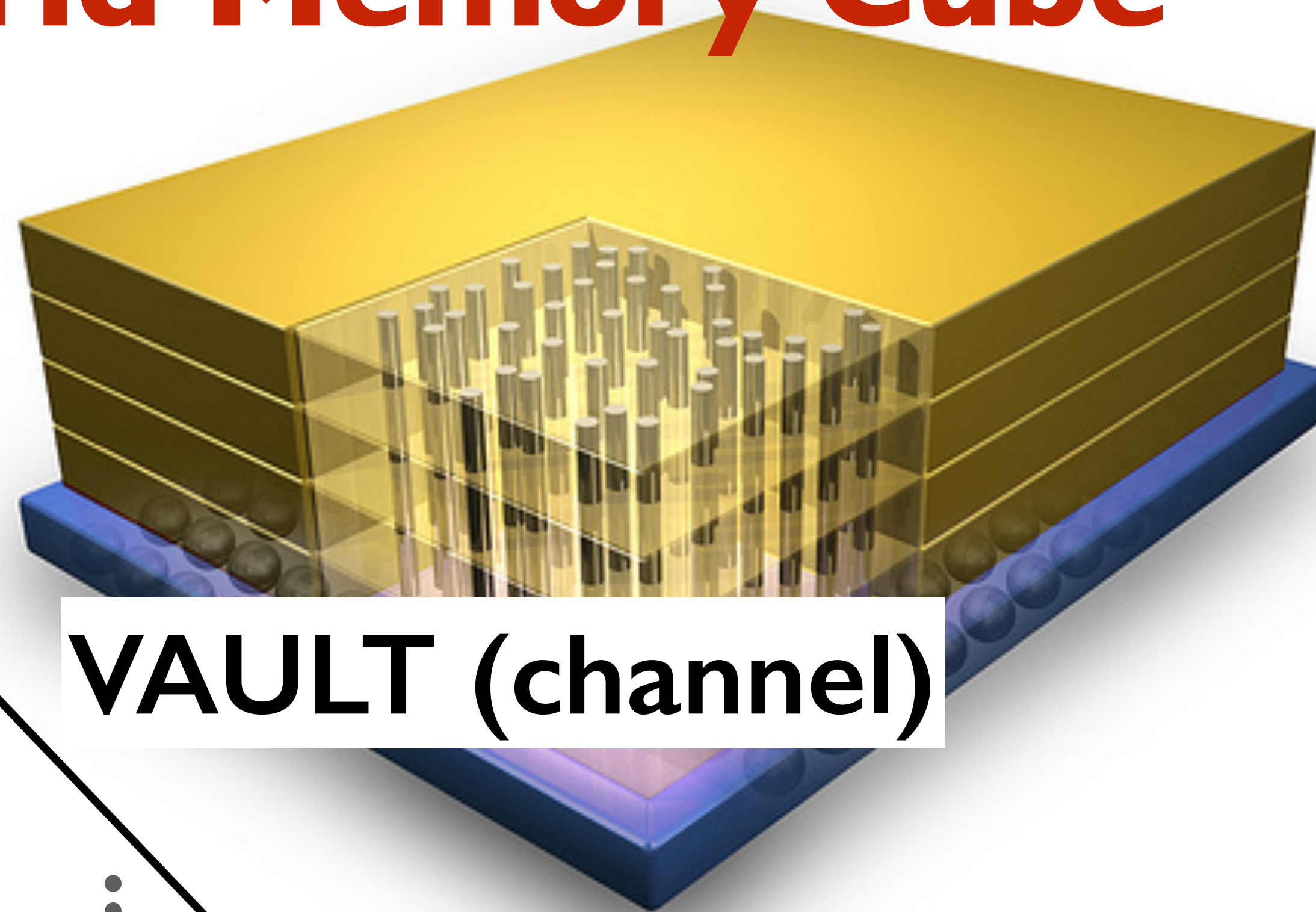
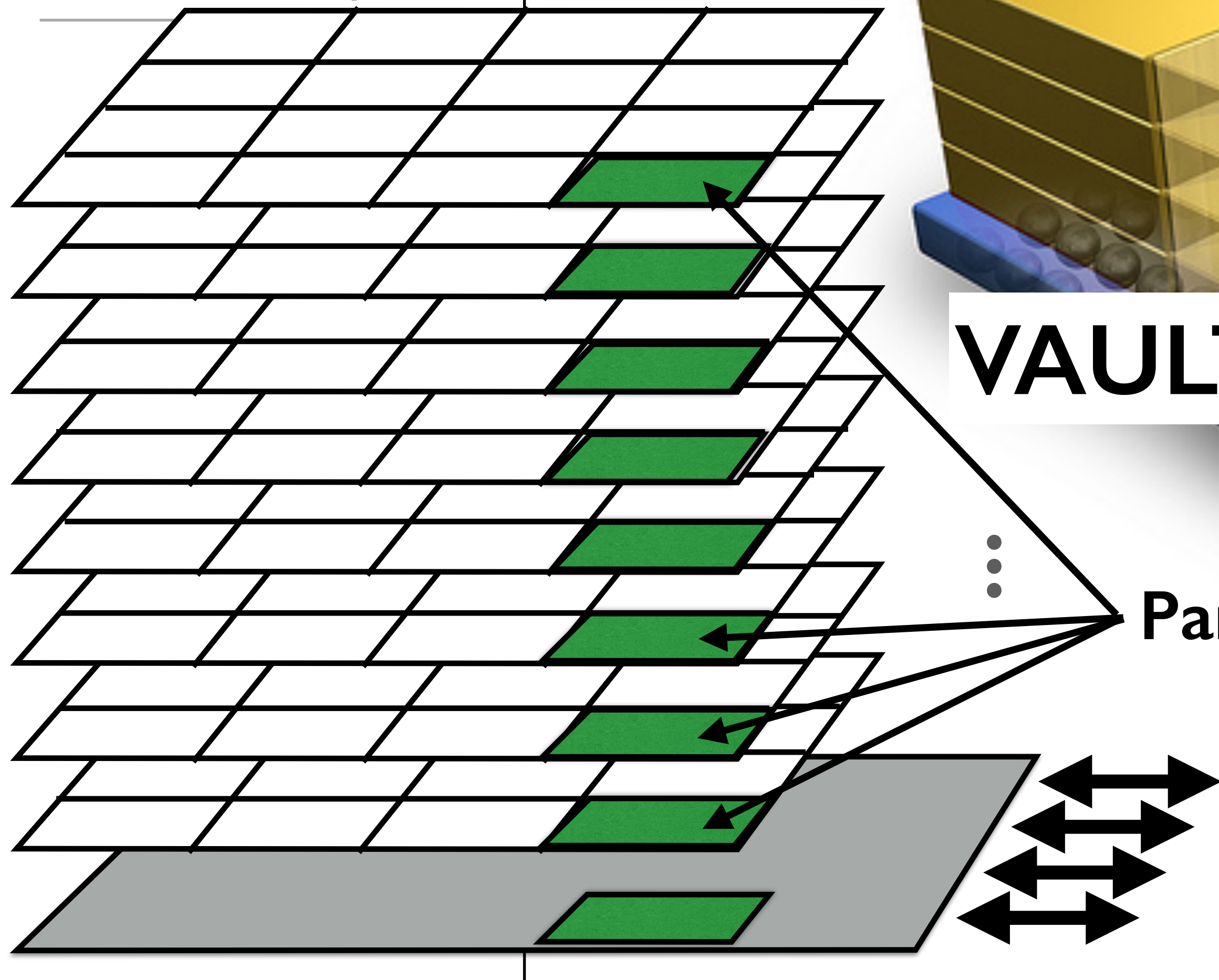


All Tomorrow's
Memory Systems

Bruce Jacob

University of
Maryland

Hybrid Memory Cube



Off-chip: high
speed SerDes
and generic
protocol

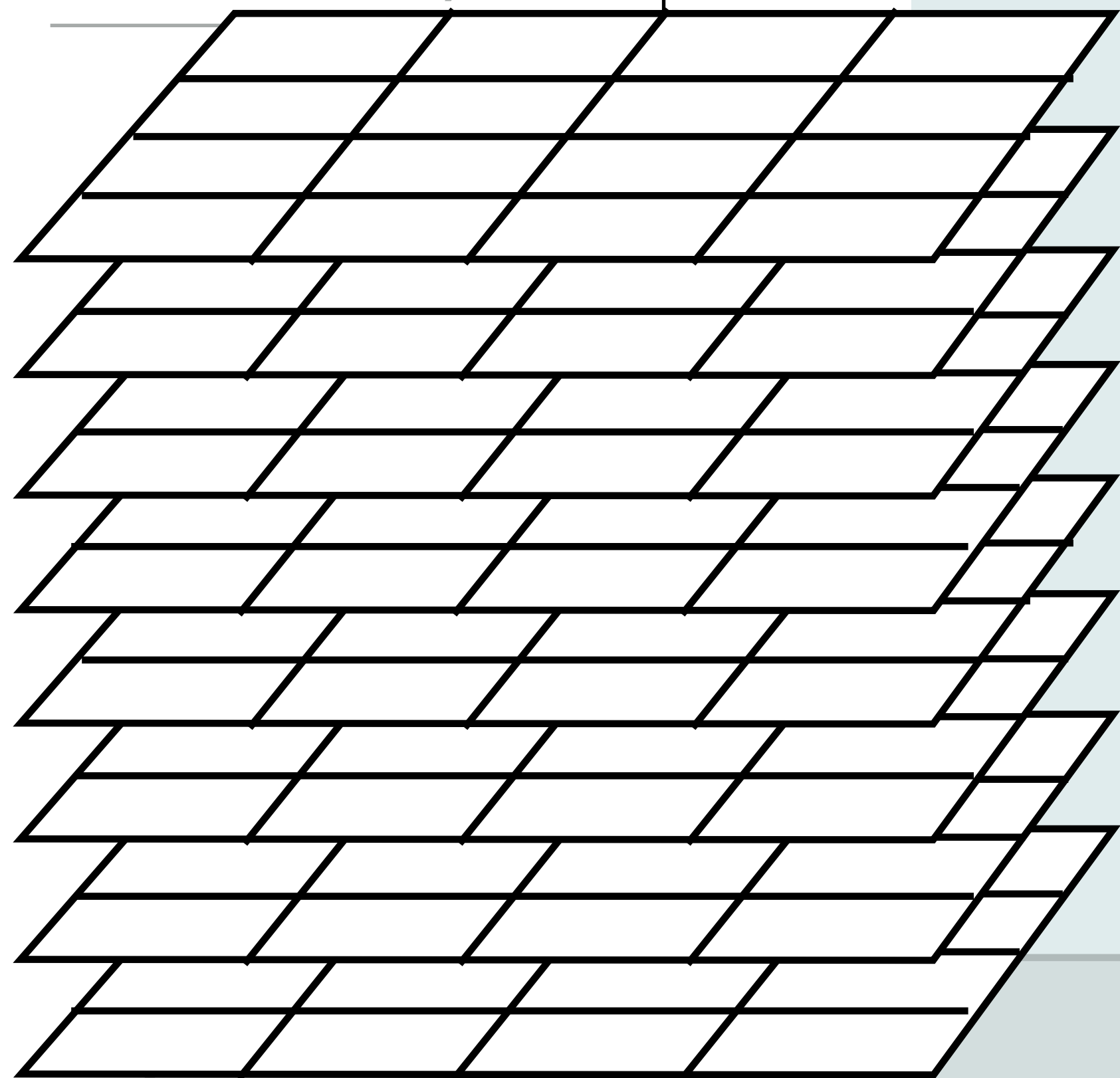
4 I/O Ports, up
to 80 GB/s each

Next gen is
160 GB/s per
(640 total)

Total conc'y =
16 x 8 x 2..8
(256–1024)

High Bandwidth Memory

Uses a simple '2.5D' instead of full 3D stacking



TSV Stack
Up to 4 or 8
DRAM dies

HBM DRAMs

1024-bit
8-Channel
Wide Interface

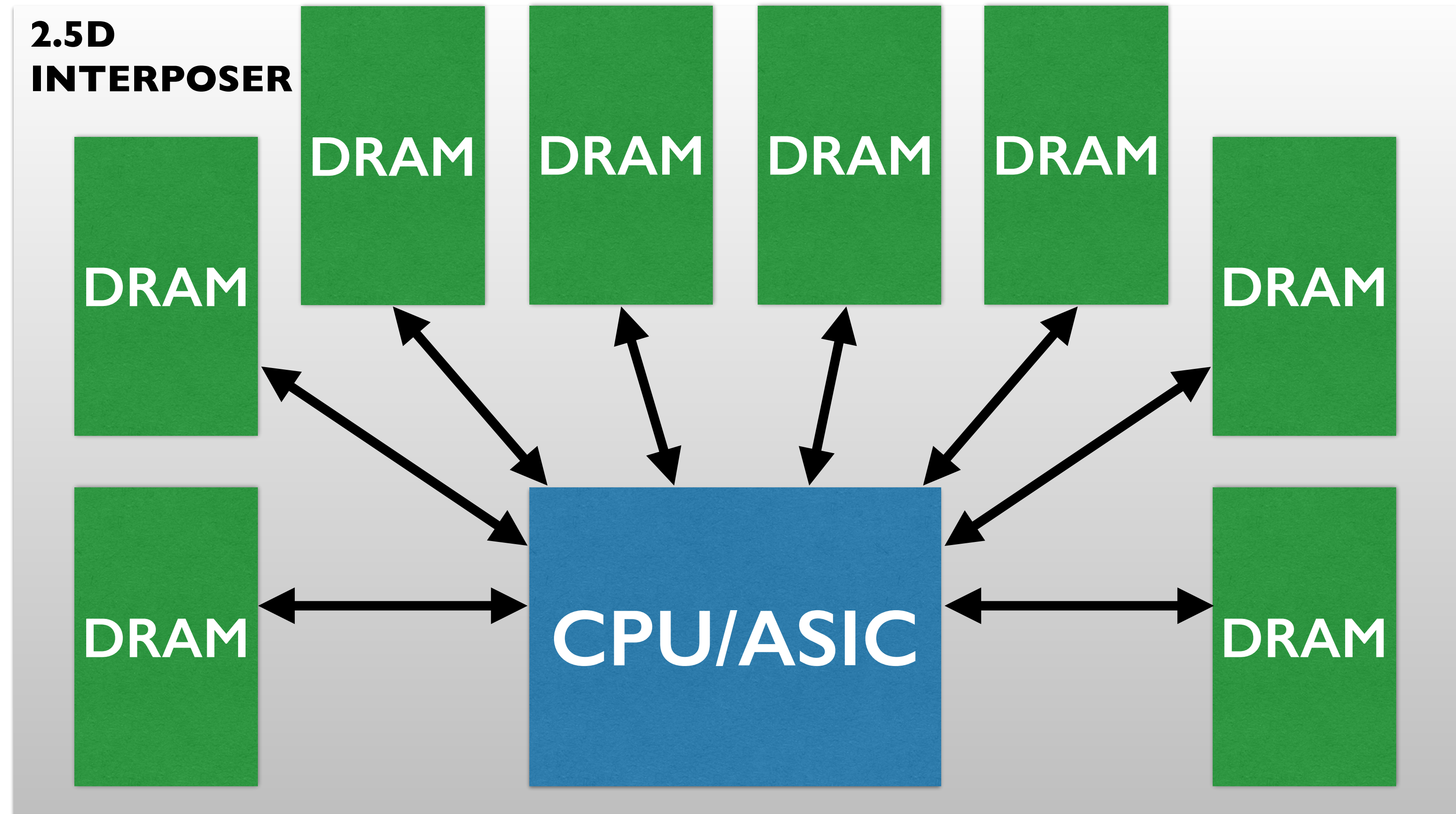
HBM
Interface

1024-bit x 2Gtps
= 256 GB/sec

GPU/CPU

TSV Interposer

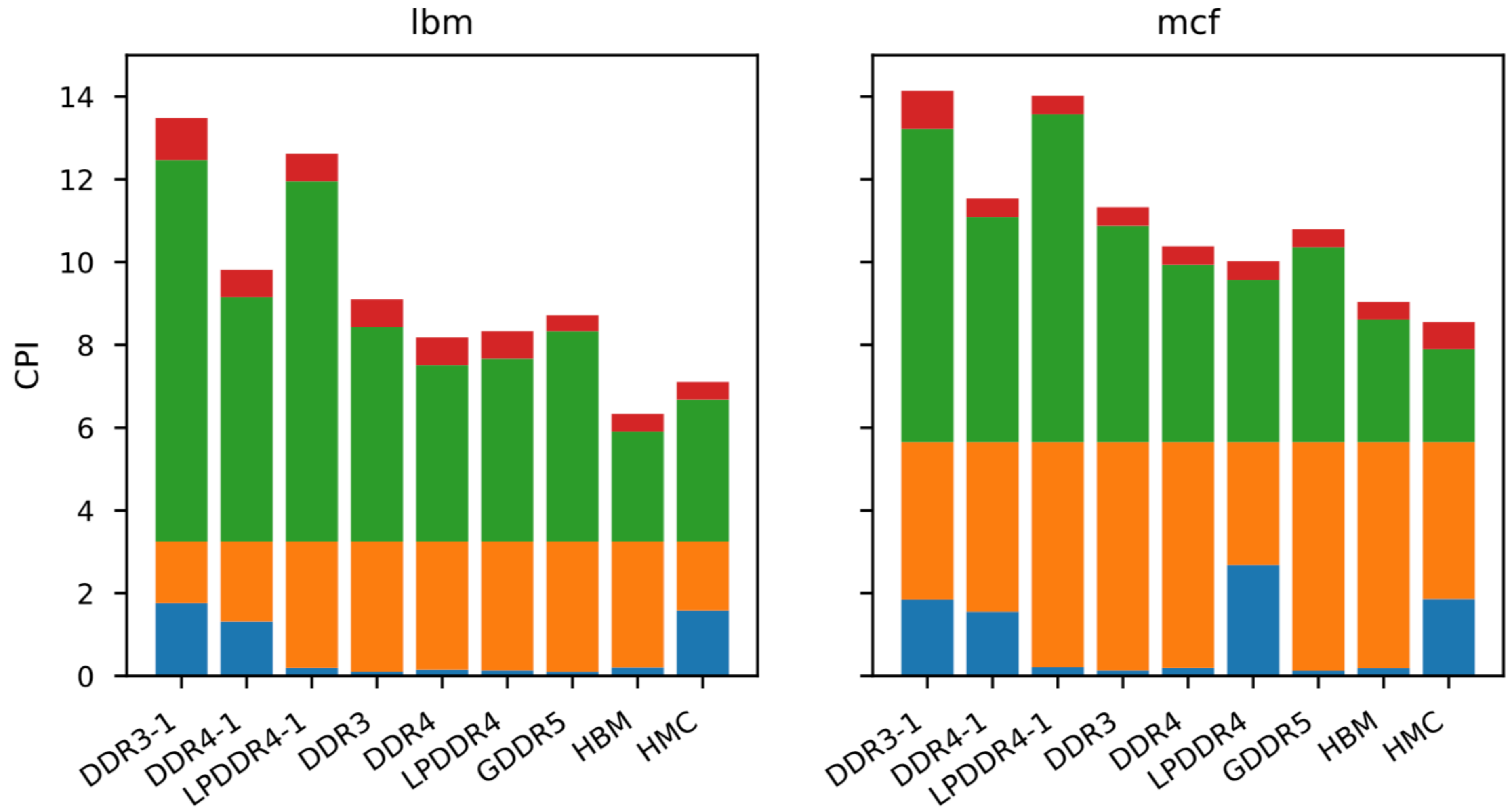
High Bandwidth Memory



Each Link is 128 Bits Wide: 1024 Total

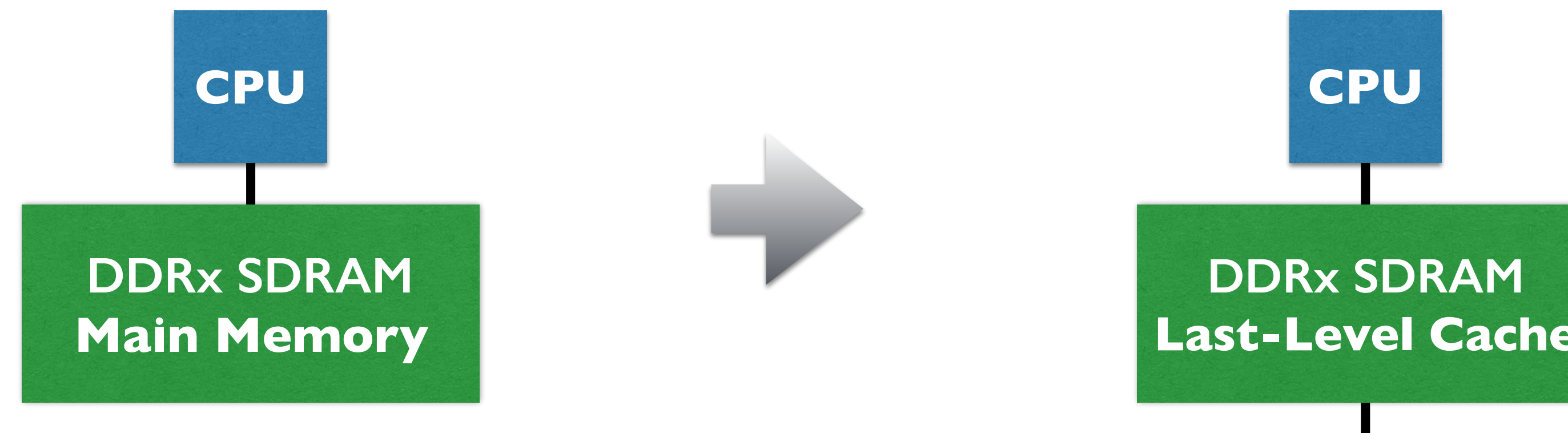
Performance Comparison

MEMSYS 2018



Non-Volatile Main Memory

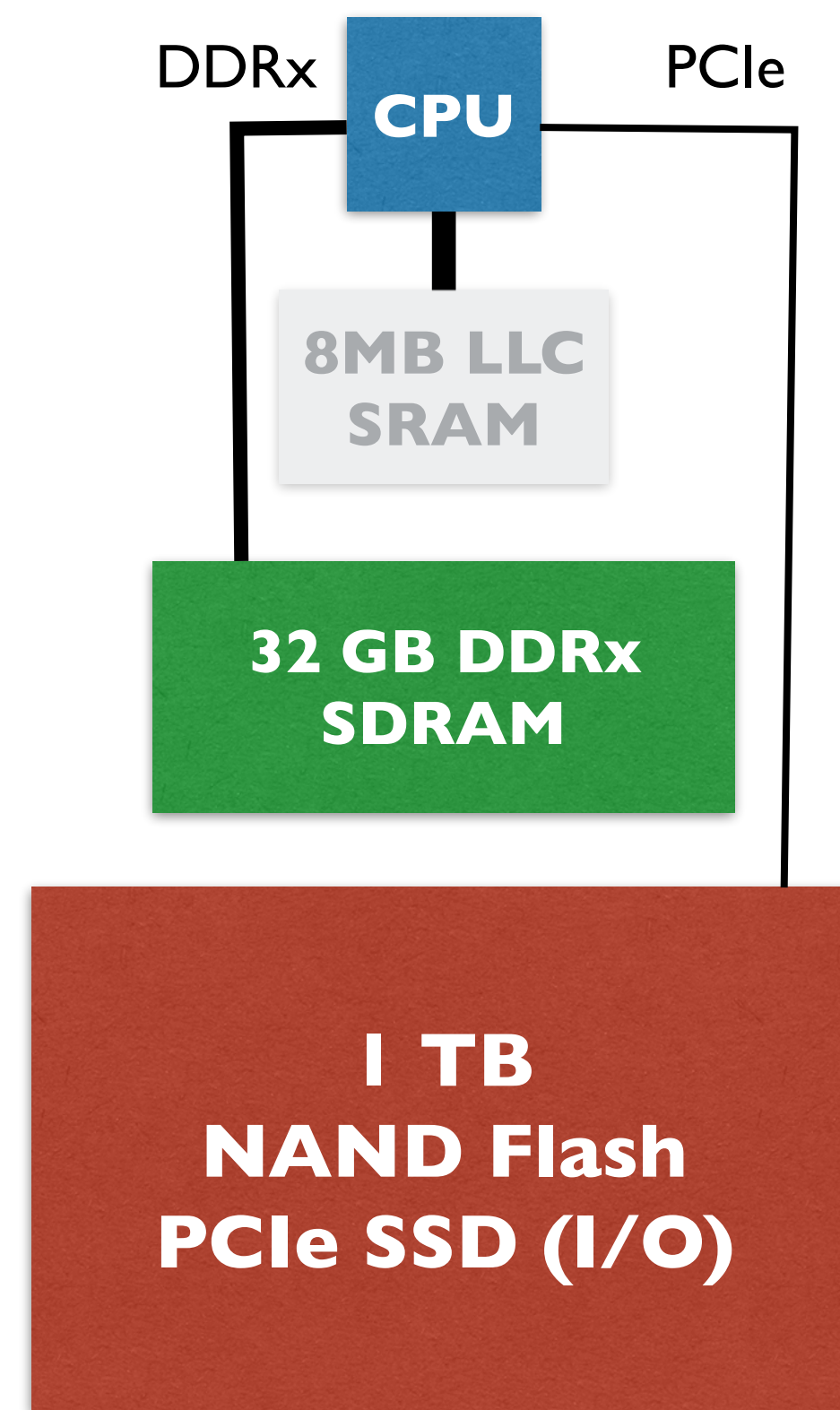
	Cost for 10 GB	Size of 10 GB	Power for 10 GB	Power per GB/s
Off-Chip SRAM	\$1,000	1 bucket	0.1–1 W	0.1 W
DDR4 SDRAM	\$100	1 DIMM	1 W	0.1 W
NAND Flash	\$10	<1 chip	0	0.1 W (?)
3D XPoint	\$40	<1 chip	0	0.1 W (?)



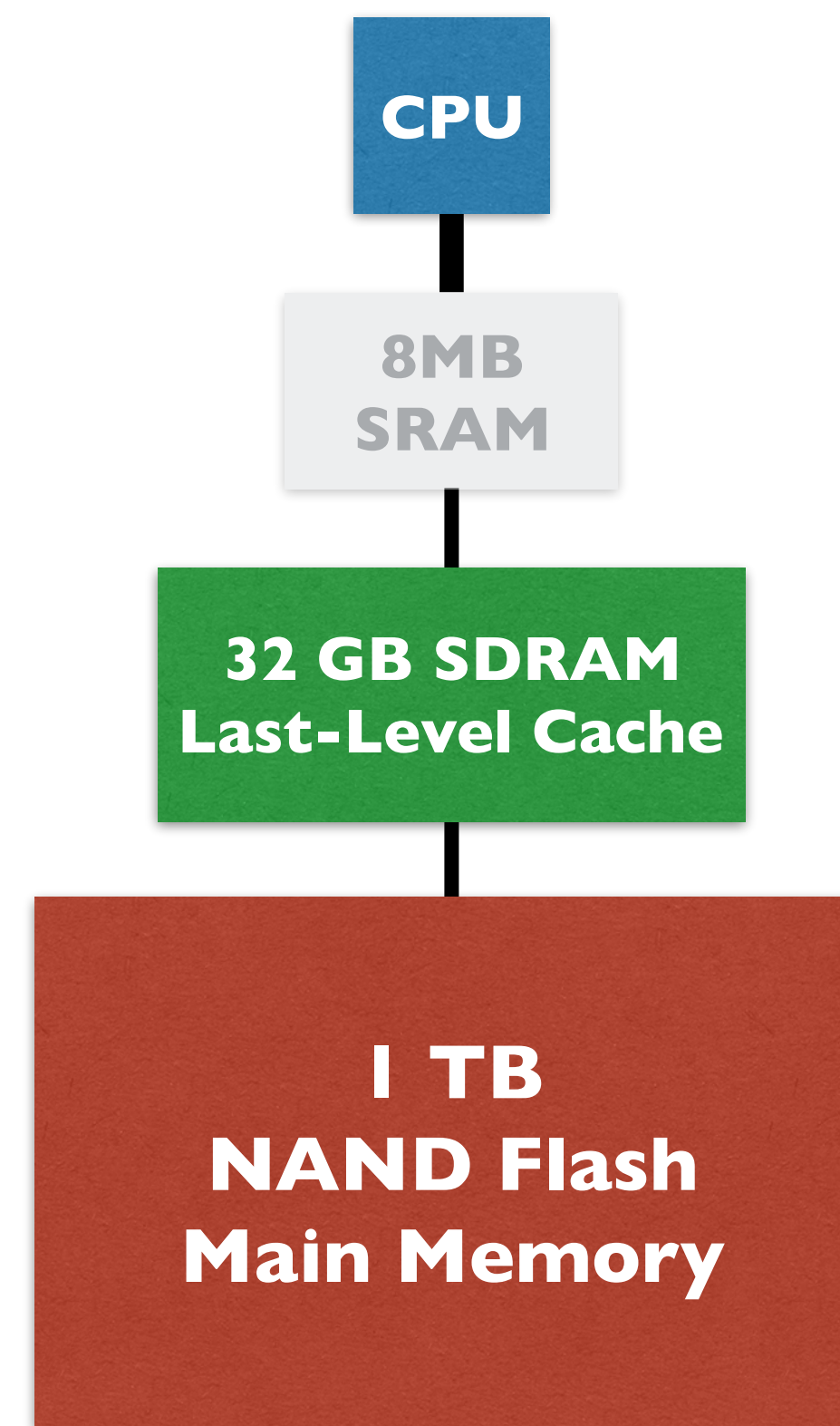
Note: wear-out mitigated by using MANY devices (thousands). A single device would wear out in under two days; therefore, 1000 devices should last for at least a year. Next, you can trade off longevity for access time and wearout: if the data need only last hours or minutes, wearout is reduced.

NAND Flash Main Memory
(... or **any** source of cheap bits)

A Tale of 3 Memory Systems



SSD
\$500 – 10W



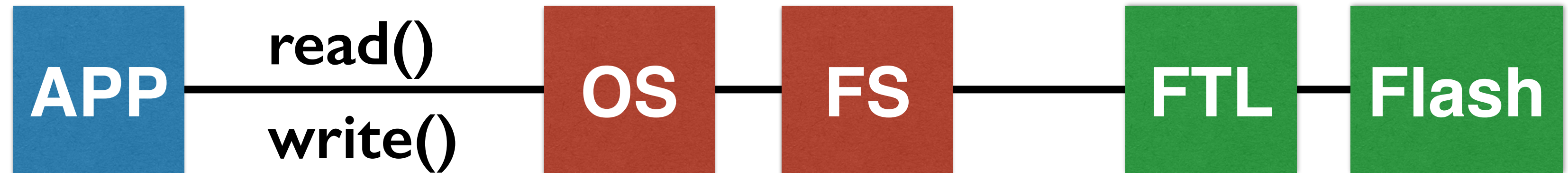
NVMM
\$500 – 10s of W



Ideal
\$10,000 – 100W

A Tale of 3 Memory Systems

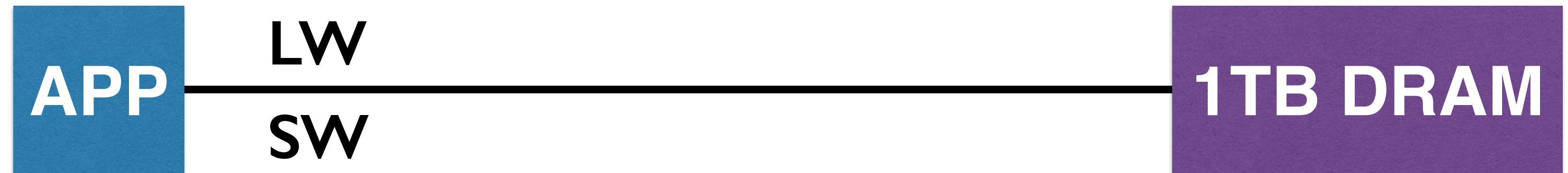
SSD



NVMM

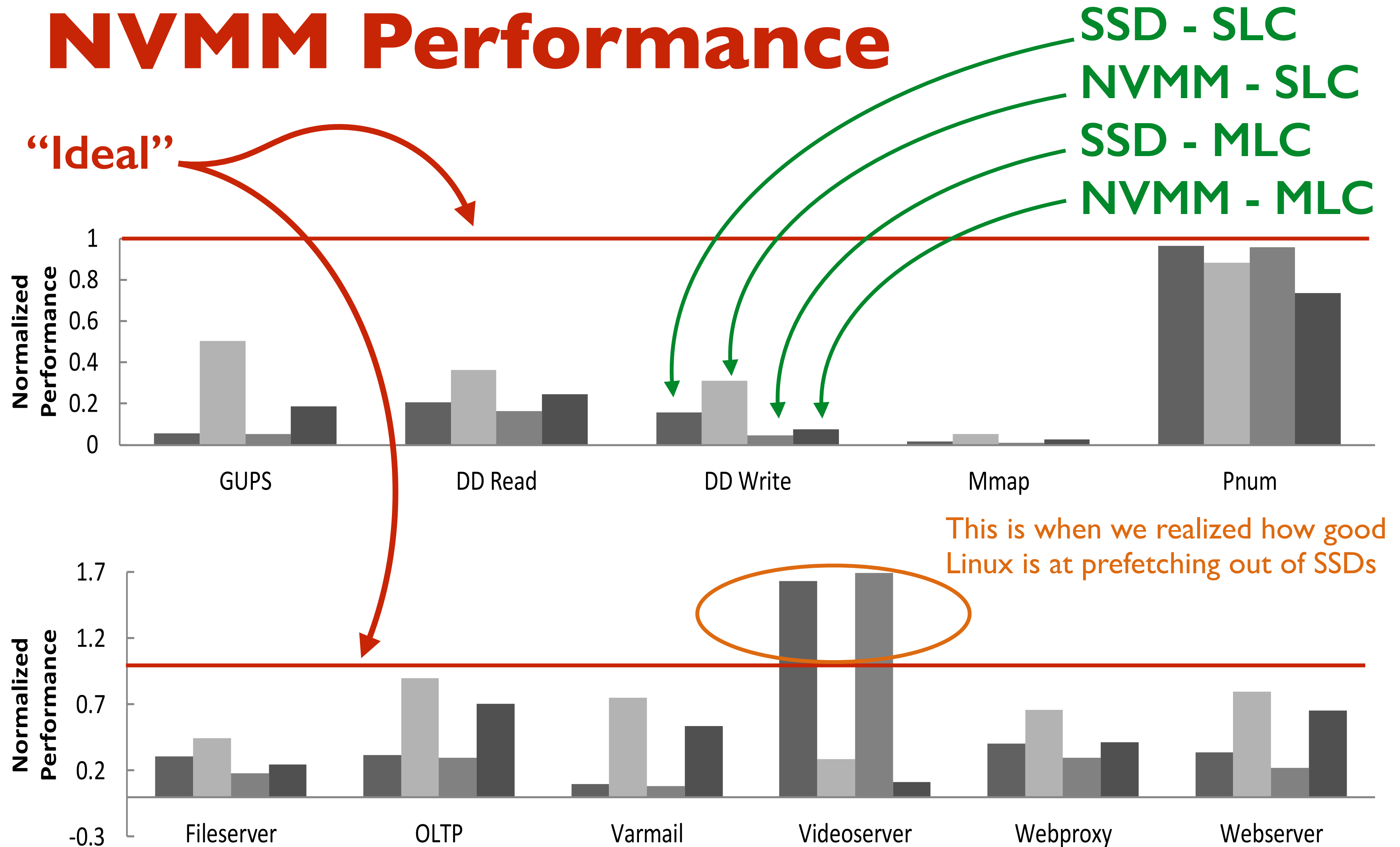


Ideal

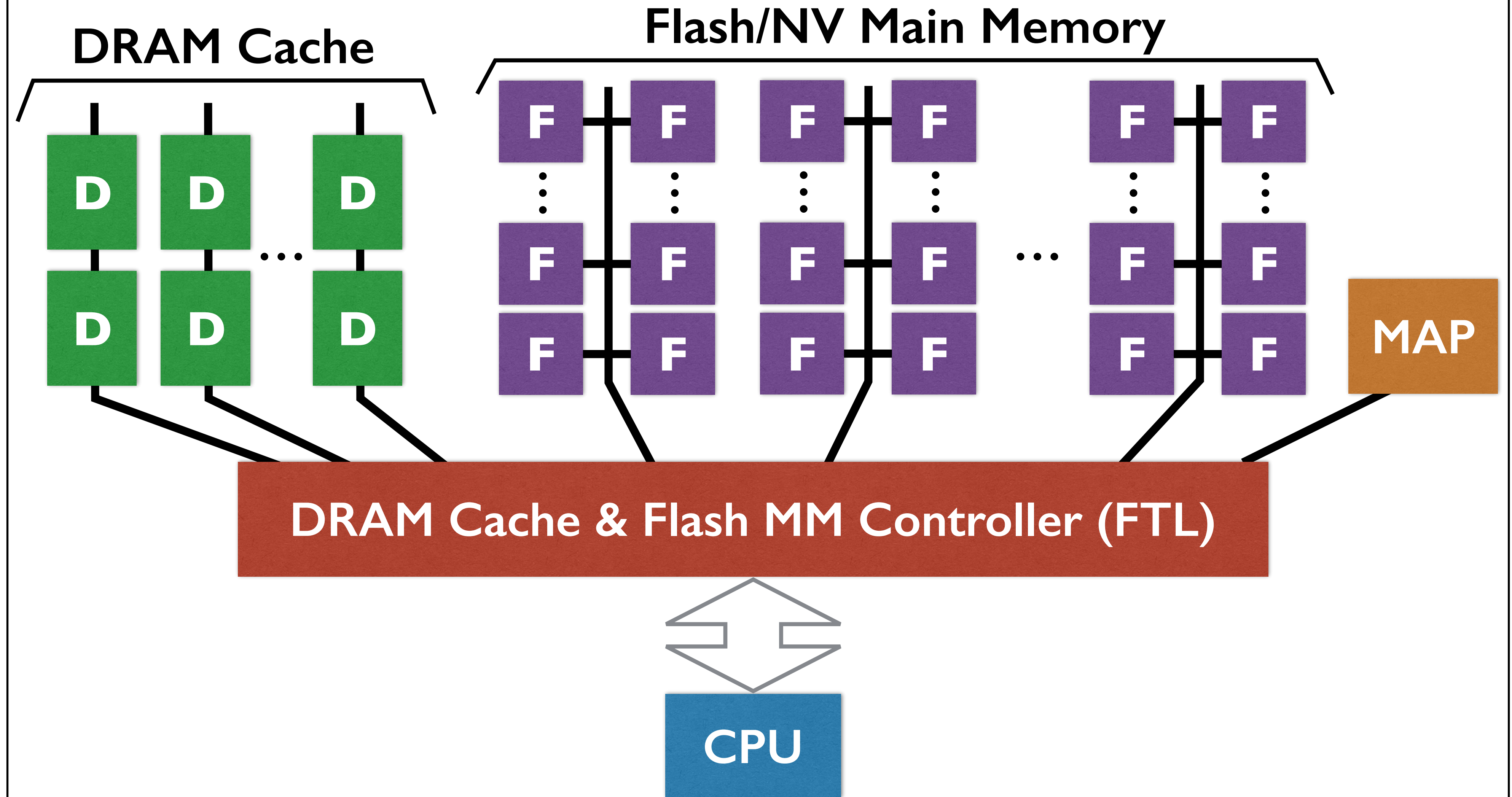


NVMM Performance

“Ideal”



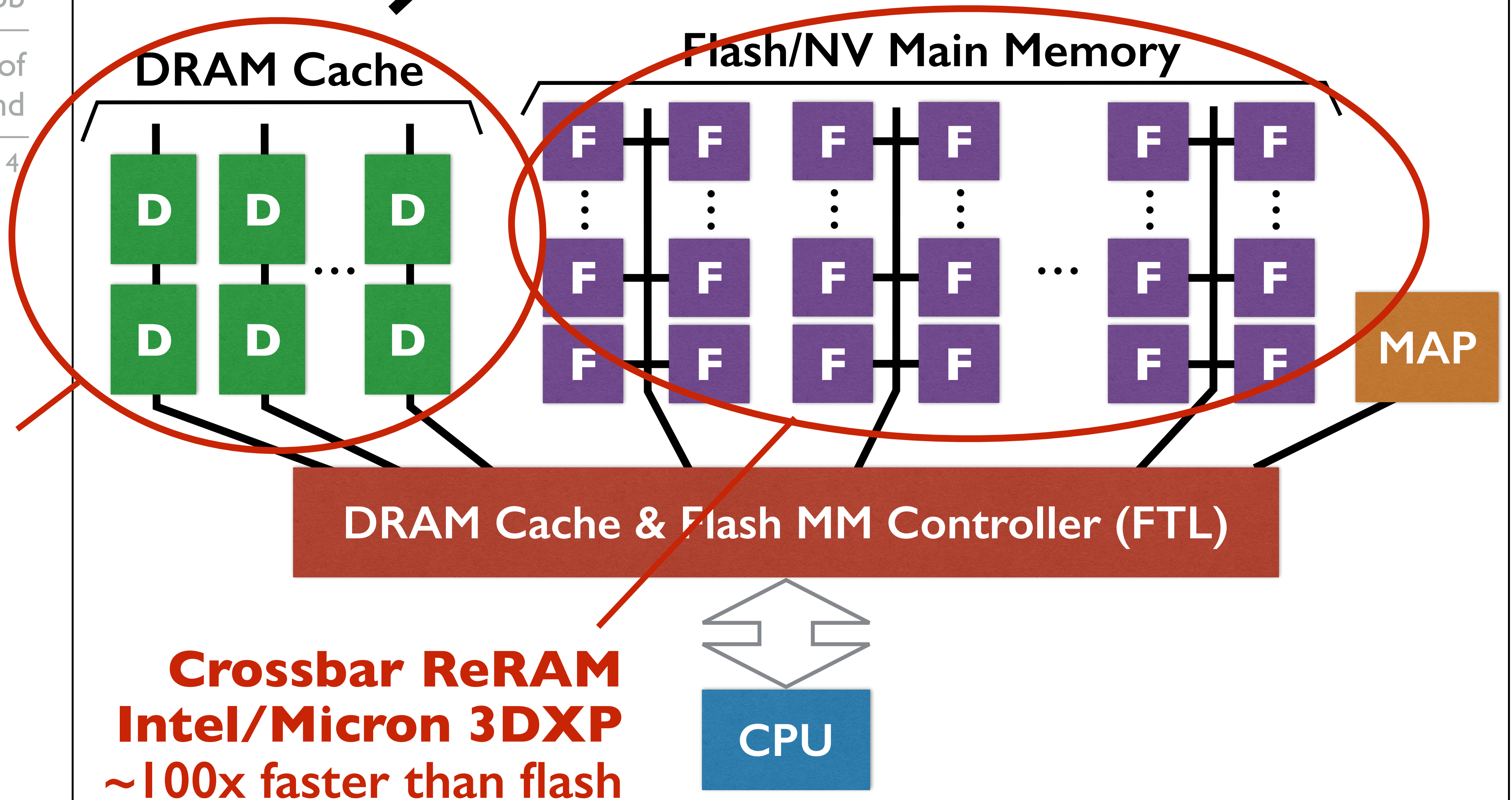
Yeah, it's a lot of engineering



it was
Yeah, it's a lot of engineering

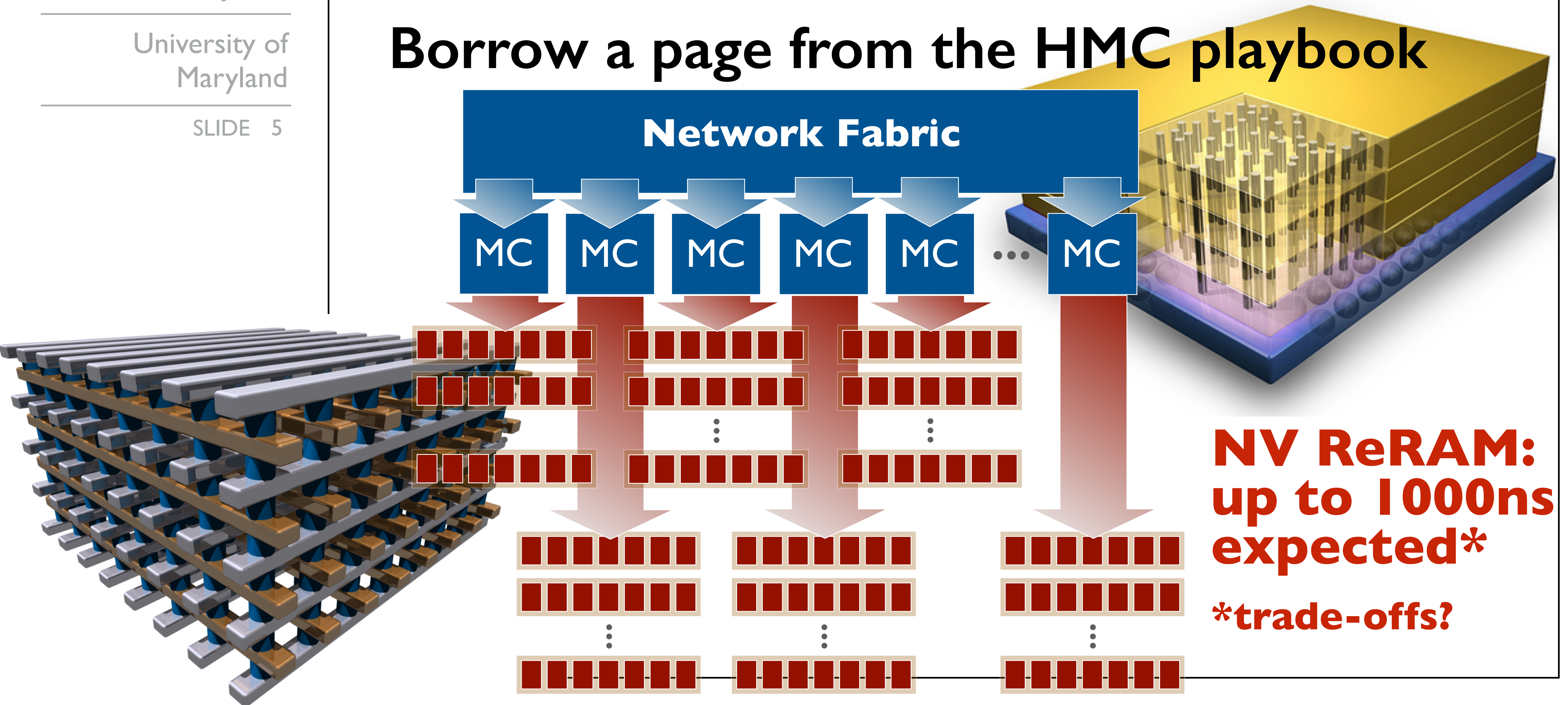
HMC:
320GB/s
16 channels
HBM:
256GB/s
8 channels

Crossbar ReRAM
Intel/Micron 3DXP
~100x faster than flash

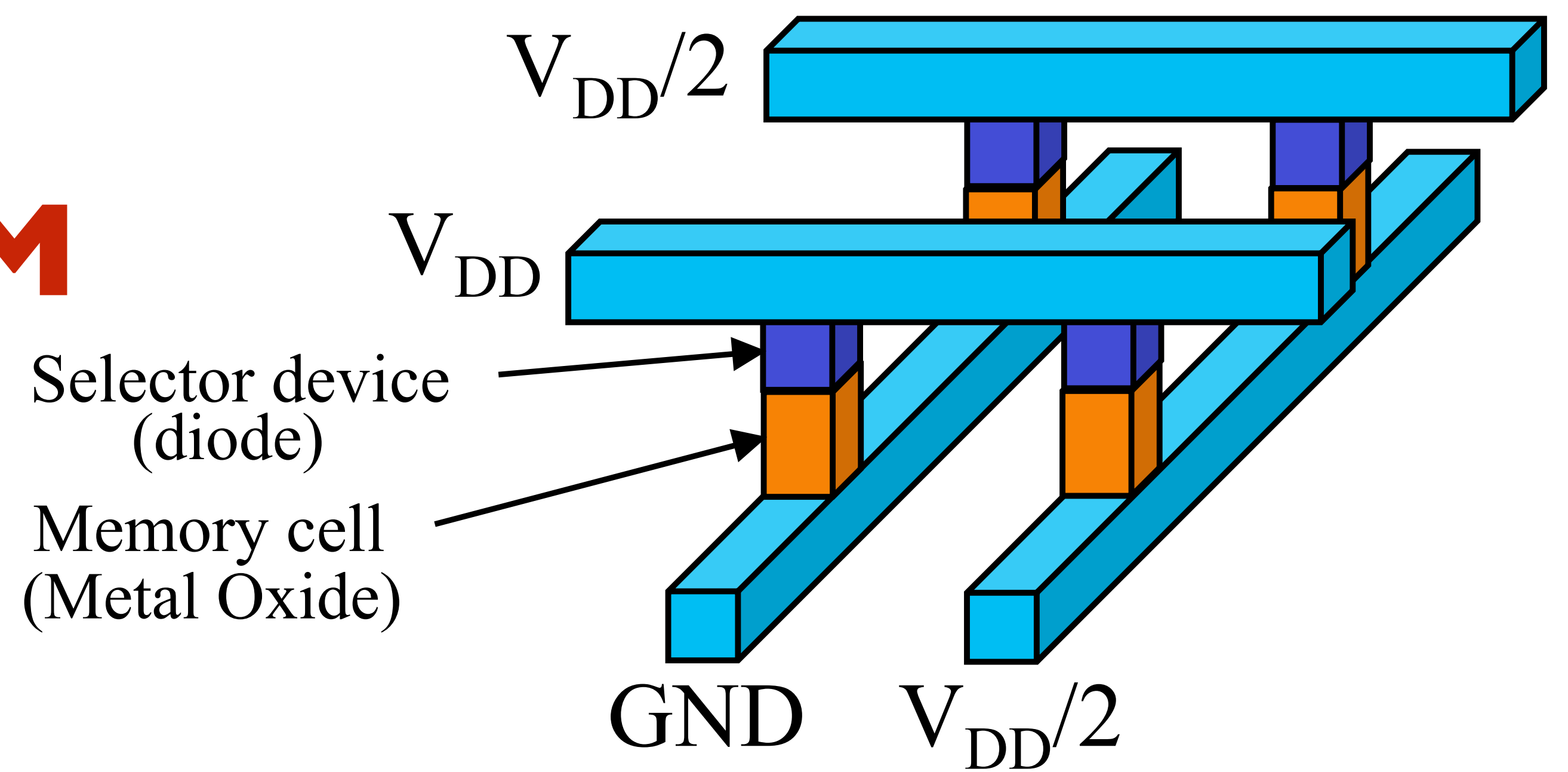
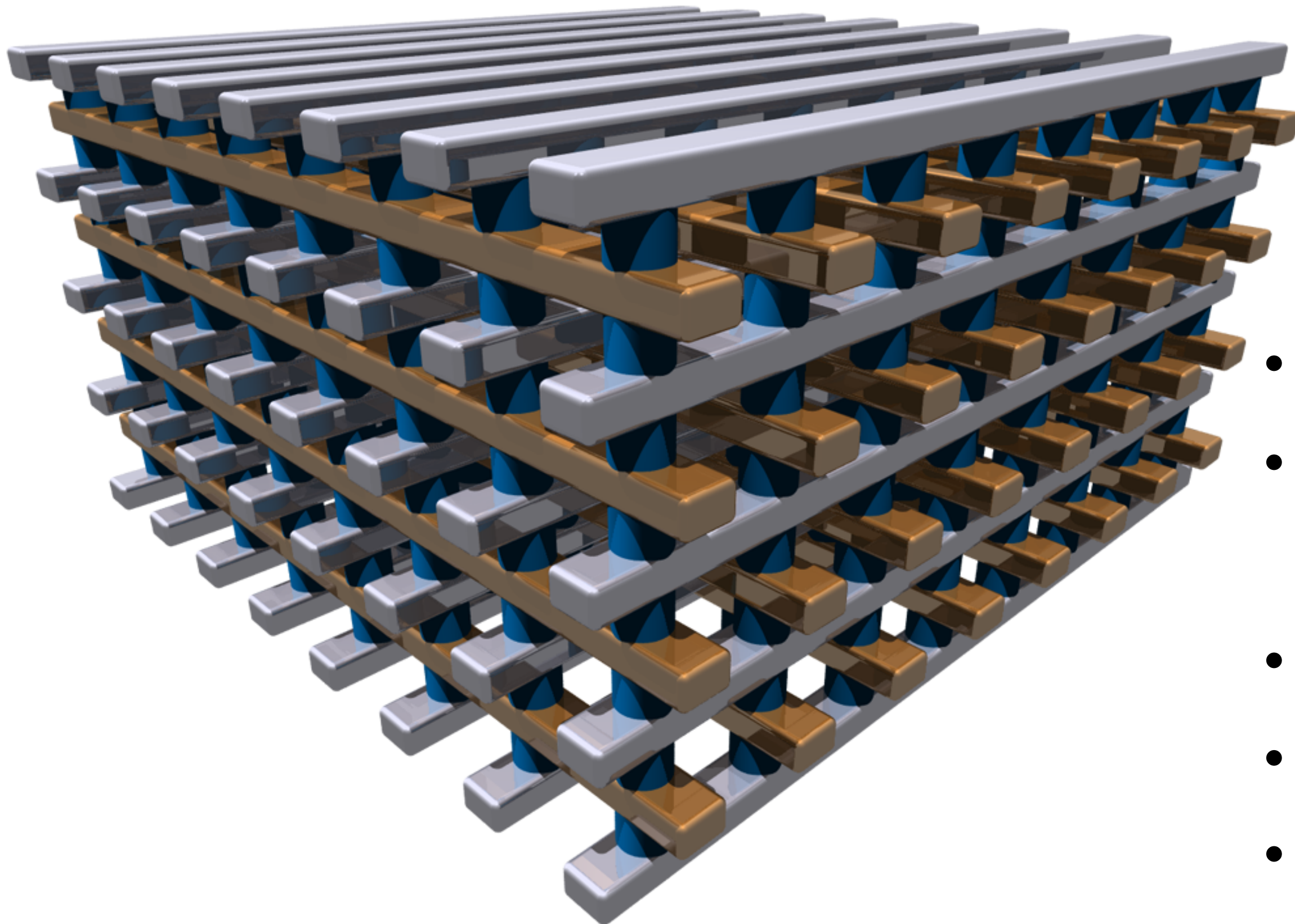


High Bandwidth Non Volatiles

Borrow a page from the HMC playbook



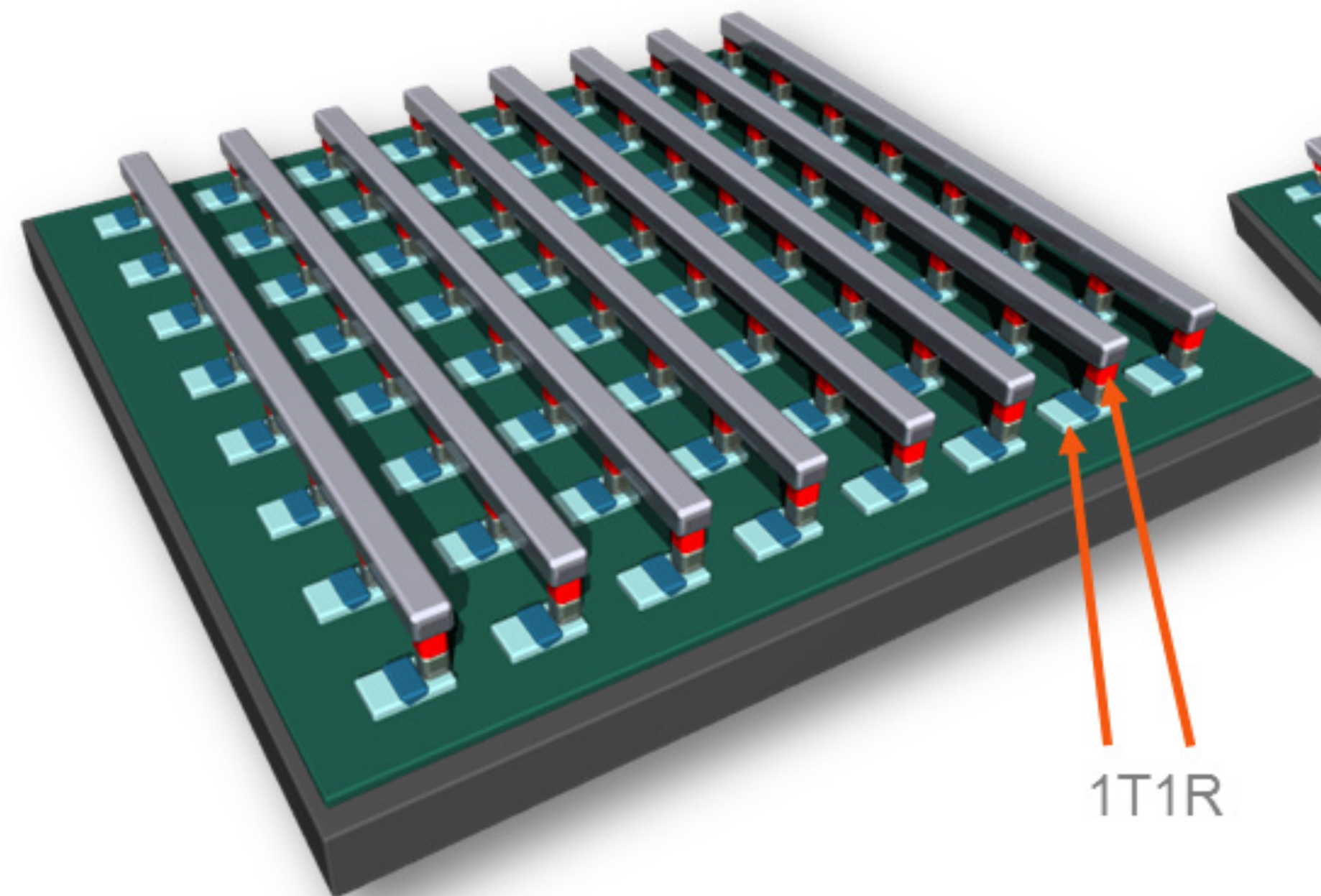
Crossbar 3D ReRAM



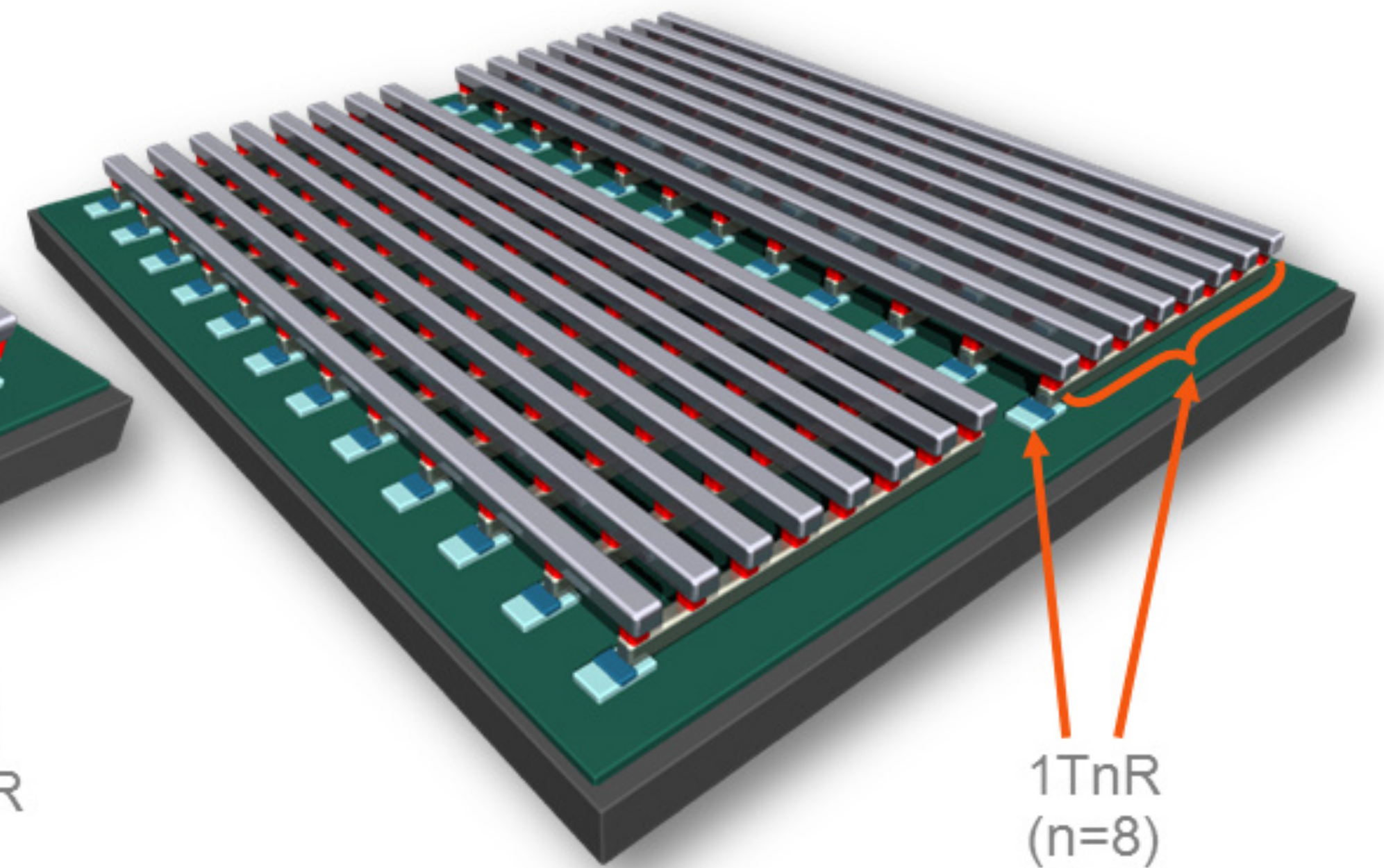
- Cells minimum area (no access transistor)
- 2-stack arrays @ 16nm, 20 x 20 mm die:
64GB of ReRAM
- 8-stack arrays => **256 GB of ReRAM**
- Stacks arbitrarily high
- No. Access. Transistor.

No Access Transistor

1T1R Memory Array
Low Latency, Low Density



1TnR Memory Array
High Performance, High Density



(n = 1 .. 2048)

Crossbar RRAM Technology

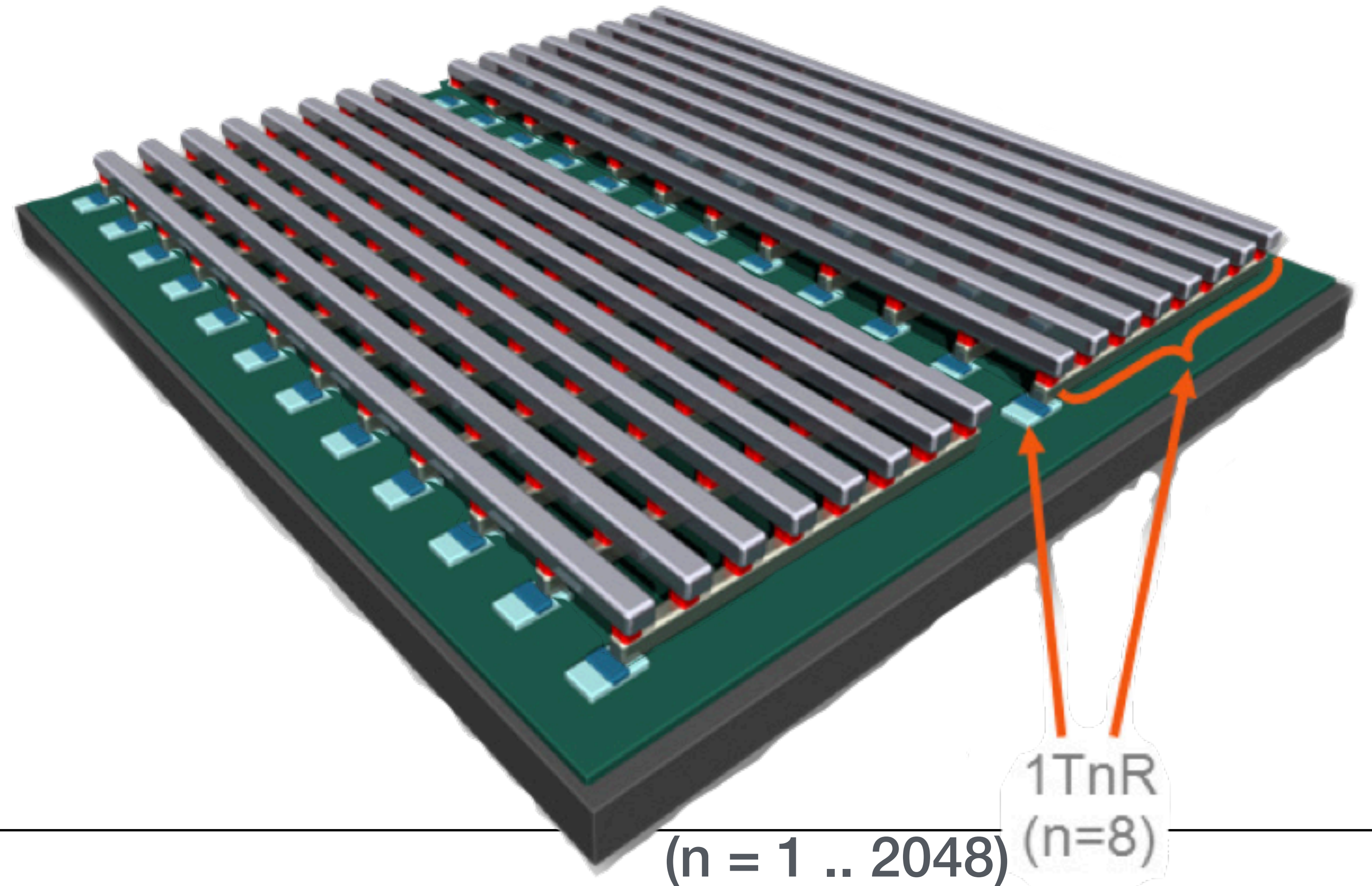
All Tomorrow's
Memories

Bruce Jacob

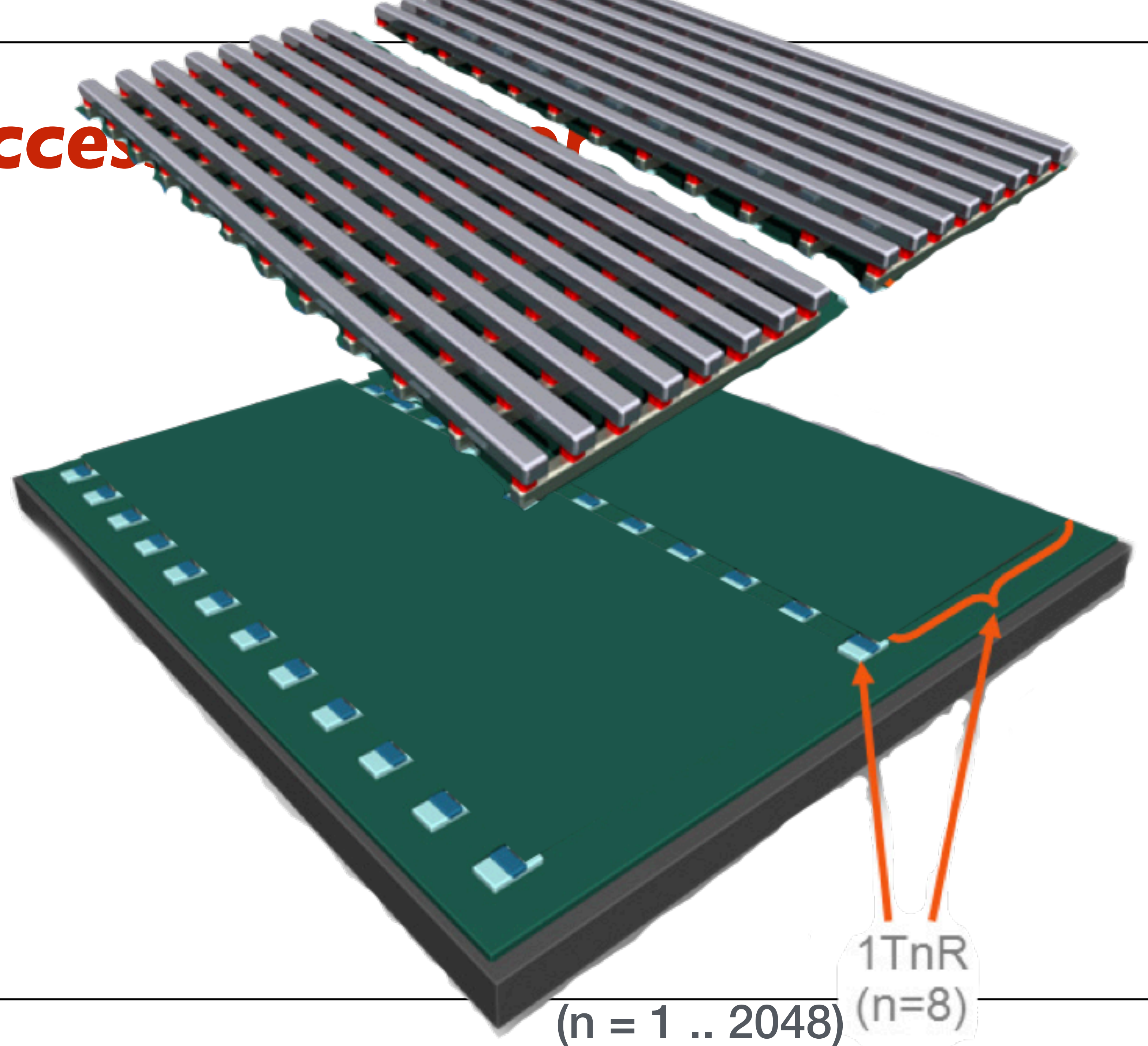
University of
Maryland

SLIDE 18

No Access Transistor



No Access



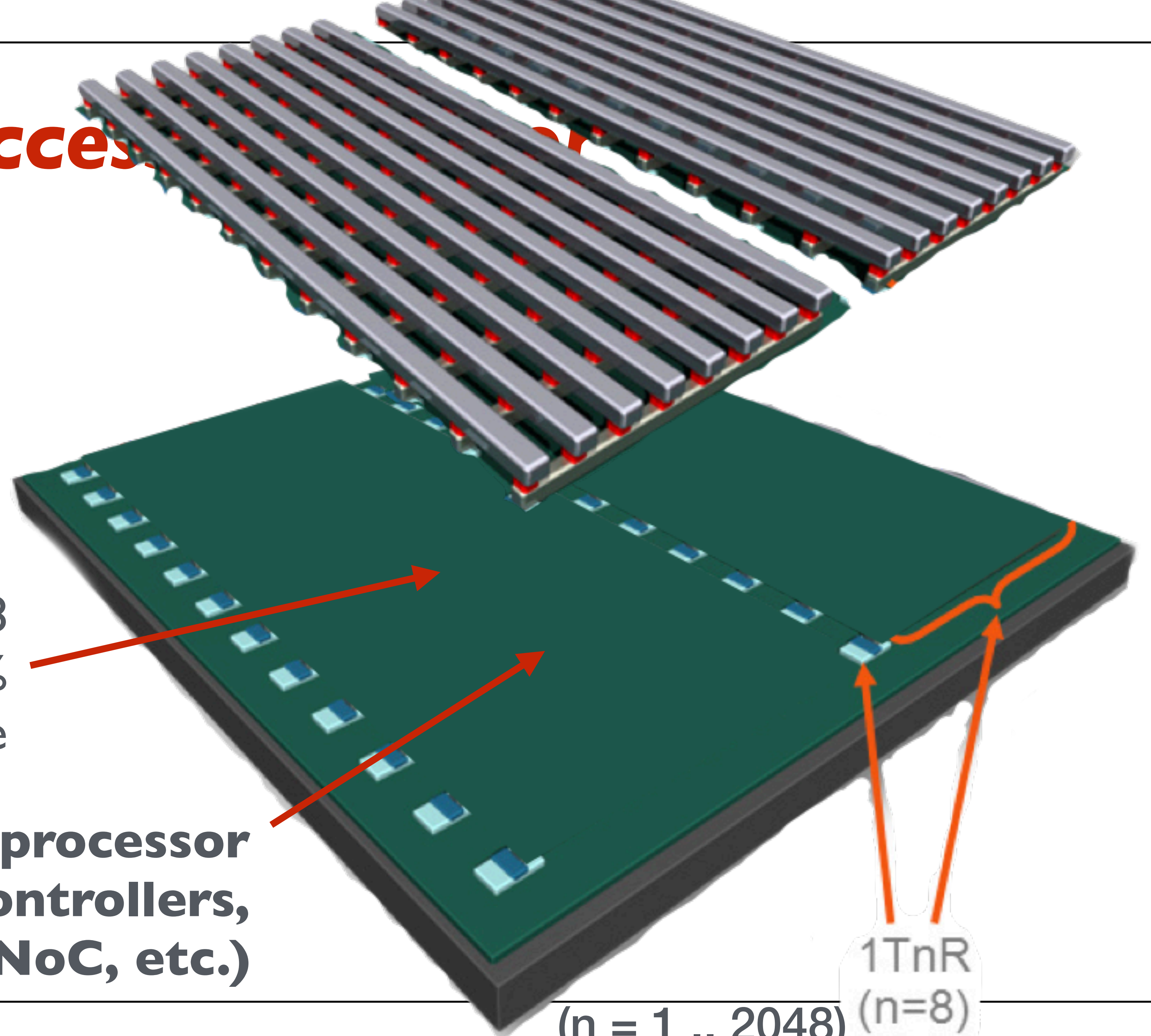
No Access

For $n = 2048$
area is $\sim 75\%$
white space

**Use for processor
(cores, controllers,
routers, NoC, etc.)**

($n = 1 \dots 2048$)

1TnR
($n=8$)



All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

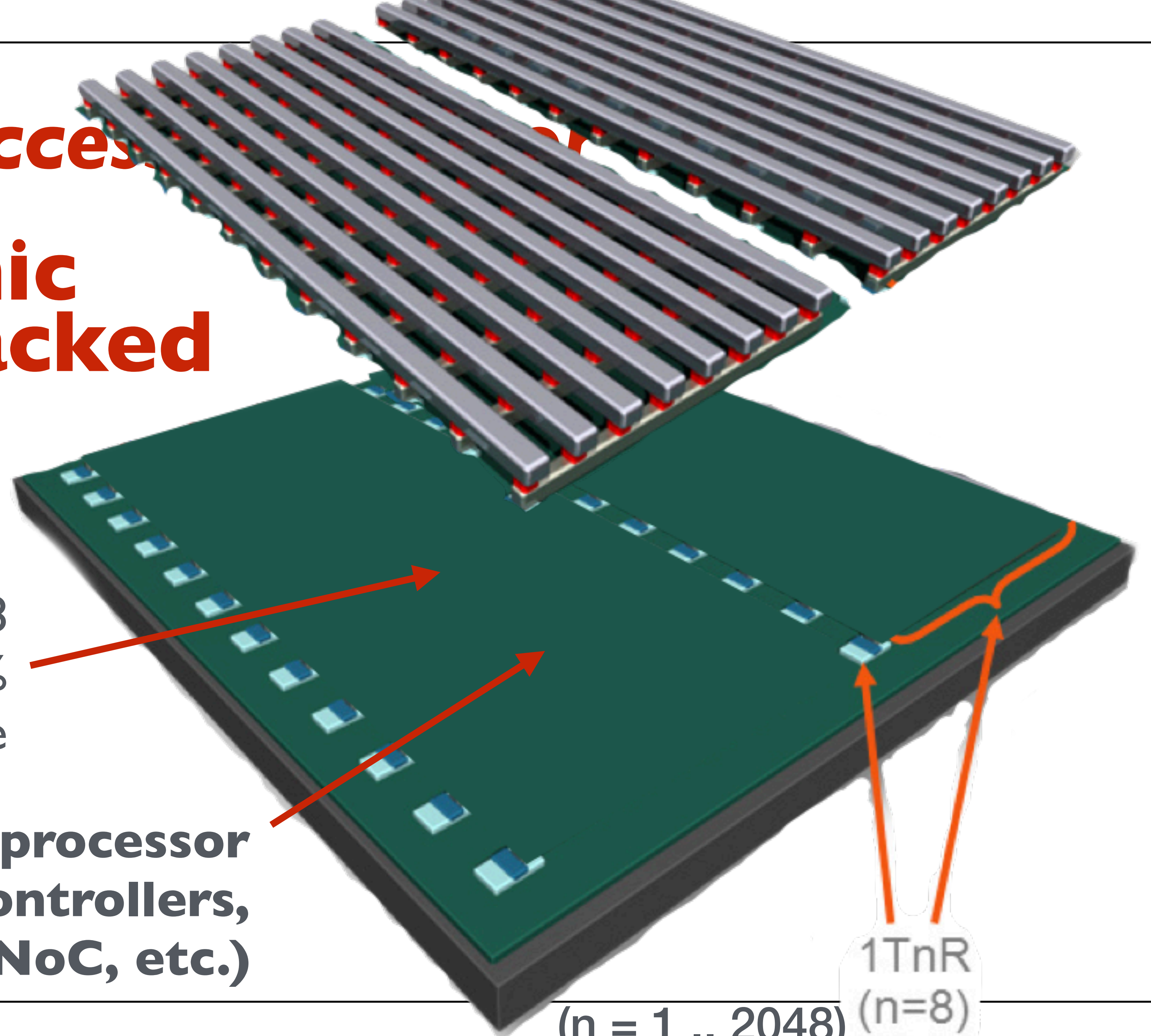
11/18

No Access **Monolithic** **Not Die-Stacked**

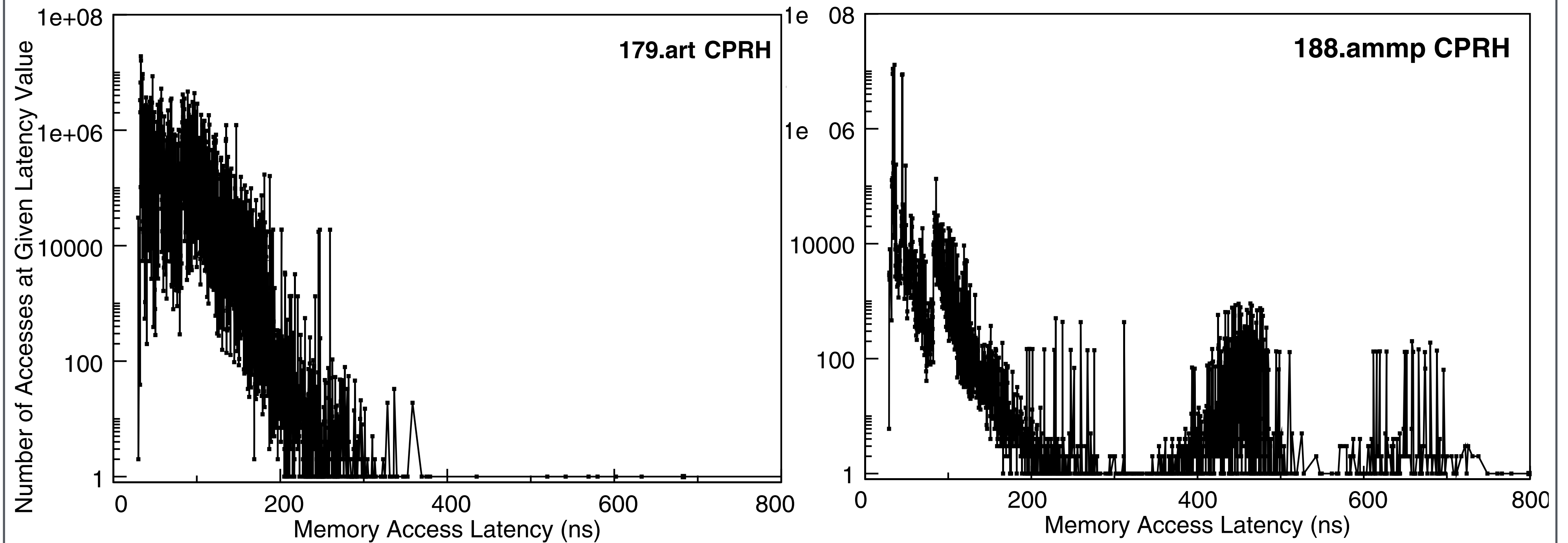
For $n = 2048$
area is $\sim 75\%$
white space

Use for processor
(cores, controllers,
routers, NoC, etc.)

($n = 1 \dots 2048$)
1TnR
($n=8$)



Recall: Real DRAM Latency Is Actually Quite Long (100ns)



This is for single core. Multicore can be much, much worse.

Example Monolithic Numbers

~64 cores, ~256GB ReRAM, ~4k banks

Assume 200ns latency for 8-byte payload:

$$\begin{aligned}\text{Bandwidth} &= 4\text{k} * 8 \text{ bytes} / 200\text{ns} \\ &= 4\text{k} * 40 \text{ MB/s} \\ &= \underline{160 \text{ GB/s}}\end{aligned}$$

e.g., 64 cores, each 4-way multithreaded,
each with 512-bit (8-longword) SIMD,
vectored & scatter-gather loads,
4-deep non-blocking => 8k

So what all does this enable?

**HBM/HMC: hugely parallel systems
(the duality of bandwidth and parallelism),
streaming applications, 2x performance**

**NVMM: massive data sets, new OS
paradigms such as merged VM+FS and
journalled main memory
(built-in checkpoint/restart)**

**HBNV: fine-grained operations on
enormous sparse data sets**

Datacenter & Cloud Issues

**Distribution at storage-level interface
simplifies application development**

Potential for significant performance

**RoCE appropriate for supercomputing?
How about RoXX?**

**At what round-trip latency does this
rival MPI as programming model?**

Expect a shake-up soon.

Nonvolatility Issues

Unified VM+FS Subsystems (OS redesign)

- ➔ *By default, data in process address space temporary, garbage-collected at exit(); **permanentify** function to keep around*
- ➔ *Possible directions:*
 - **Persistent objects (e.g. Mneme, POMS)**
[failed only due to reliance on disk]
 - **Named regions**
- ➔ *Journalled main memory w/ checkpointing*

Capacity Issues

Rethink Protection & Translation

➔ *TLB overhead is ~20%*

- **So get rid of it already!**
- **BUT: need protection, authentication**

➔ *Why not waste bits? Simplify both sharing and translation by eliminating much of VM*

➔ *OS/HW co-design needed: e.g., sharing via $vaddr$ instead of $paddr$, language support?*

Recall: Nonvolatile main memories ~TB per node

Bottom Line

It's going to happen. :)

- **Combined VM+FS subsystems**
- **Journalled main memory**
- **Persistent Object Store work from 80s**
- **OS: Simpler design, fewer potential bugs**
- **VM arguably a way better abstraction to distribute than the FS**
- **Monolithic = good for many applications**

Bottom Line

It's going to happen. :)

- Combined VM+FS
- Journal
- ...

... and the storage guys
are showing us the way!

from 80s

wer potential bugs

a way better abstraction

route than the FS

Monolithic = good for many applications

All Tomorrow's
Memory Systems

Bruce Jacob

University of
Maryland

SLIDE 26

Shameless Plug

www.memsys.io

Washington DC Oct 2019

Call For Papers

www.memsys.io

Call For Papers

MEMSYS 2018

The International Symposium on Memory Systems ❖ October 1–4, Washington DC

Important Dates

Memory-device manufacturing, memory-architecture design, and the use of memory technologies by application software all profoundly impact today's computing systems, in terms of their performance, predictability, power dissipation, and cost. Existing memory technologies are seen as limiting in terms of power, capacity, and performance. Emerging memory technologies offer the potential to overcome these and design-related limitations to answer the requirements of new applications. Our goal is to bring together researchers, students, and others interested in this exciting and rapidly evolving field to get together on the latest state of the art, to exchange ideas, and to address current challenges. Visit memsys.io for more information.

Accepted Papers

Accepted papers containing significant novel ideas and contributions are solicited. Papers focusing on system, software, and hardware concepts, outside of traditional conference scopes, will be preferred over others (e.g., the desired focus is away from pipeline design, processor cache design, prefetching, data prediction, etc.). Symposium topics include, but are not limited to, the following:

- Memory-system design from both hardware and software perspectives
- Memory failure modes and mitigation strategies
- Memory-system resilience, especially at large scale
- Memory and system security issues
- Operating system design for hybrid/nonvolatile memories

(double-blind), blind submission (no authors listed), up to 16 pages long

Organizers

Bruce Jacob, U. Maryland
Kathy Smiley, Memory Systems
Rajat Agarwal, Intel
Abdel-Hameed Badawy, NMSU

RAM, 3DXP, memristors, etc.
Languages, optimization
memory technologies
hardware, software, mitigation
cache/memory/accelerators
management techniques
hardware and software,
applications
datacenter applications
em-memory machines

technologies to support them,
and heterogeneous memories
beyond traditional
processors.

ideas that
groups—to
people,
people and
accepted
with papers, and each



accepted submission is given a 20-minute presentation time slot.
All accepted papers will be published in the ACM Digital Library.



The IEEE International Symposium on Memory Systems 2018

Jishen Zhao, UC San Diego

All Tomorrow's
Memories

Bruce Jacob

University of
Maryland

SLIDE 27

Thank You!

Bruce Jacob

blj@umd.edu

www.ece.umd.edu/~blj

