# arm Research

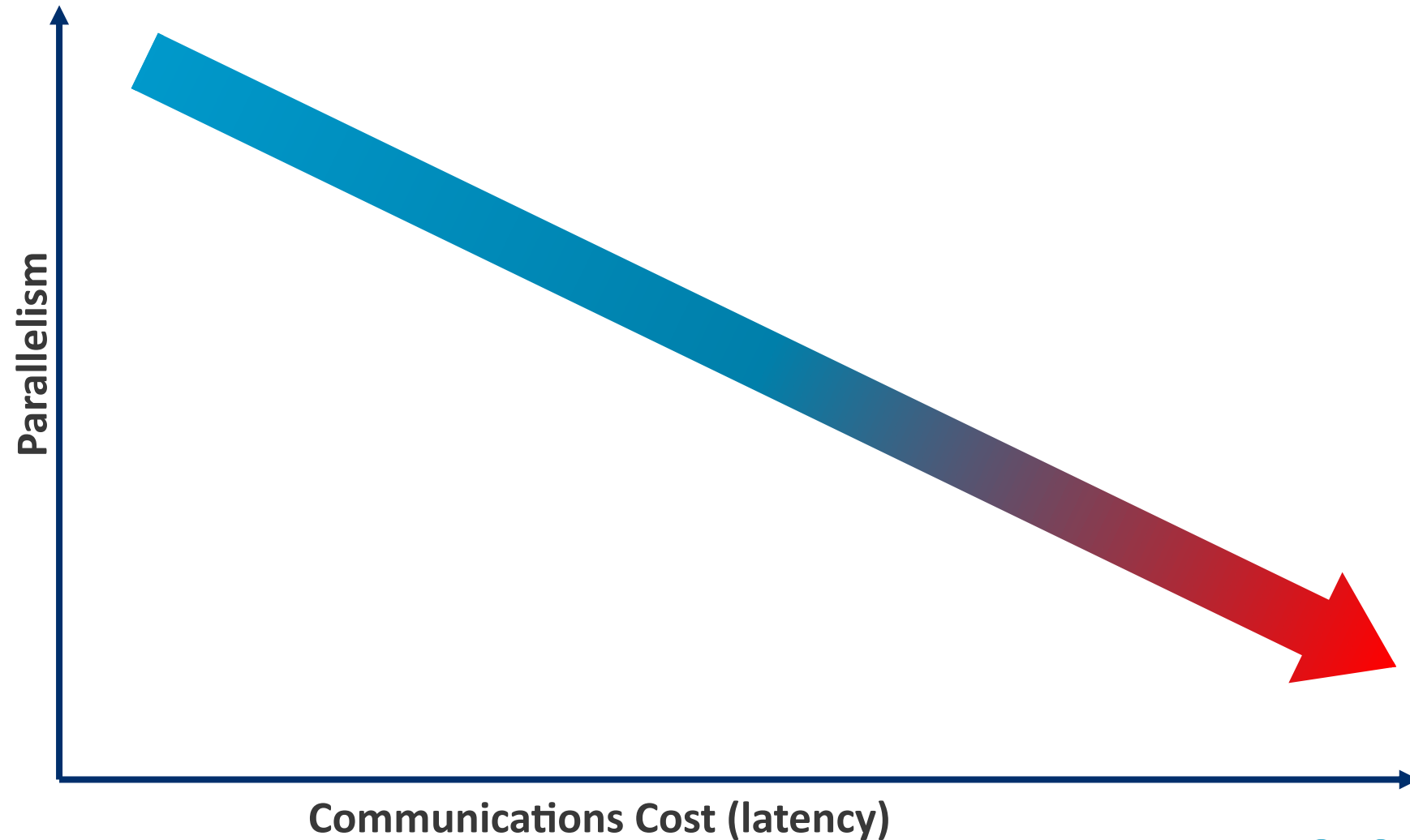# New and cool memory technologies

SYSTEMS

Jonathan Beard

2 October 2019
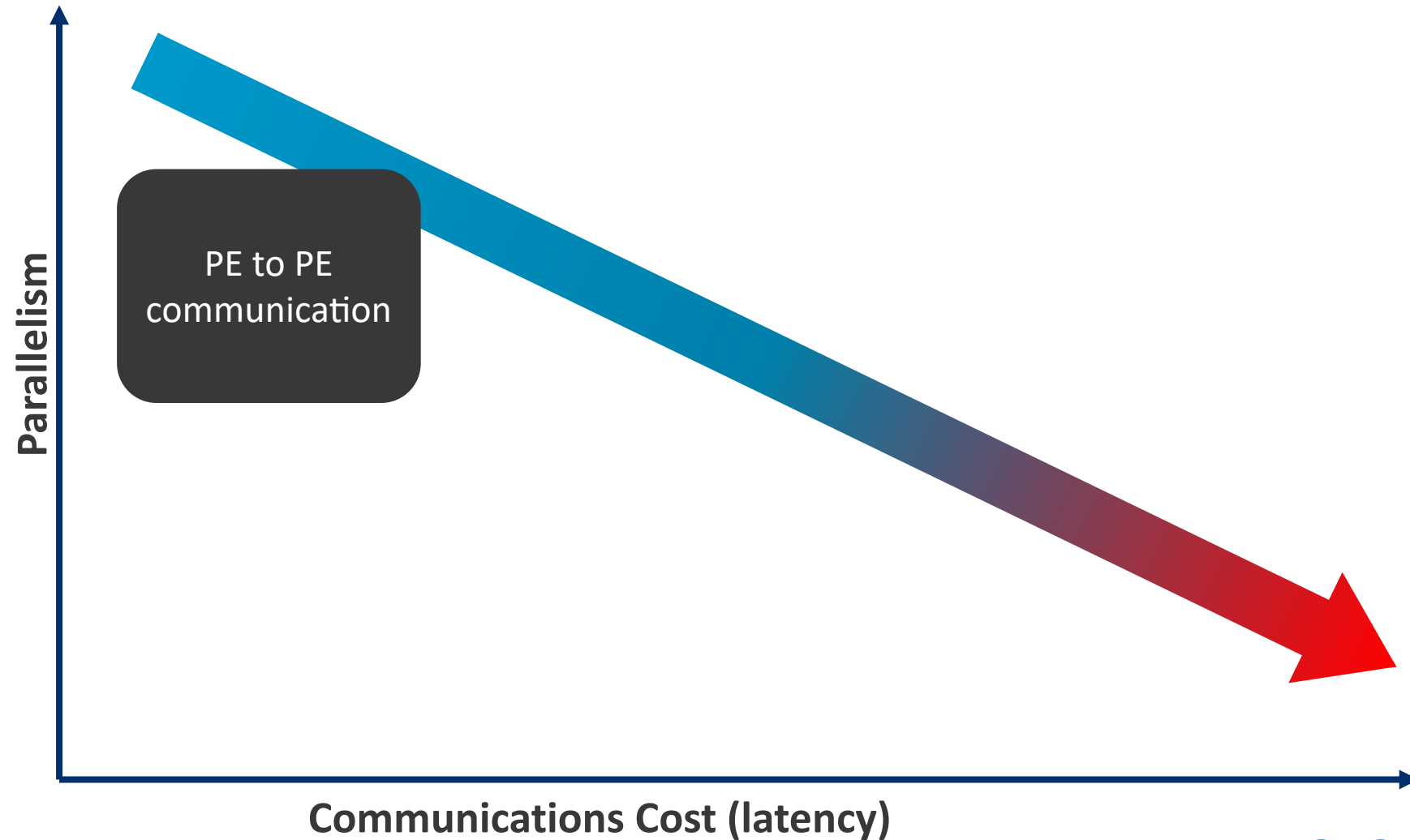
# Compute efficiency post-Moore

Data movement dominates, parallelism is critical
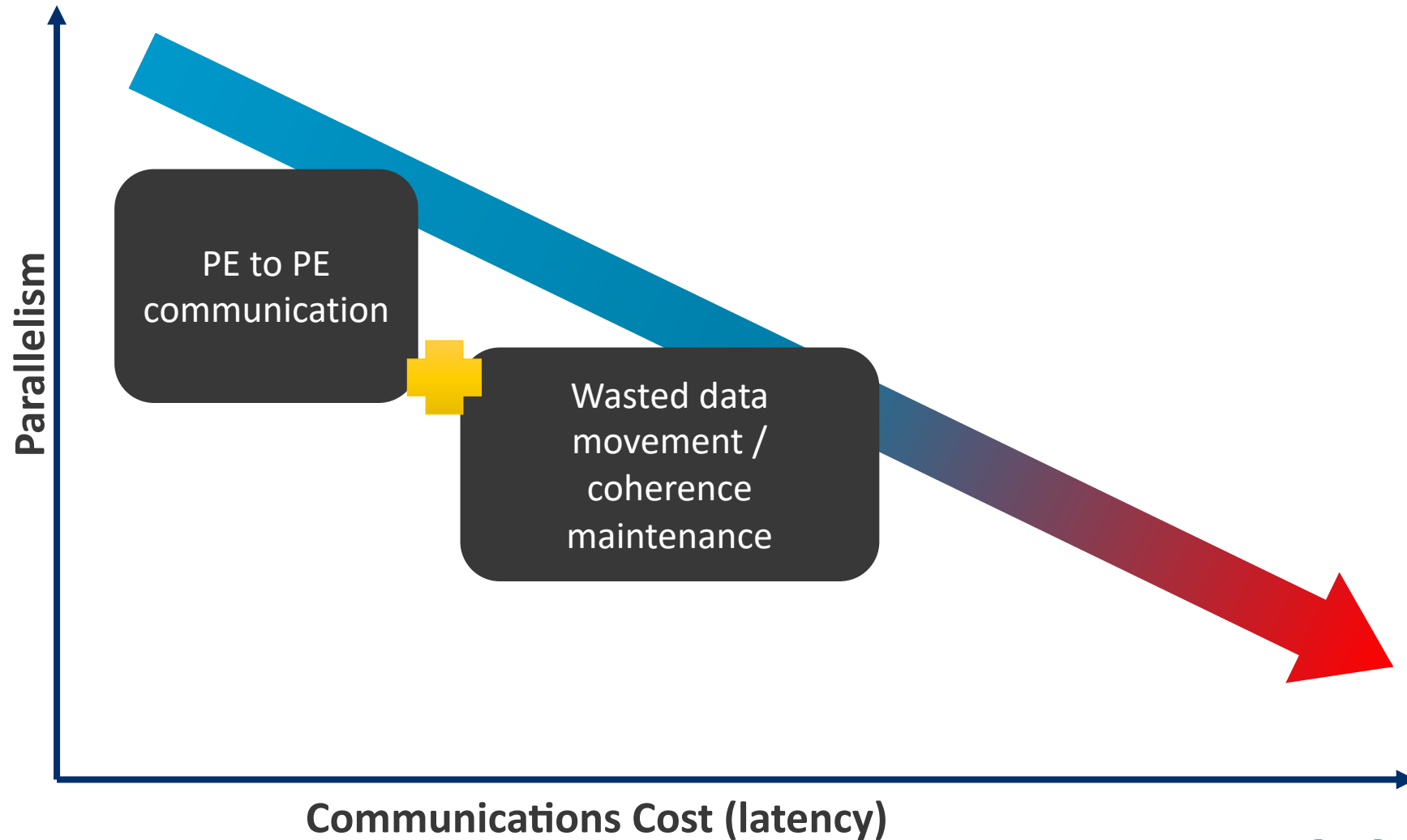


**Parallelism** (y-axis)

**Communications Cost (latency)** (x-axis)

**arm** Research

# Compute efficiency post-Moore

Data movement dominates, parallelism is critical

**Parallelism** (y-axis)

PE to PE communication

**Communications Cost (latency)** (x-axis)

**arm** Research

# Compute efficiency post-Moore

Data movement dominates, parallelism is critical



**Parallelism**

PE to PE communication

Wasted data movement / coherence maintenance

**Communications Cost (latency)**

**arm** Research

# Compute efficiency post-Moore

Data movement dominates, parallelism is critical



**Parallelism**

PE to PE communication

Wasted data movement / coherence maintenance

Context creation/dispatch time

**Communications Cost (latency)**

**arm** Research

# Compute efficiency post-Moore
## Data movement dominates, parallelism is critical

**Fine-grained parallelism**

Parallelism

PE to PE communication

Wasted data movement / coherence maintenance

Context creation/dispatch time

**Coarse-grained parallelism**

**Communications Cost (latency)**

**arm** Research

# Compute efficiency post-Moore

Data movement dominates, parallelism is critical



**Fine-grained parallelism**

PE to PE communication

Wasted data movement / coherence maintenance

Context creation/dispatch time

BTW, we must make this programmable

**Coarse-grained parallelism**

**Communications Cost (latency)**

**arm** Research

Most software designed for this

5 years ago

Today

Today – Mobile/Client
0-3 years

3-10 years

Towards Extreme Heterogeneity

Low          Cost to build/adopt/run          Extreme

*Adapted/modified from original figure courtesy of Dilip Vasudevan*

arm Research

# Bottom line up front

Just in case I get kicked off the stage before I finish…..

- Problem we're really trying to solve is always data movement

  - Context to PE

  - PE communicating results to PE

  - Creating new context from parent PE

  - PE storing results

  - Aligning instruction/command data with input data

- PINM is just another accelerator, but not one we should tackle first.

- We have to face our inconvenient truths……

**arm** Research

# Bottom line up front
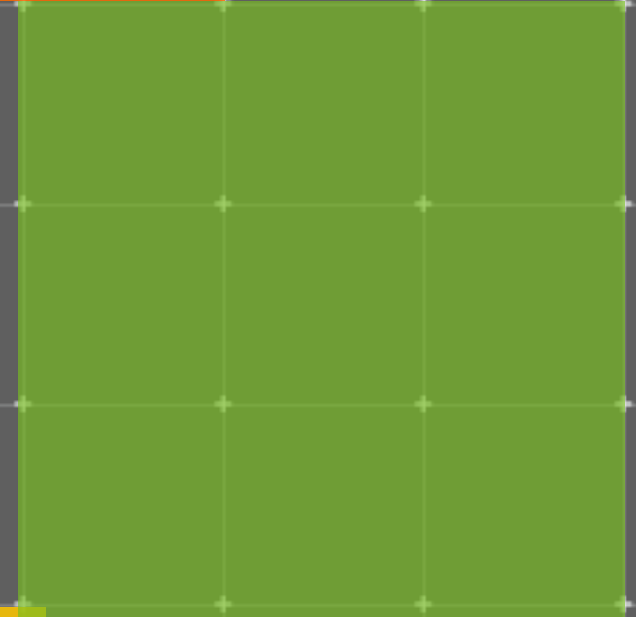Just in case I get kicked off the stage before I finish…..

- Most PINM solutions often have issues with

    - VA -> PA Translation / interleaving (bank/channel/etc.)

    - Programmability

    - Cache maintenance operations? Where?

    - In-NVM compute, what happens when cells die? Interaction with wear-leveling??

    - Exceptions, error handling?

    - Synchronization: between PINM units and with host cores

    - Working set size vs. device size….thread migration is needed (some have solutions, do others?)

- There are no magic memories

    - If it sounds too good to be true, it usually is.

arm Research

# Bottom line up front

Just in case I get kicked off the stage before I finish…..

- Previous slide is a tad depressing…

- Let's talk about some easier opportunities….

© 2019 Arm Limited

**arm** Research

# What is memory?

**arm** Research

"the faculty by which the mind stores and remembers information"

- Apple Dictionary

**arm** Research

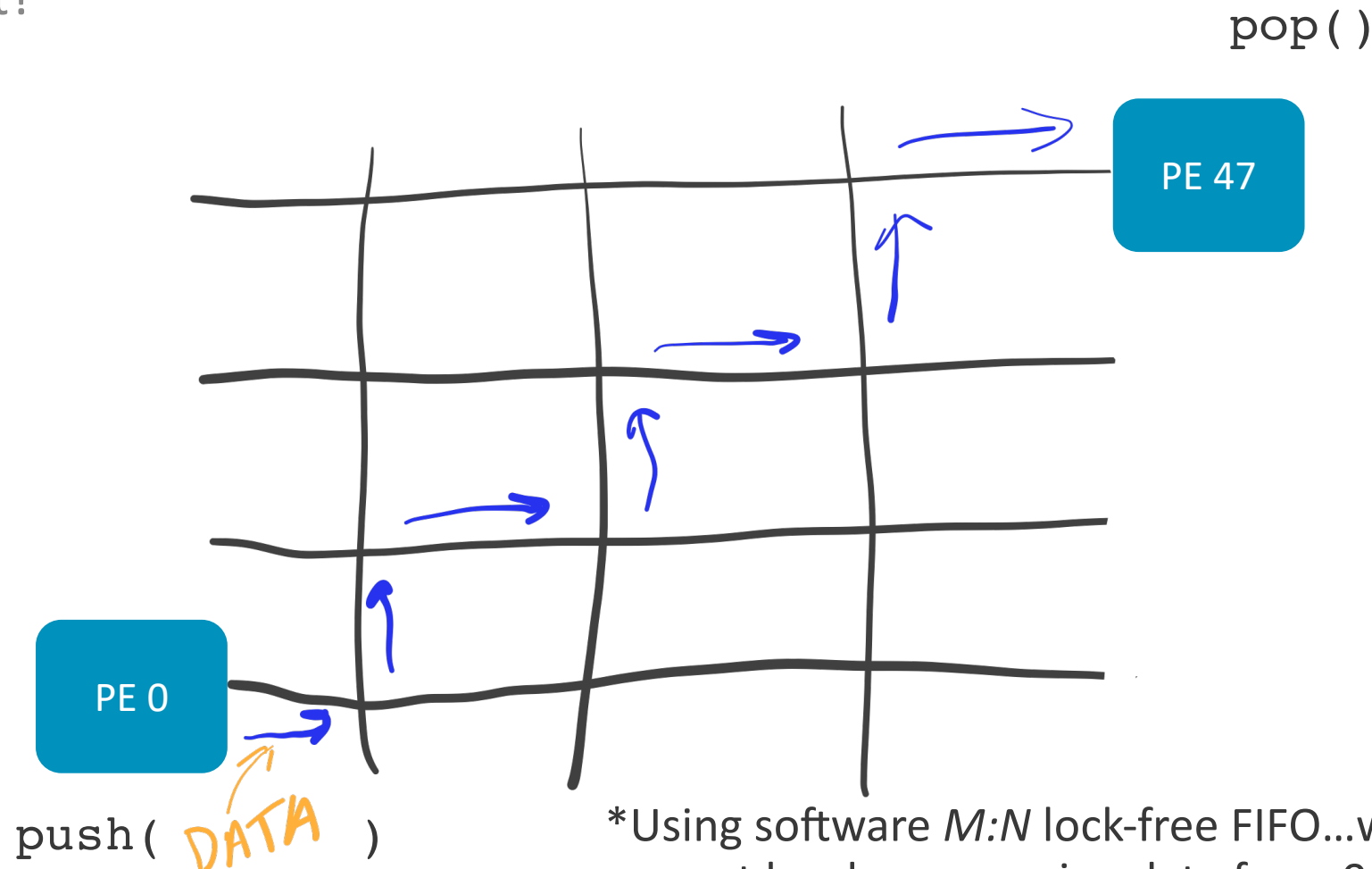"the faculty by which the mind stores and remembers information"

- Apple Dictionary

"the faculty by which the
application ~~mind~~ stores and remembers information"

- Apple Dictionary, edited ☺

**arm** Research

# Why don't we consider the interconnect as memory too?

Maybe...it should be, it is right?

pop()

- It stores memory right? Even if only a few cycles at a time.

- Interconnect is way cheaper than DRAM (energy/latency)

- Keep data within interconnect when possible.

- Why is this so hard to do? Some networking cores have, why not general purpose?
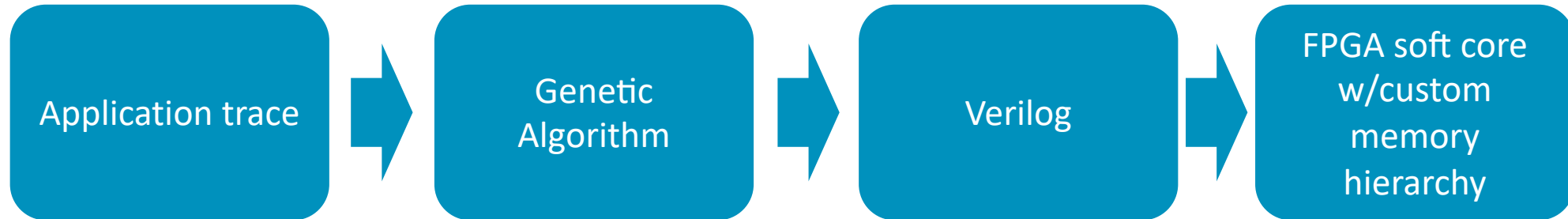
PE 47

PE 0

push( DATA )

*Using software *M:N* lock-free FIFO...with current hardware moving data from 0 to 47 takes about **500** cycles ☹. Could take as few as **10-100**.
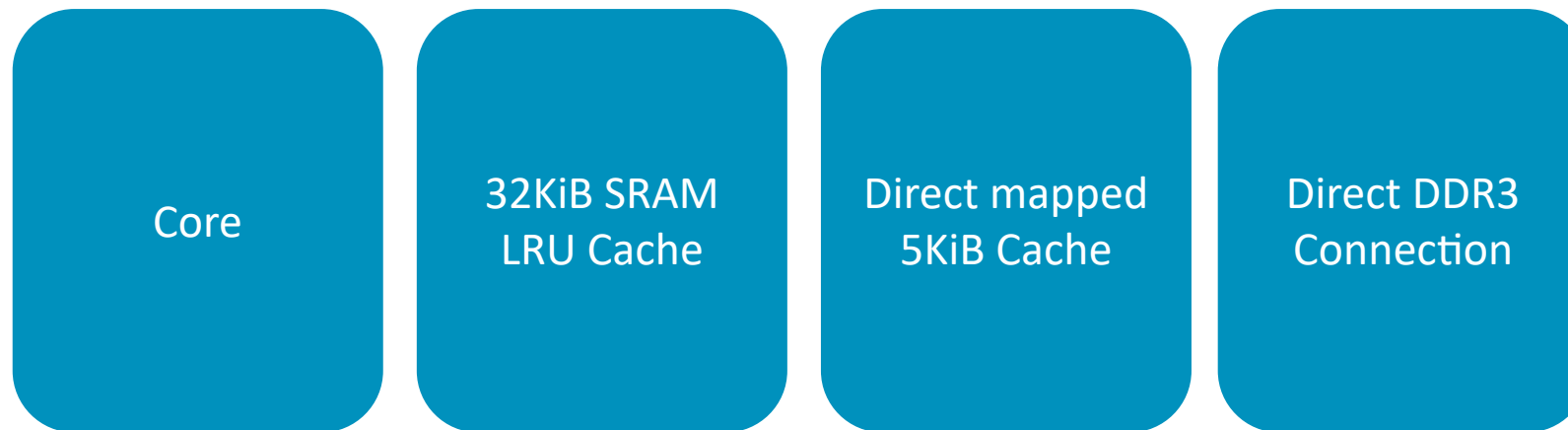
© 2019 Arm Limited

arm Research

# Accelerator != Just Logic

**arm** Research

# Customized memory hierarchy

| Application trace | → | Genetic Algorithm | → | Verilog | → | FPGA soft core w/custom memory hierarchy |

**There are some strange combinations**

| Core | 32KiB SRAM LRU Cache | Direct mapped 5KiB Cache | Direct DDR3 Connection |

**But…2-10x speedup on FPGA, up to 100x when built as ASIC (higher clock rates)**

© 2019 Arm Limited

**arm** Research

# Customized memory hierarchy

Application trace

Genetic
~~~~~~

FPGA soft core
/custom
memory

Core

~~KiB~~
LRU Cache

Direct mapped
5KiB Cache

Direct DDR3
Connection

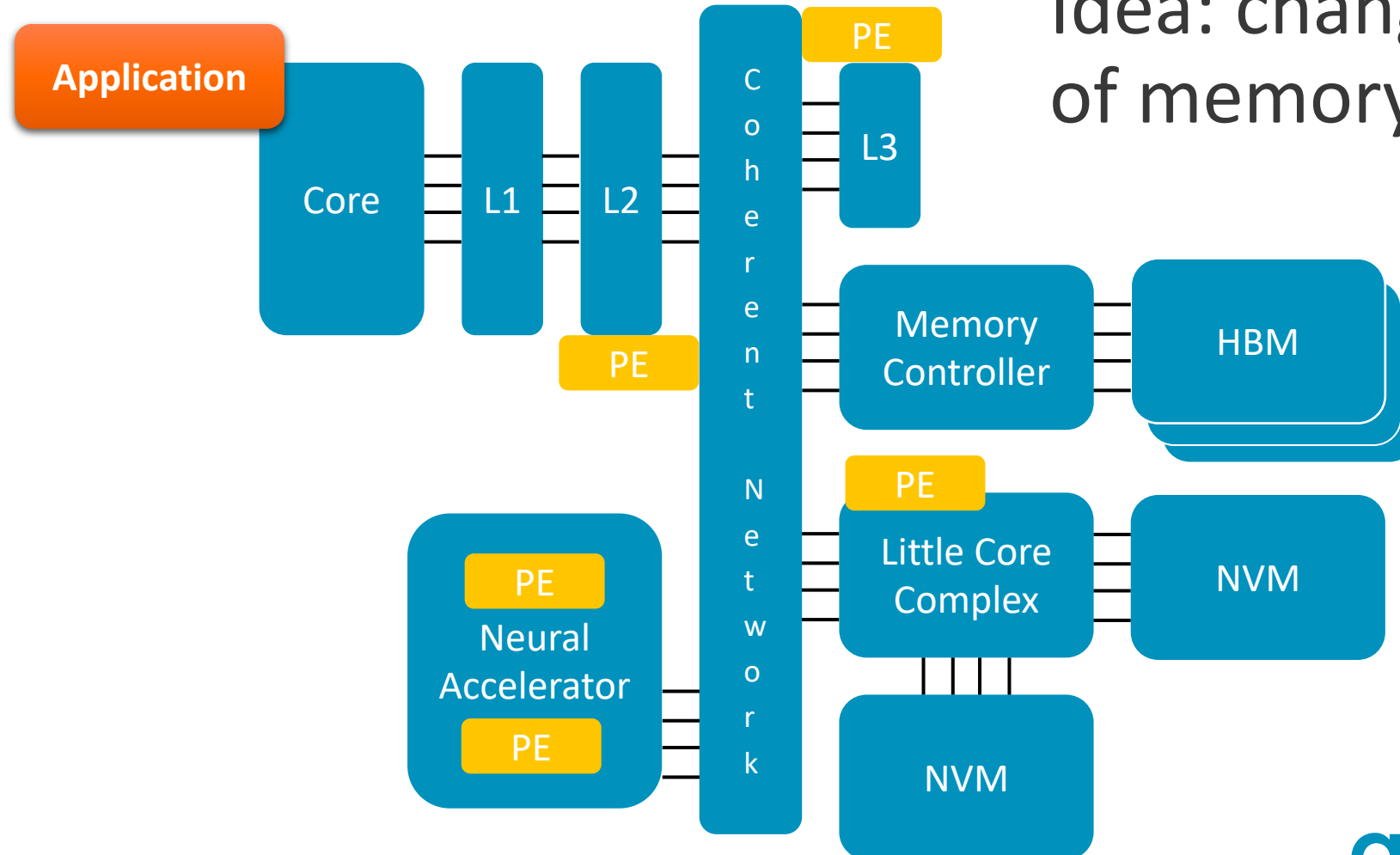**Not really practical**

arm Research

# Relativity

Locality is, from a certain point of view.

**Observation:** Manipulating distance by moving processor can expand locality and decrease latency for some workloads

L3

NVM

PE

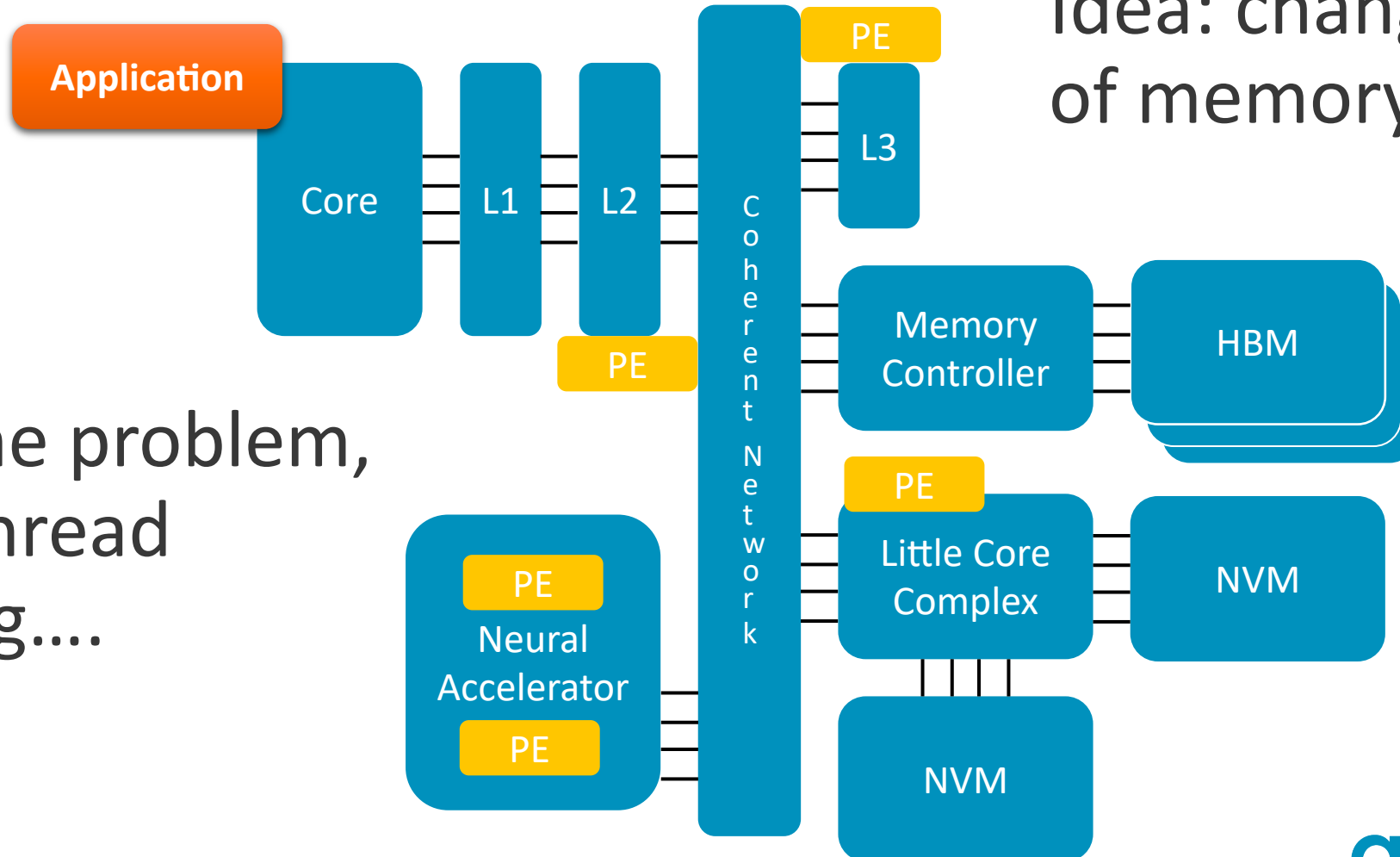Data X

L2

**~100 cycles**

`load x`

L1

**184 cycles**

**45% reduction in latency if no reuse**

Core

`load x`

**arm** Research

# Flip the script

Instead of changing memory hierarchy, add processors (PEs) everywhere



Idea: change "view" of memory hierarchy

© 2019 Arm Limited

arm Research

# Flip the script

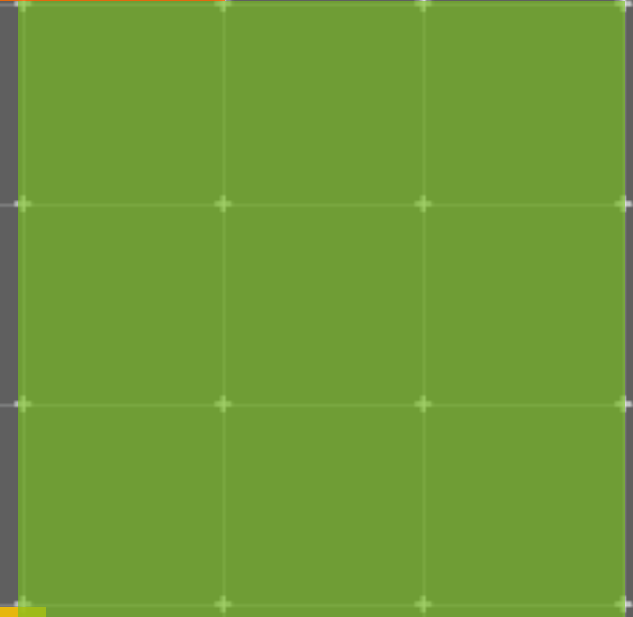Instead of changing memory hierarchy, add processors (PEs) everywhere

Idea: change "view" of memory hierarchy

Change the problem, it's now thread scheduling….

Application

Core

L1

L2

Coherent Network

PE

L3

PE

Memory Controller

HBM

PE

Little Core Complex

NVM

PE

Neural Accelerator

PE

NVM

**arm** Research

# Do you really want to program that??

**arm** Research

# Just one problem among many..

Function

OS

Host Core

Controller

**?**

Where to find the data?

Data

Controller

**Memory Device**

| | | | |
|---|---|---|---|
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |
| Core | Core | Core | Core |

arm Research

# The interface

**arm** Research

# Meet "Bob"….

**arm** Research

HALF FOODS
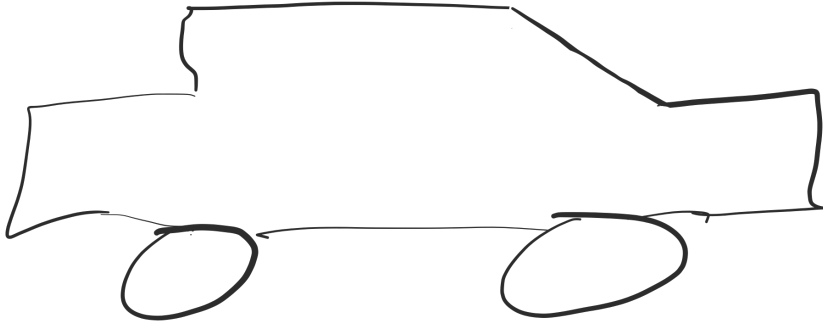
24 hour MART

Bob's Son "Bobby"

**arm** Research

# Interfaces

**arm** Research

# Interfaces

# Interfaces

arm Research

# Interfaces

　© 2019 Arm Limited

# But for programming computer hardware....



Active Memory

cnFET Acc
PoP

CMOS Acc 3

CMOS Acc 2

L1 Cache

Neuromorphic MRAM

Bus

Stacked SRAM

VM Interface

PCIe-G6

NVM Interface

Interface

CPU

CPU

**Build Up (3D)**

Compute or Memory depending on use case

**Build Over (multi-chip modules)**

**arm** Research

# We need intuitive, productive interfaces….

# But….

Pragmas / specialized APIs

RDMA

NUMA Page Move

mlock

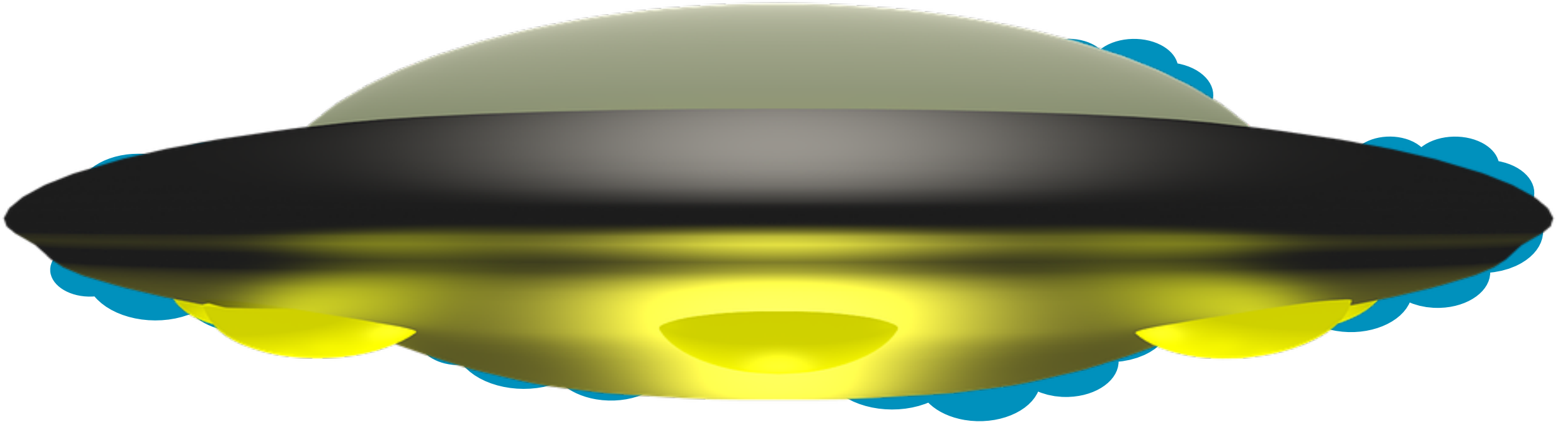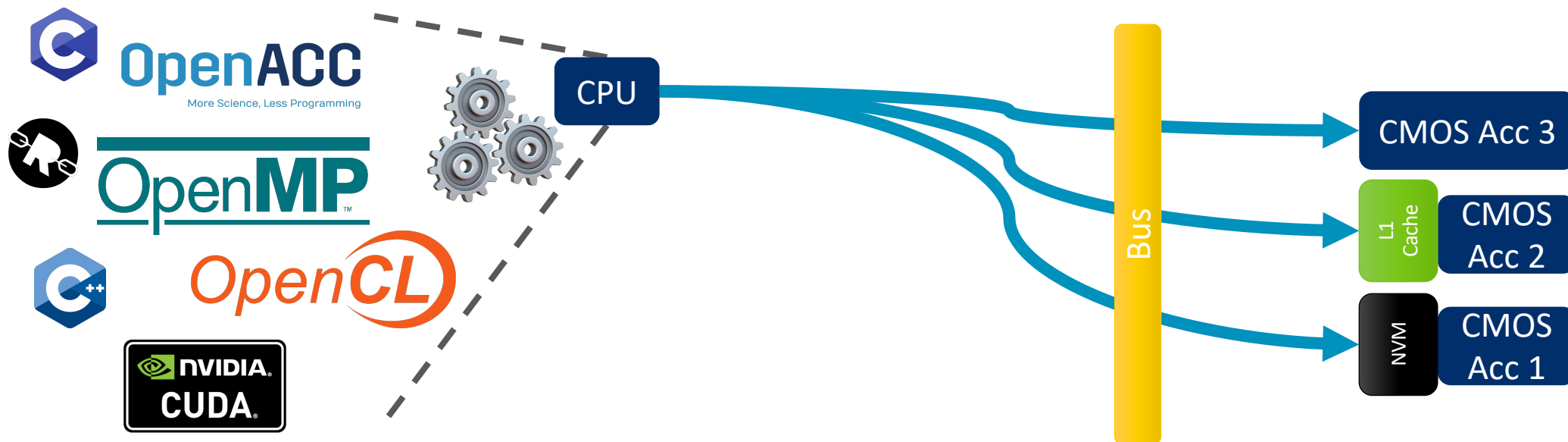Set Affinity

intrinsics

arm Research

# But….



**Do we need super advanced aliens to help us build programs??**

(after all, they did help build the pyramids right?? ;) )

*(disclaimer: aliens didn't actually build the pyramids…)*

arm Research

Most "start-up" hardware vendors fail…because of software, or lack of.

arm Research

# Mature Software Ecosystem

**OpenACC** — More Science, Less Programming

**OpenMP**™

**OpenCL**

C++

**nVIDIA CUDA**®

Market success isn't based on having the "best" hardware, it's in having a broad software user-base....more simply, it's the software!!

CPU

Bus

CMOS Acc 3

L1 Cache — CMOS Acc 2

NVM — CMOS Acc 1

arm Research

# What if....

Accelerators built to common interface....



© 2019 Arm Limited

arm Research

# What if....

**Mature Software Ecosystem**

**Accelerators built to common interface....**

OpenACC — More Science, Less Programming

OpenMP™

C++  OpenCL

NVIDIA CUDA

CPU → Common Interface (middleware) → Bus → CMOS Acc 3

L1 Cache — CMOS Acc 2

NVM — CMOS Acc 1

arm Research

# What if….

**Mature Software Ecosystem**

**Accelerators built to common interface….**



- **Build interface to make it easier to determine data locality, do fast dispatch, etc.!**
- **Lower offload latency == more tasks / unit time**

arm Research

# Rethinking OS/Hardware interfaces

Application

Virtualization Layers

Learning Agent (s)

ML / AI / Heuristic / Analytic Model

Context Management

**arm** Research

# Rethinking OS/Hardware interfaces

# Rethinking OS/Hardware interfaces

Application

Virtualization Layers

Learning Agent (s)

ML / AI / Heuristic / Analytic Model

Context Management

Security

Data Management

**arm** Research

# The AI/ML Assisted System

(shh, it's not really aliens)



$$\frac{\lambda^2}{\mu(\mu-\lambda)} \quad \frac{\alpha(\alpha+1)\beta^2\lambda^2}{2(1-\alpha\beta\lambda)}$$

**Injected Models**

**Learned Models**

**Drive**

**Models**

**Control Systems**

**Storage**

**Memory**

**Application Execution**

CMOS ACC 0 | CMOS ACC 1 | GPGPU
**Hardware**

© 2019 Arm Limited

**arm** Research

# Parting thoughts…

- Application performance is all about the data, and how fast you can access it

- Can we build better communications primitives to reduce main memory access?

- Memory is an accelerator, memory is compute, it is not an accessory.

- We shouldn't need advanced aliens to come down to help us program our systems… ;)

**arm** Research

# Parting thoughts...

- Application performance is all about the data, and how fast you can access it

- Can we build better communications primitives to reduce memory access?

- Memory is an accelerator, it is compute, it is not an accessory.

- We shouldn't need advanced aliens to come down to help us program our systems... ;)
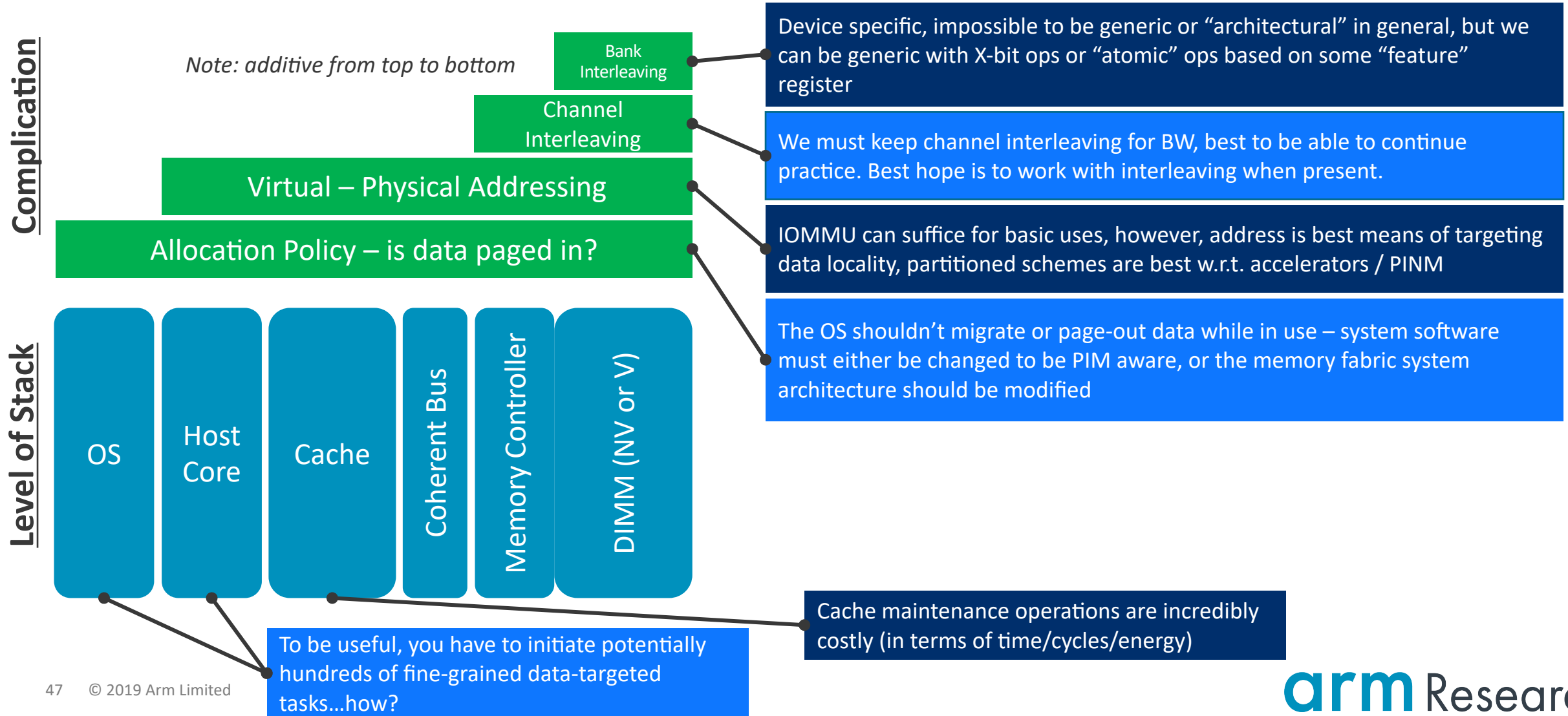
# Thanks!

arm Research

The Arm trademarks featured in this presentation are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere.  All rights reserved.  All other marks featured may be trademarks of their respective owners.

www.arm.com/company/policies/trademarks

**arm** Research

# Issues with PINM

Systems integration – translation / interleaving

**Complication**

*Note: additive from top to bottom*

**Bank Interleaving** — Device specific, impossible to be generic or "architectural" in general, but we can be generic with X-bit ops or "atomic" ops based on some "feature" register

**Channel Interleaving** — We must keep channel interleaving for BW, best to be able to continue practice. Best hope is to work with interleaving when present.

**Virtual – Physical Addressing** — IOMMU can suffice for basic uses, however, address is best means of targeting data locality, partitioned schemes are best w.r.t. accelerators / PINM

**Allocation Policy – is data paged in?** — The OS shouldn't migrate or page-out data while in use – system software must either be changed to be PIM aware, or the memory fabric system architecture should be modified

**Level of Stack**

| OS | Host Core | Cache | Coherent Bus | Memory Controller | DIMM (NV or V) |

To be useful, you have to initiate potentially hundreds of fine-grained data-targeted tasks…how?

Cache maintenance operations are incredibly costly (in terms of time/cycles/energy)

arm Research

# Issues with PINM

Integration storage – translation / interleaving

**Complication**

*Note: additive from top to bottom*

What happens when cells go bad?

All interleaving stuff from prior slide

Internally drives are further divided into channels / banks / vaults / arrays (terminology dependent on technology)

VA to Logical Block Address

Within the device, the code (VA) must be translatable to a logical block address..doable, just not easy to do w/o more new stuff.

Virtual – Physical Addressing

MMAP to INODE – Block Address

**Level of Stack**

OS

Host Core

Cache

Coherent Bus
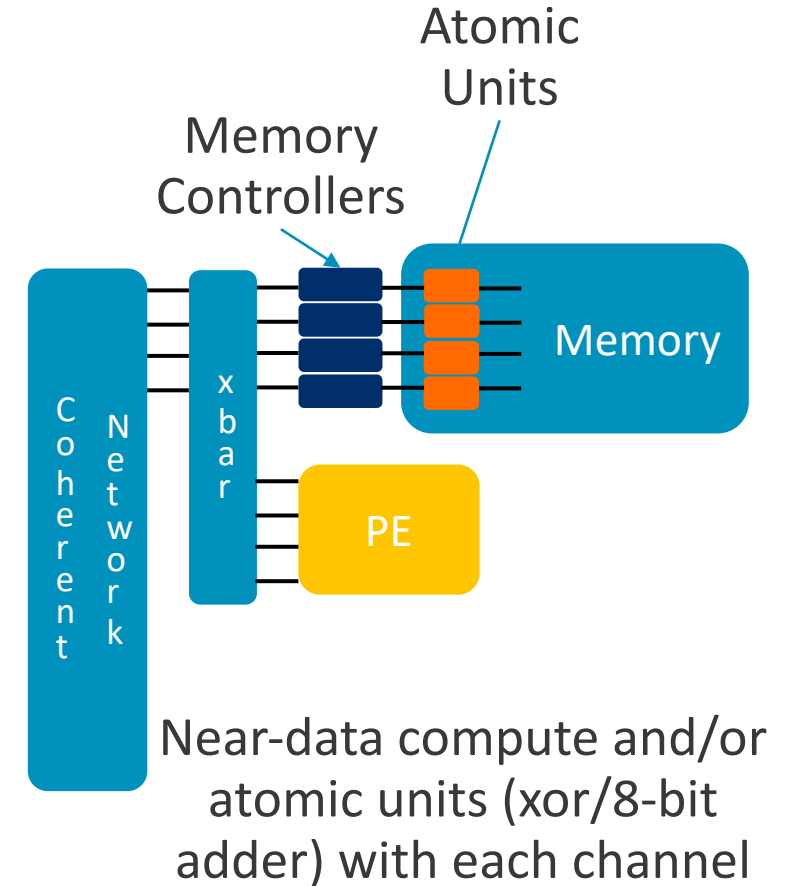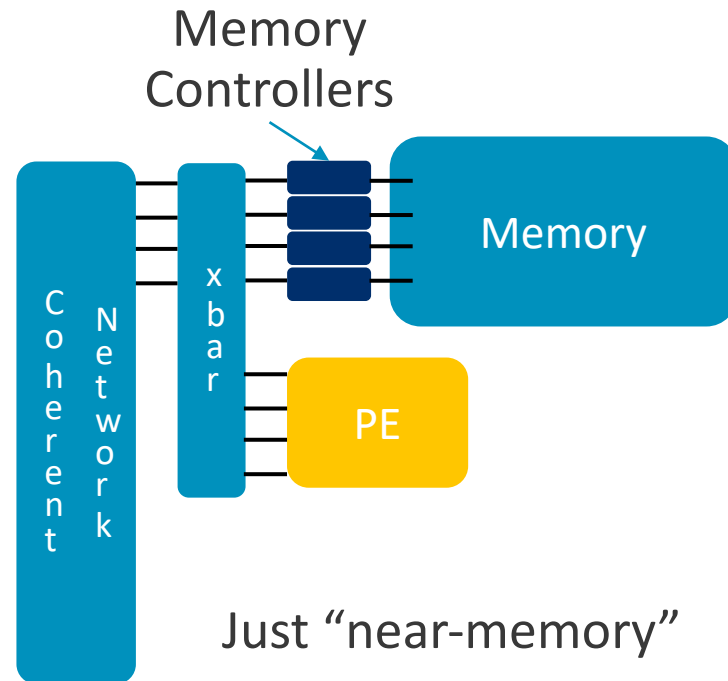
SATA or PCIe

SSD M2 / PCIe or SATA Block Device

To be useful, you have to initiate potentially hundreds of fine-grained data-targeted tasks...how?

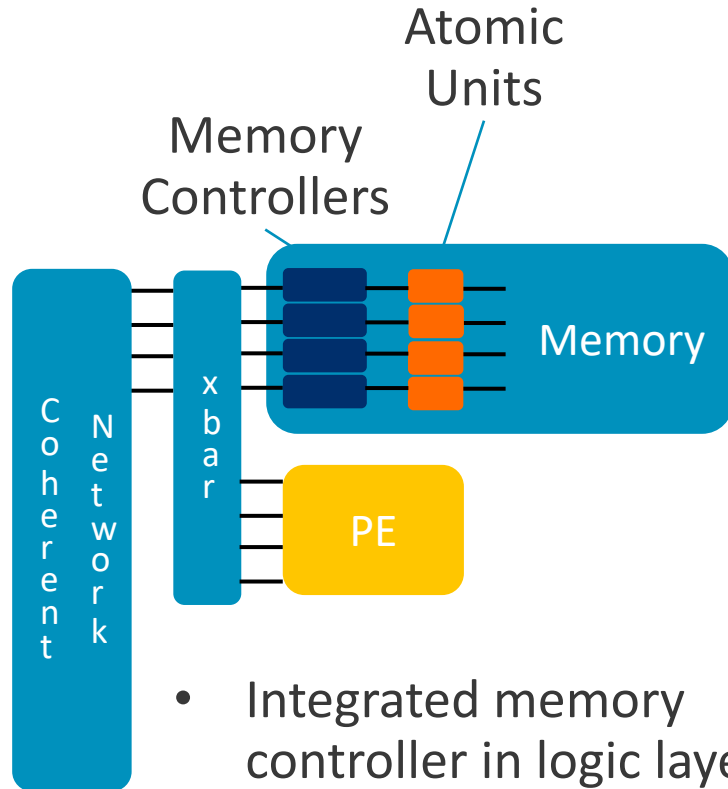Cache maintenance operations are incredibly costly (in terms of time/cycles/energy)

arm Research

# Types of PINM Accelerators

Near-memory / In-memory / In-NVM



Just "near-memory"

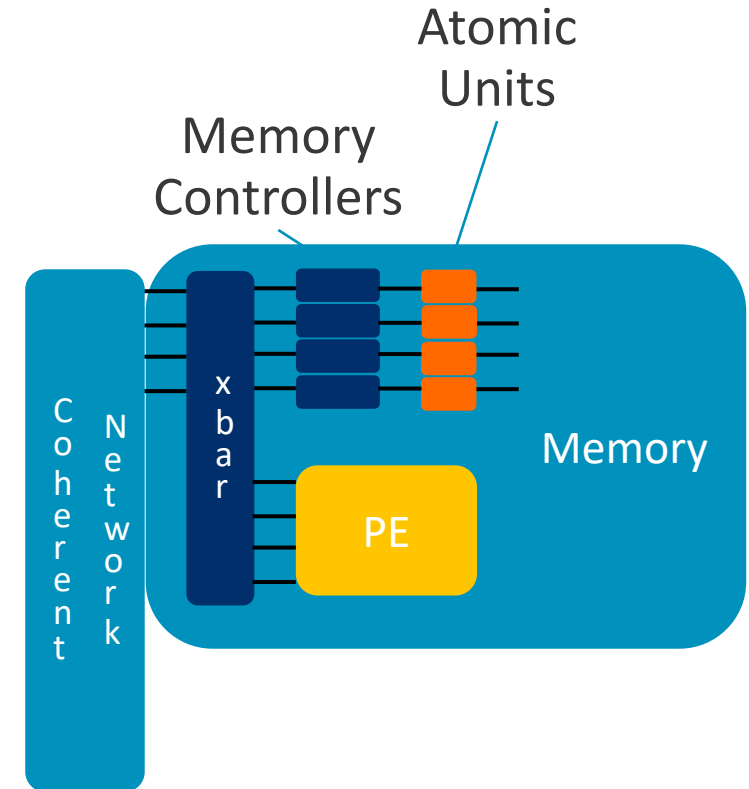Near-data compute and/or atomic units (xor/8-bit adder) with each channel

**arm** Research

# Types of PINM Accelerators

Near-memory / In-memory / In-NVM



- Integrated memory controller in logic layer…
- It doesn't make sense to move controller for this config.

- Integrated xbar + full cores near-memory in logic layer…

arm Research