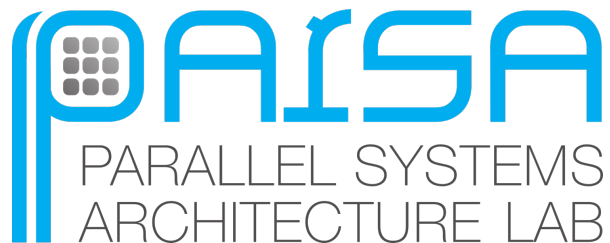


# WHAT'S HOT?

## POST-MOORE DATACENTER ARCHITECTURE

Babak Falsafi



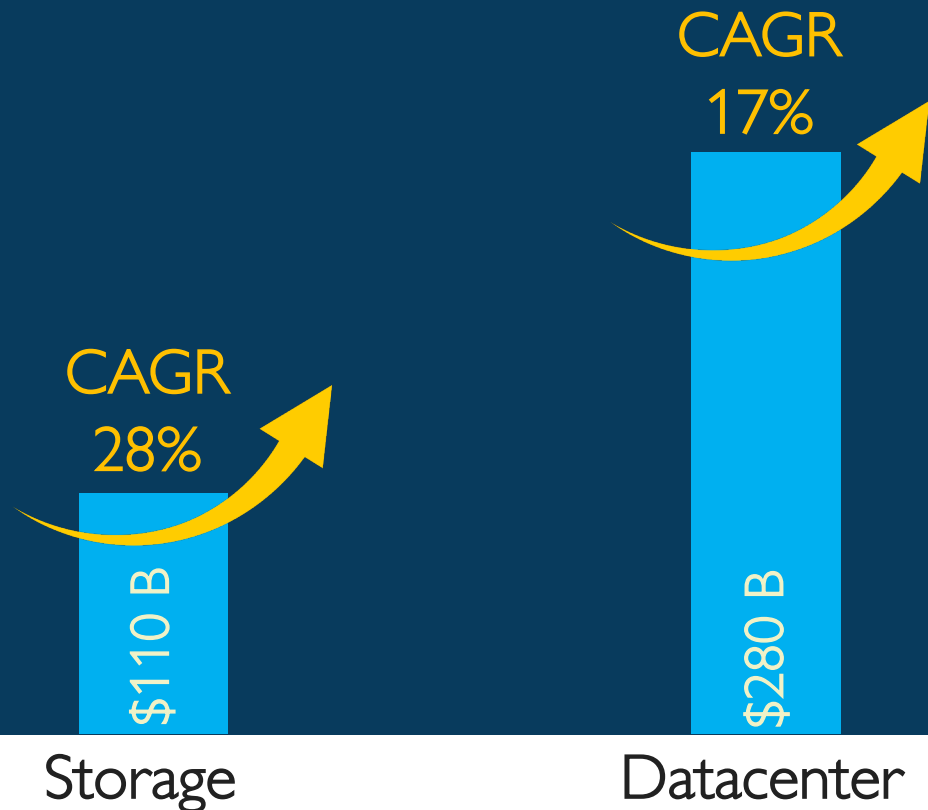
[parsa.epfl.ch](http://parsa.epfl.ch)





# DATACENTER GROWTH

Market Growth 2018-2023  
[Technavio, IDC]



- Data → fuel for digital economy
- Exponential demand for digital services
- Many apps (e.g., AI) with higher exponential demand



# DATACENTERS ARE BACKBONE OF CLOUD

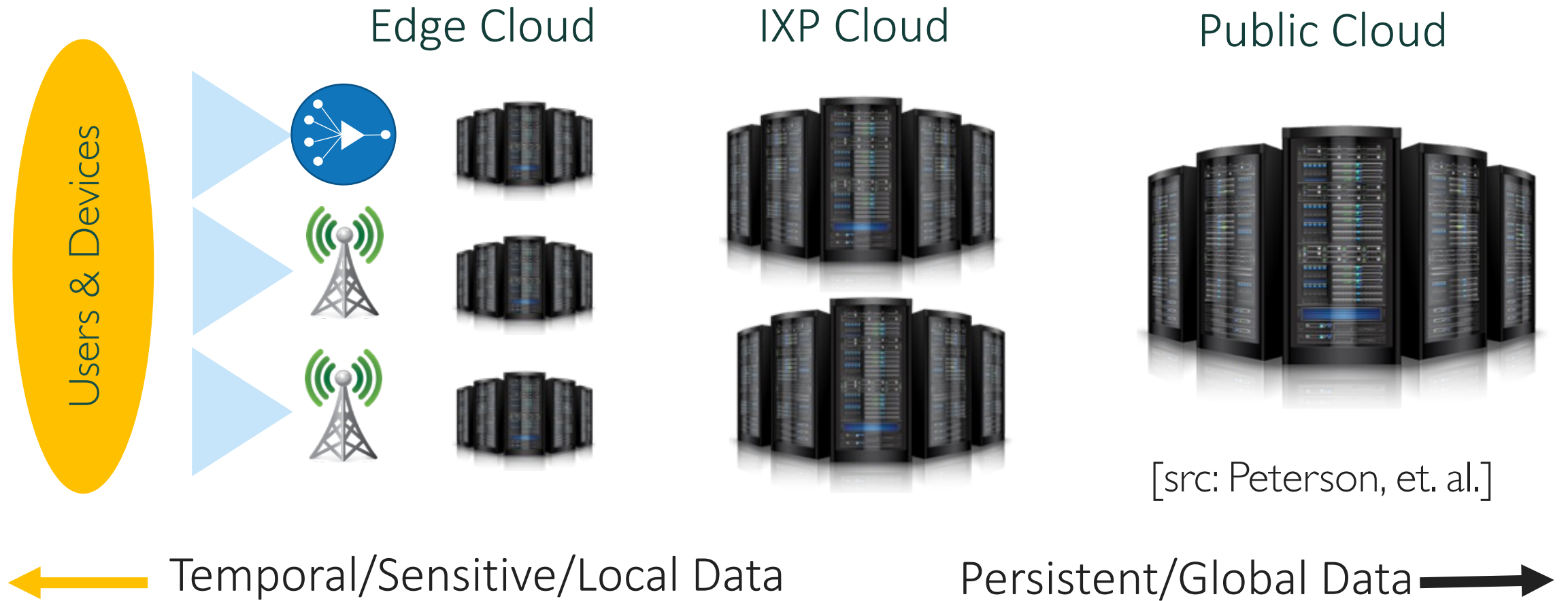
- 100s of 1000 of commodity or home-brewed servers
- Centralized to exploit economies of scale
- Network fabric w/  $\mu$ -second connectivity
- Often limited by
  - Electricity
  - Network
  - Cooling



350MW, Boydton



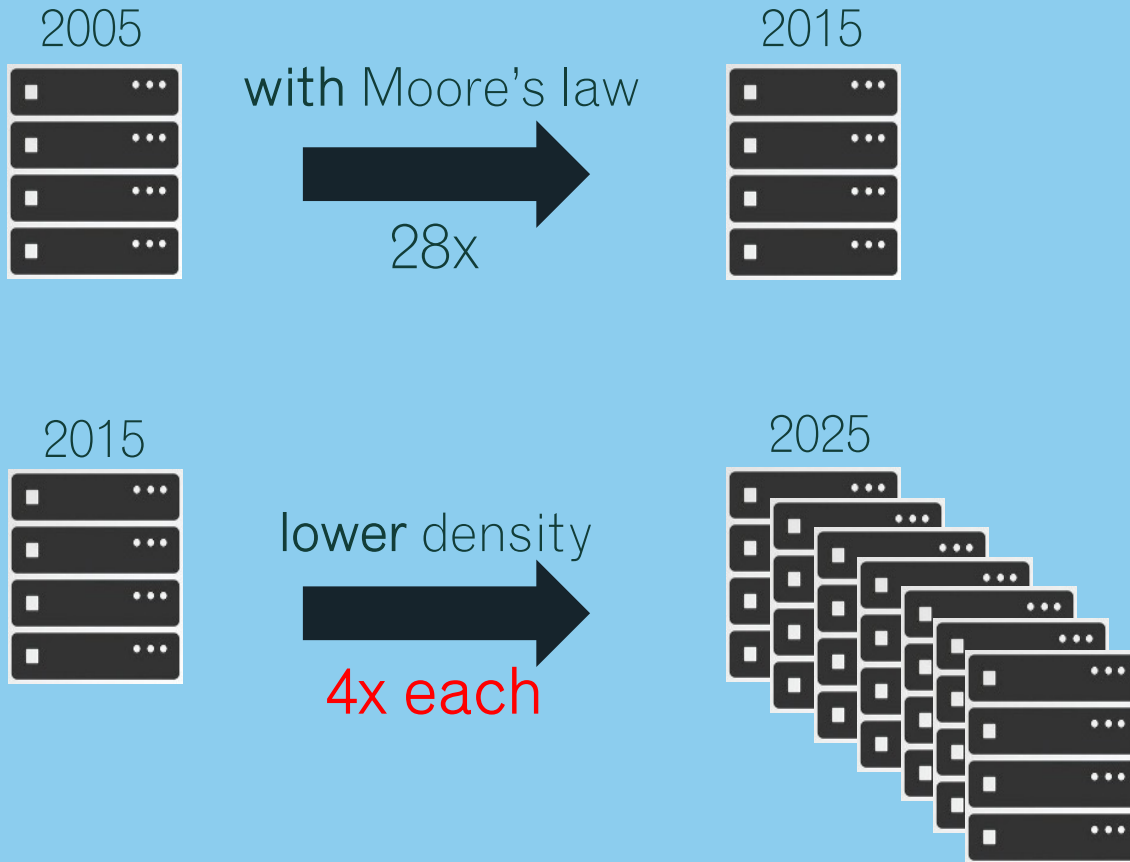
# CLOUDS AT VARIOUS SCALES





# DATACENTERS NOT GETTING DENSER

W/o Moore, building more



## End of Moore's Law (of Silicon)

- Five decades of doubling density
- Recent slowdown in density
- Chip density limited by physics

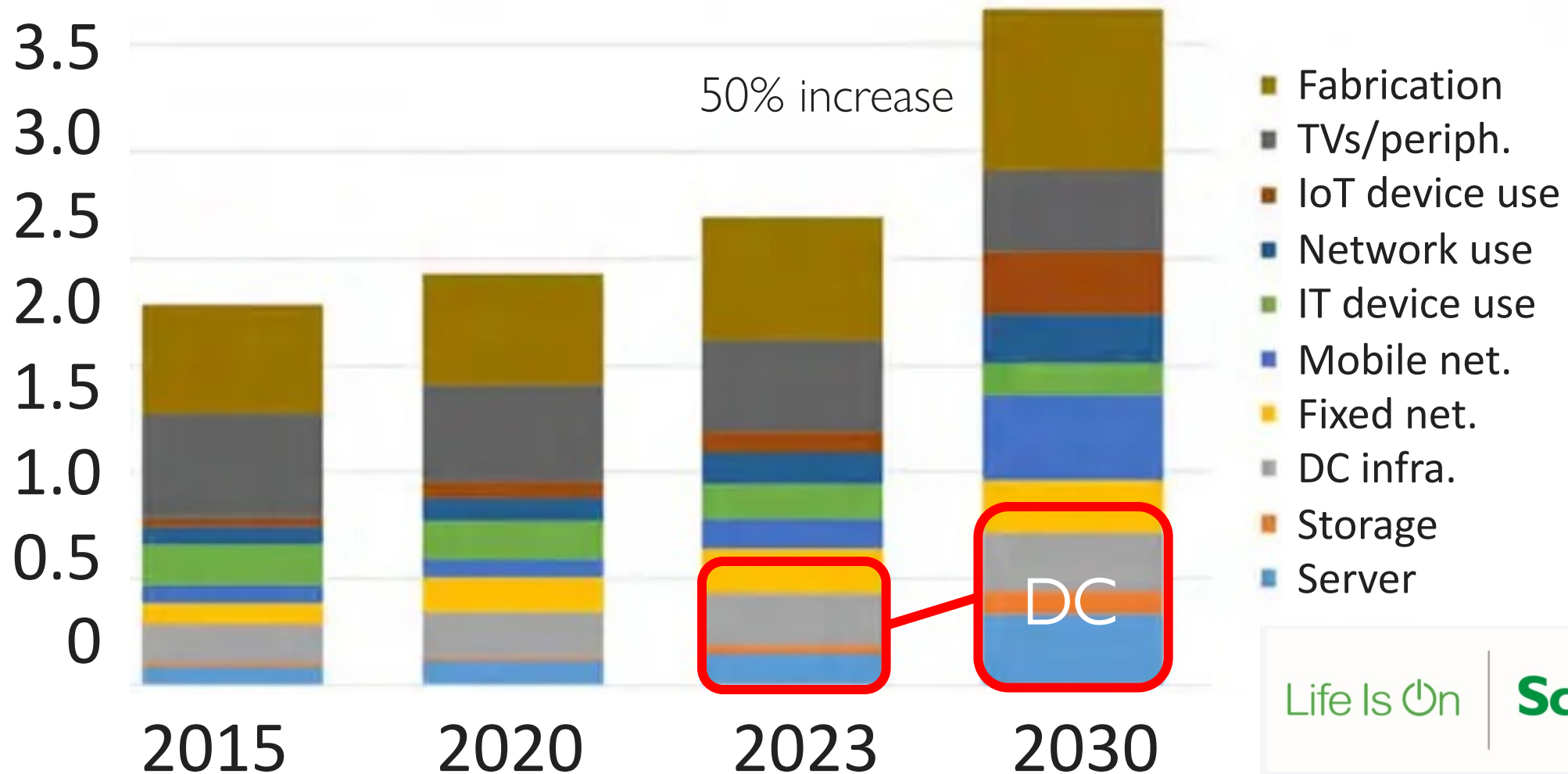
## Growth means building more

- 41%/year  $\rightarrow$  28x in ten years
- At 15%/year  $\rightarrow$  7x more DCs



# ELECTRICITY IN 1000TWH

[Digital economy & climate impact, May 2022]



Life Is On

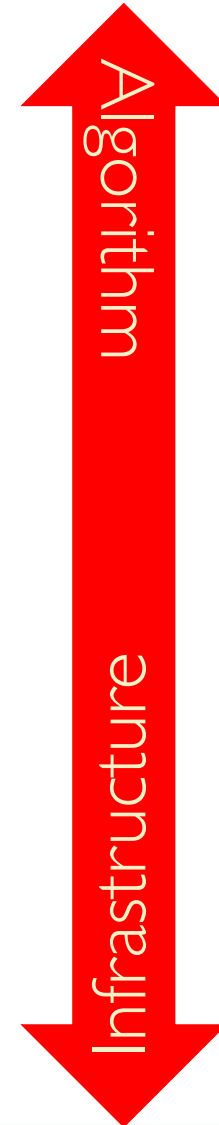
**Schneider**  
Electric



# POST-MOORE DATACENTERS

## Design for “ISA”

- Integration
  - Move data less frequently
  - Move data less distance
- Specialization
  - Customize resources
  - Less work/computation
- Approximation
  - Adjust precision





# OUTLINE

- ~~Overview~~

- Post-Moore servers

- Today's servers
- ISA opportunities

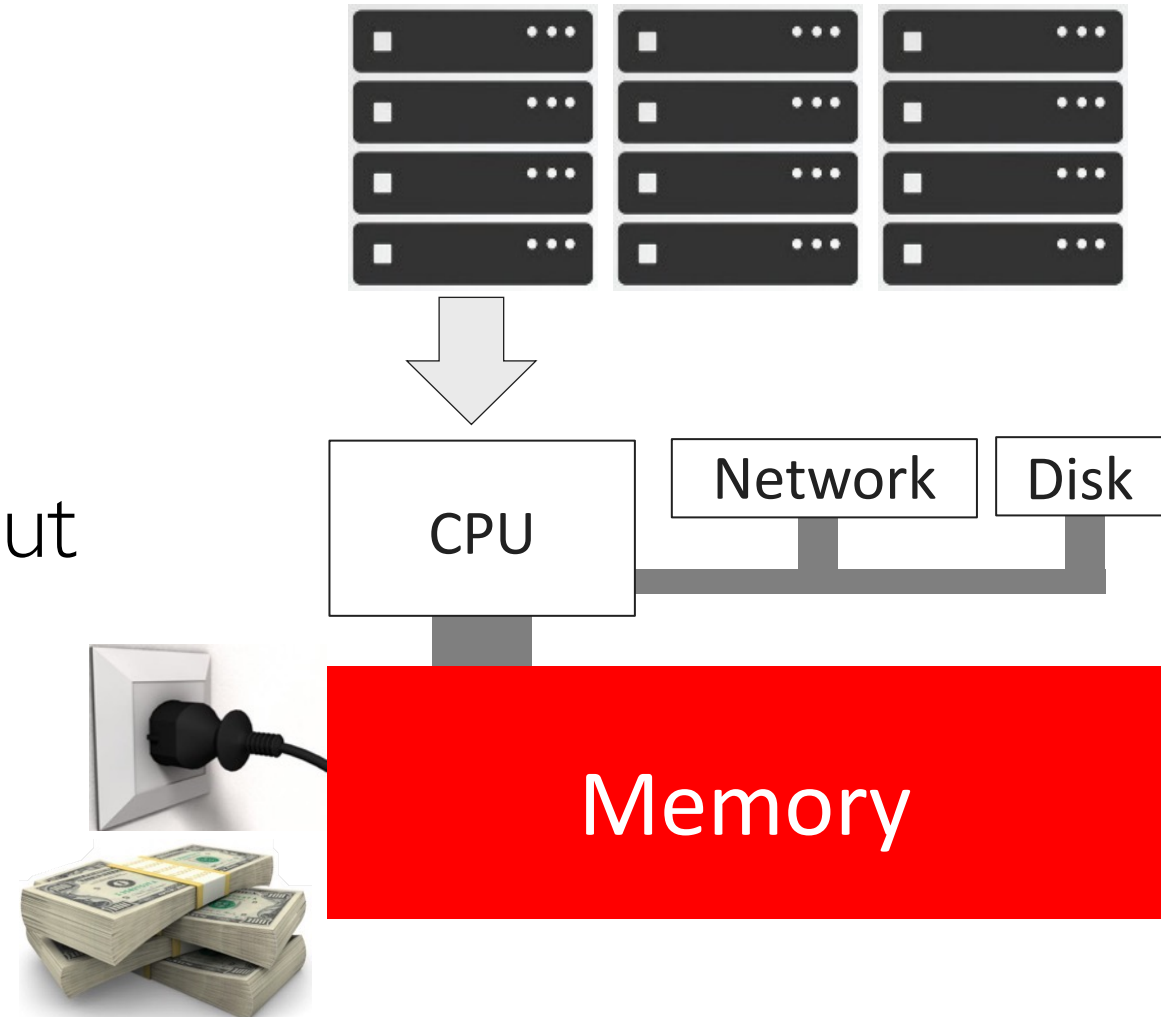
- Datacenter sustainability

- Summary



# SCALE-OUT DATACENTERS

Cost is the primary metric  
Online services hosted in memory  
Divide data up across servers  
Design server for low cost, scale out  
👉 Memory most precious silicon





# TODAY'S SERVERS

Today's platforms are PCs of the 80's

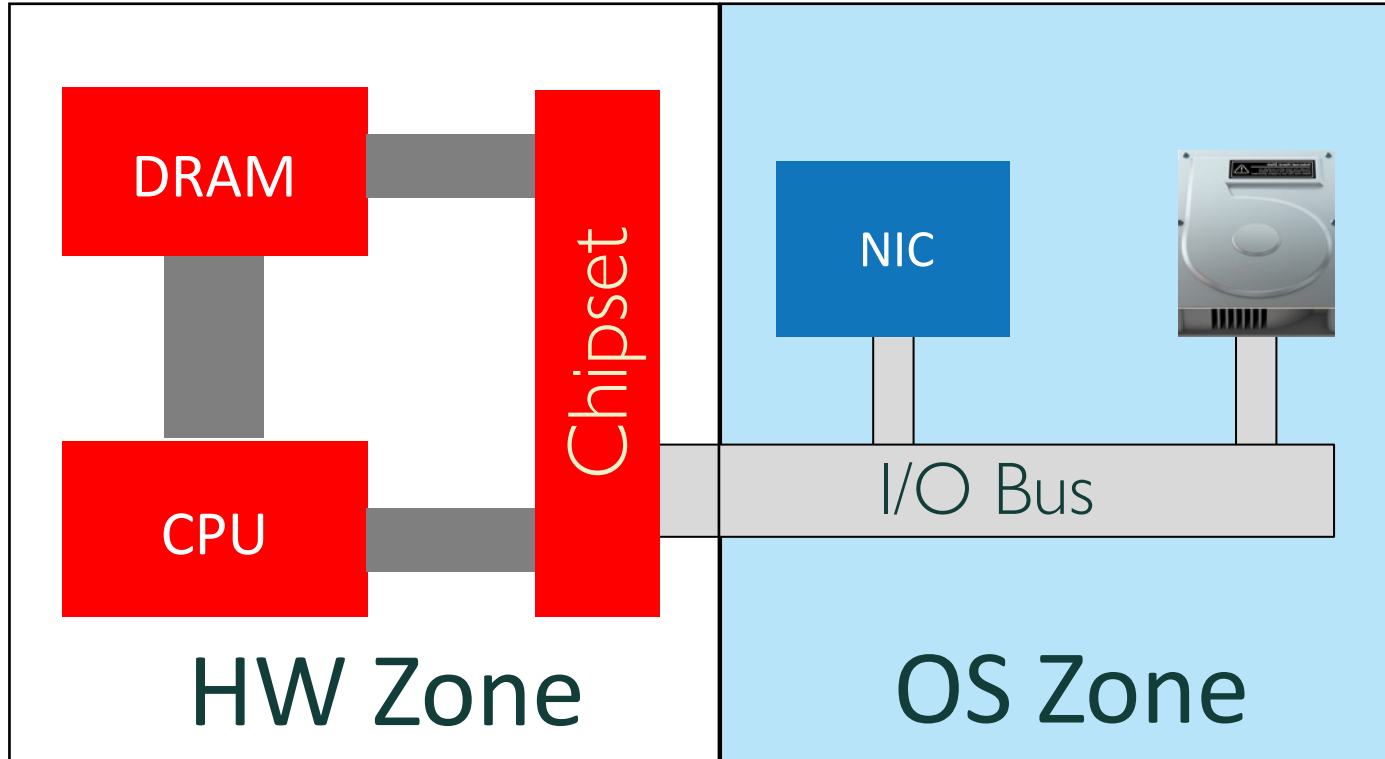
- CPU “owns” and manages memory
- OS moves data back/forth from peripherals
- Legacy interfaces connecting the CPU/mem to outside
- Legacy POSIX abstractions

Fragmented logic/memory:

- Manycore network cards w/ own memory
- Flash controllers with embedded cores and memory
- Discrete accelerators with own memory



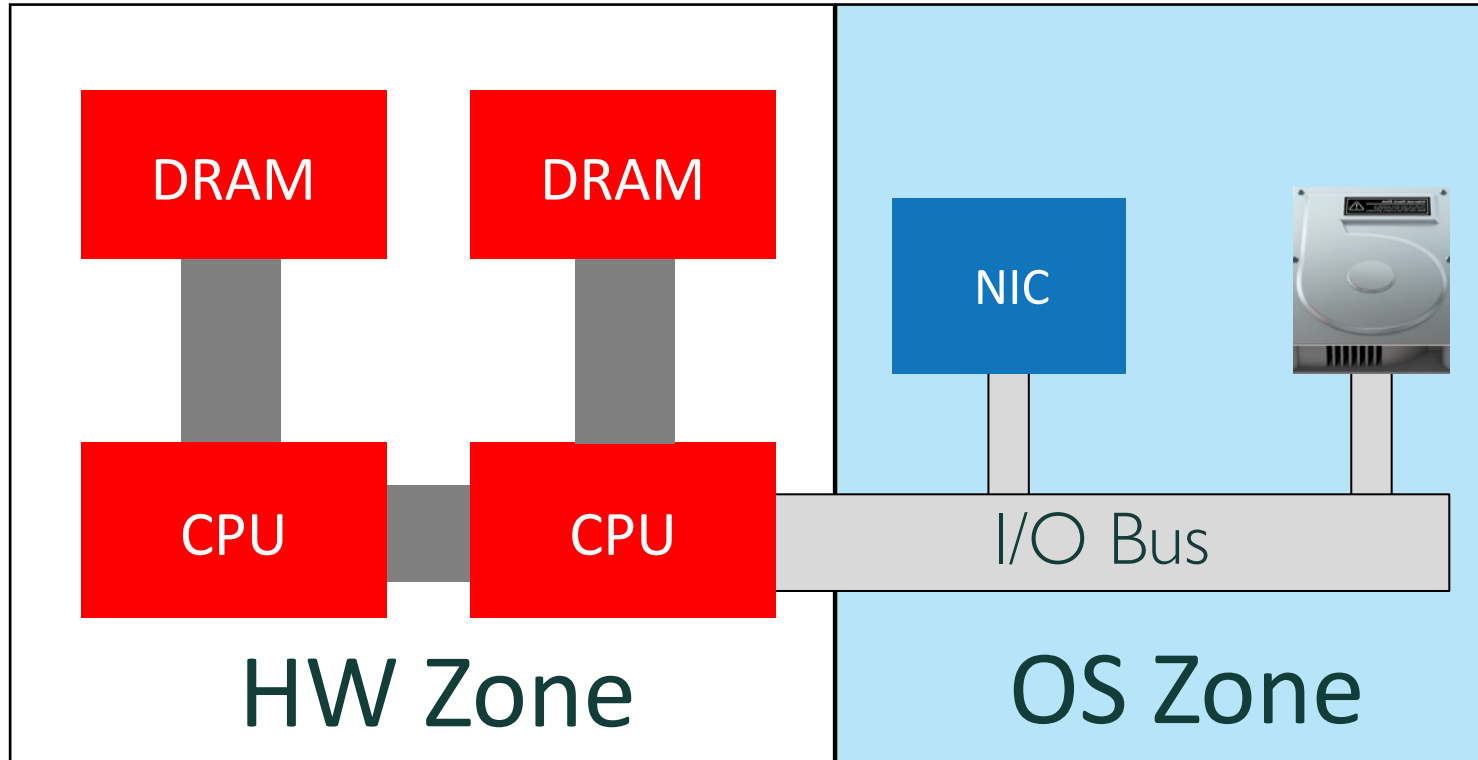
# 80S' DESKTOP



- 33 MHz 386 CPU, 250ns DRAM
- OS: Windows, Unix BSD (or various flavors)
- Focus: multiprogrammed in-memory compute



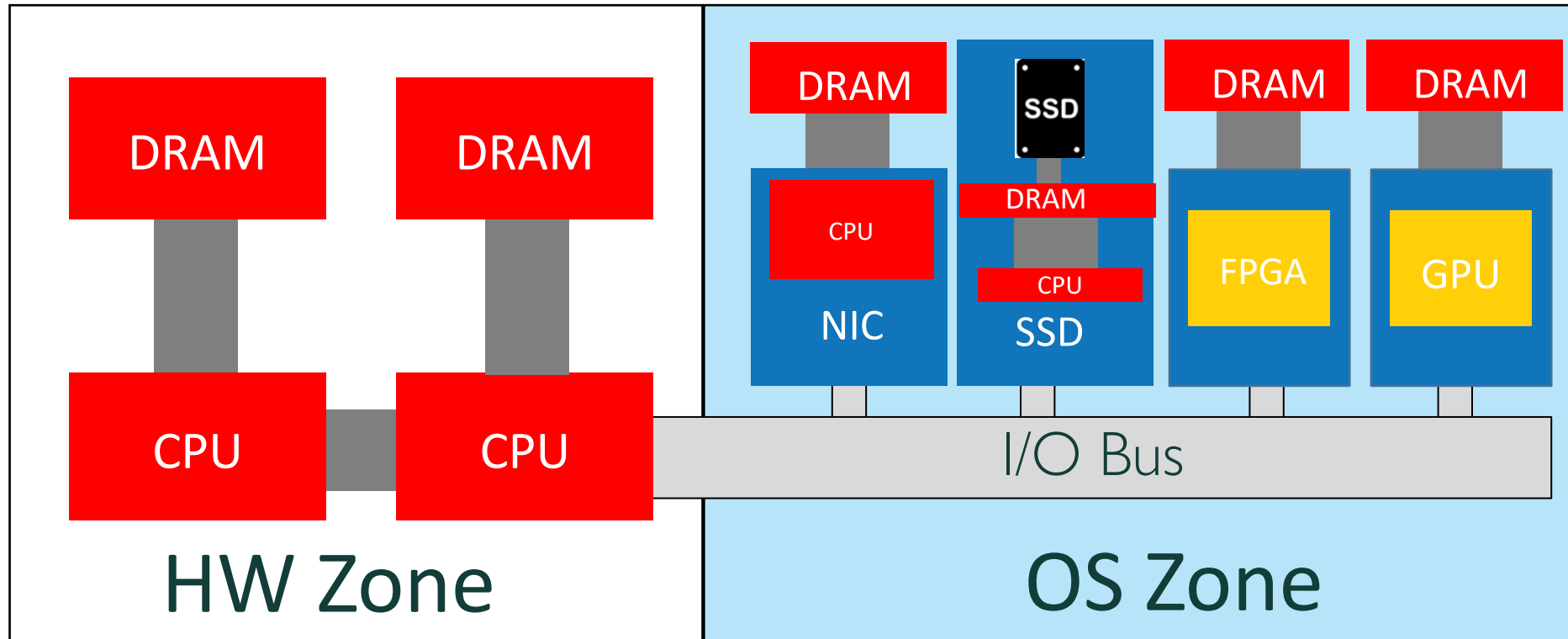
# TODAY'S SERVER: 80S' DESKTOP



- Dual 2GHz CPU's, 50ns DRAM
- OS: Linux (and various distributions)



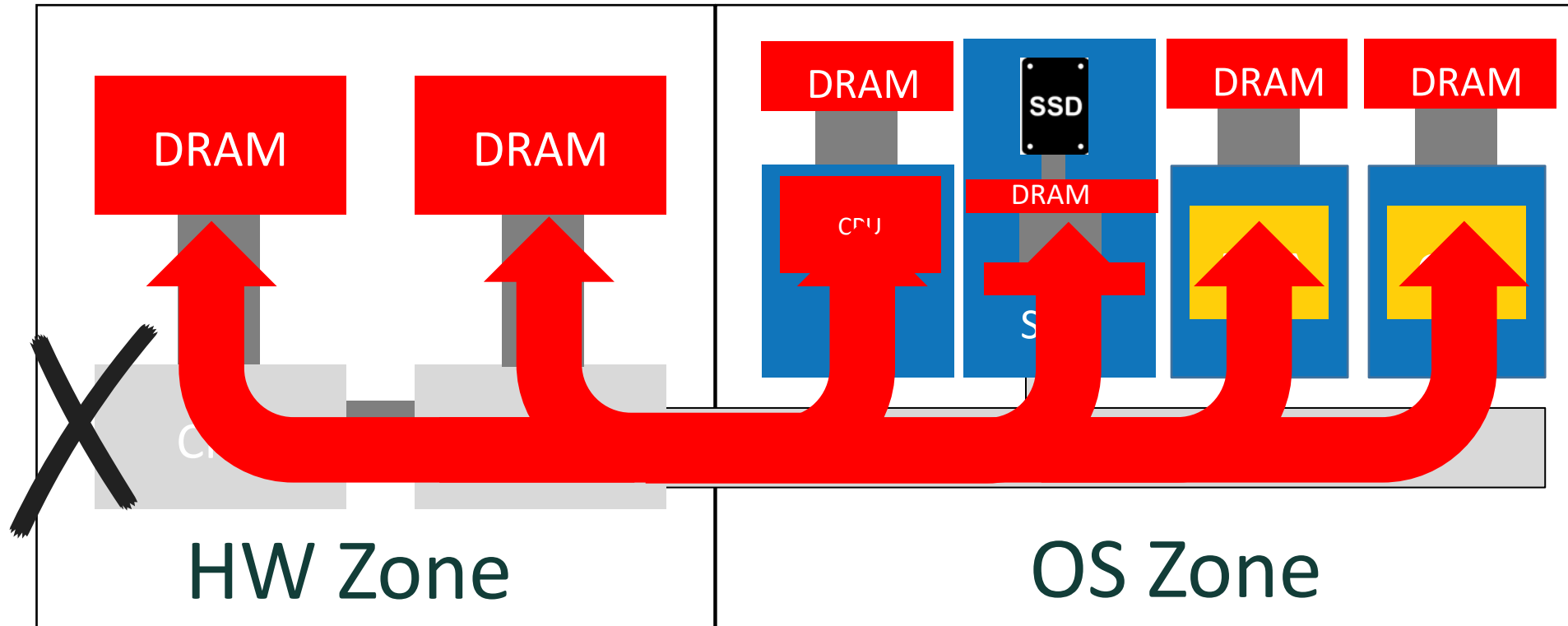
# TODAY'S SERVER: 80S' DESKTOP



- Dual 2GHz CPU's, 50ns DRAM, Linux
- Bottlenecked by legacy interfaces
- Fragmented silicon



# TODAY'S SERVER: 80S' DESKTOP



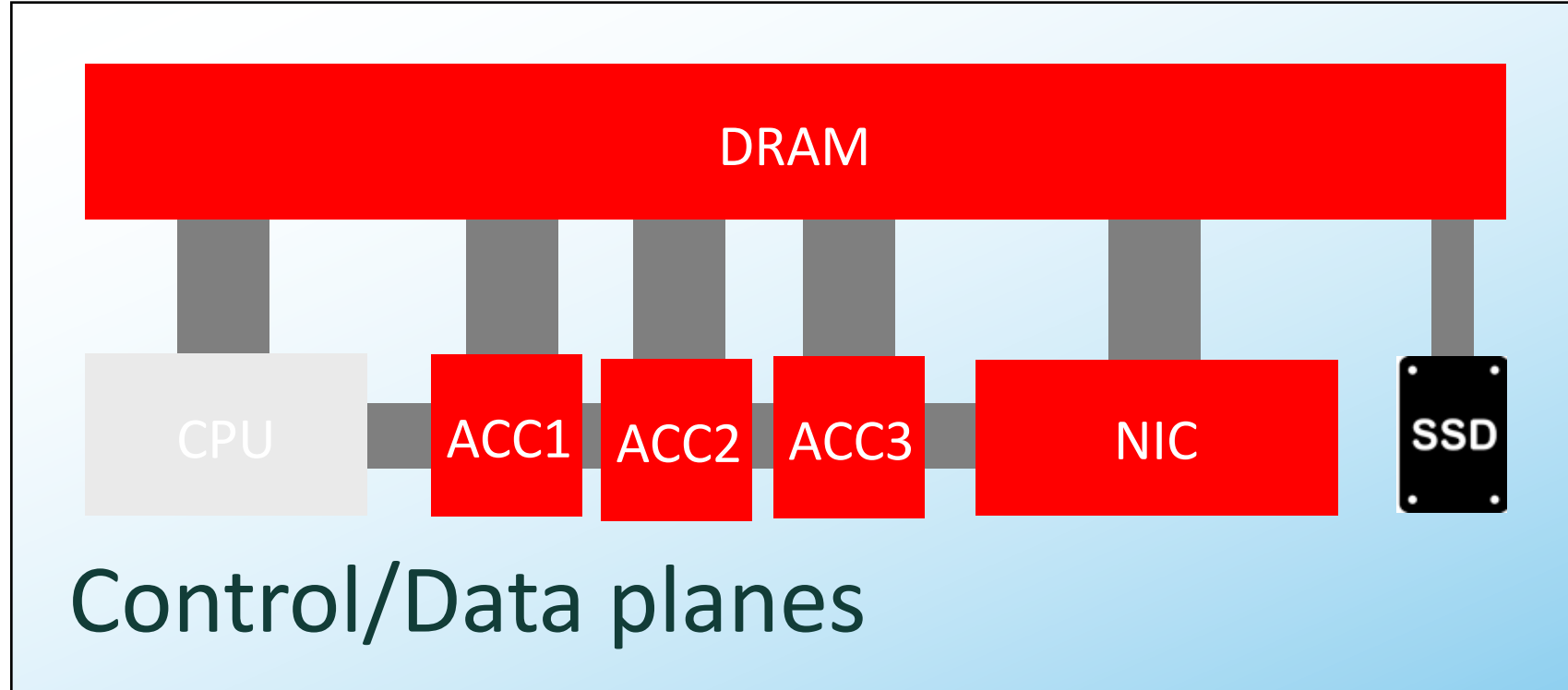
- Dual 2GHz CPU's, 50ns DRAM, Linux
- Bottlenecked by CPU, OS & legacy interfaces
- Fragmented silicon



~~Red Hat~~



# IDEAL POST-MOORE SERVER



- Think of the server as a network
- Control plane: set up via CPU & OS
- Data plane: protected access to memory
- Eliminates silicon fragmentation



# OUTLINE

- ~~Overview~~

- Post-Moore servers

  - ~~Today's servers~~

  - ISA opportunities

    - CPU/Memory/Storage/Network

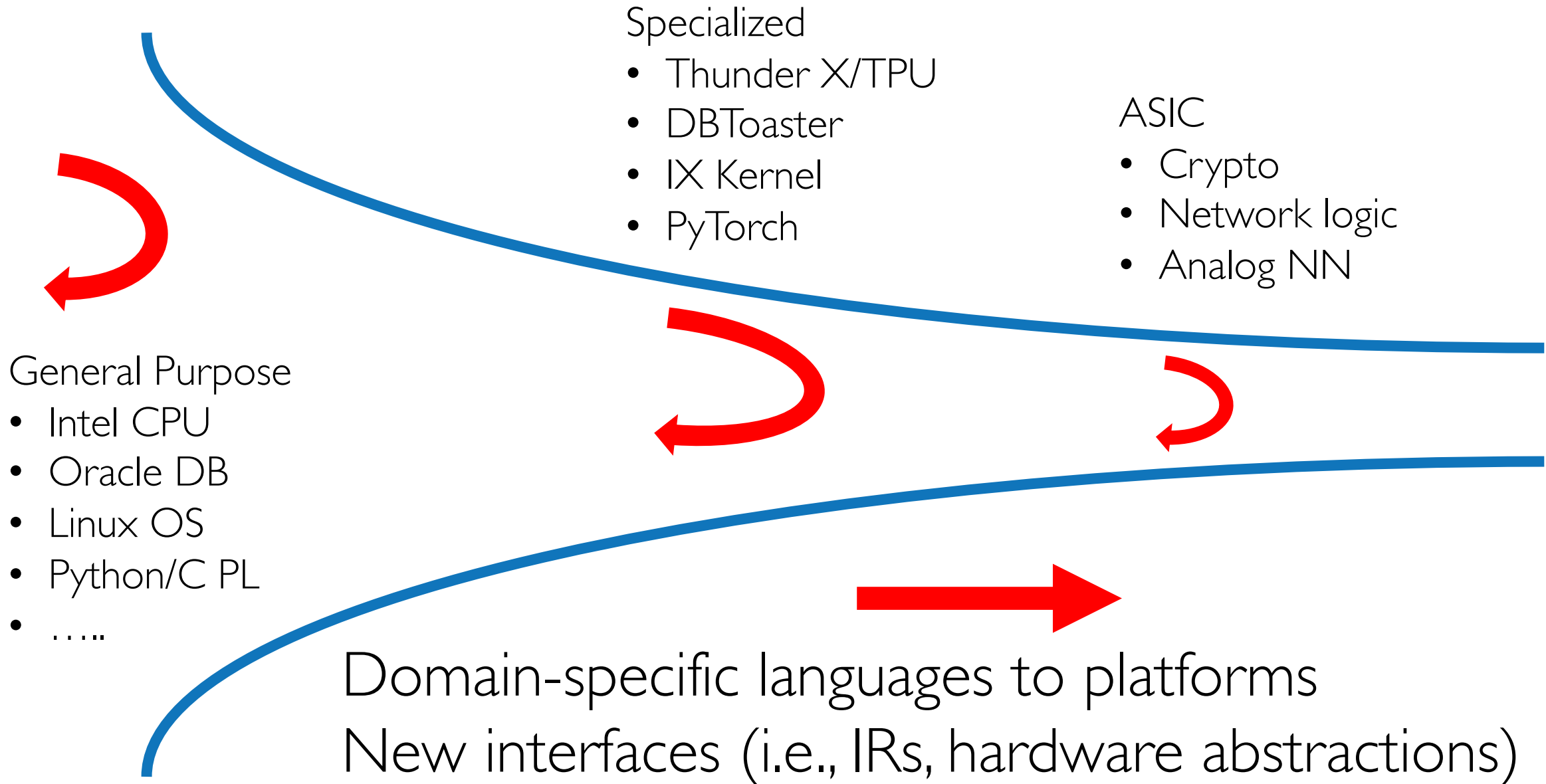
    - AI

- Datacenter sustainability

- Summary



# THE SPECIALIZATION FUNNEL



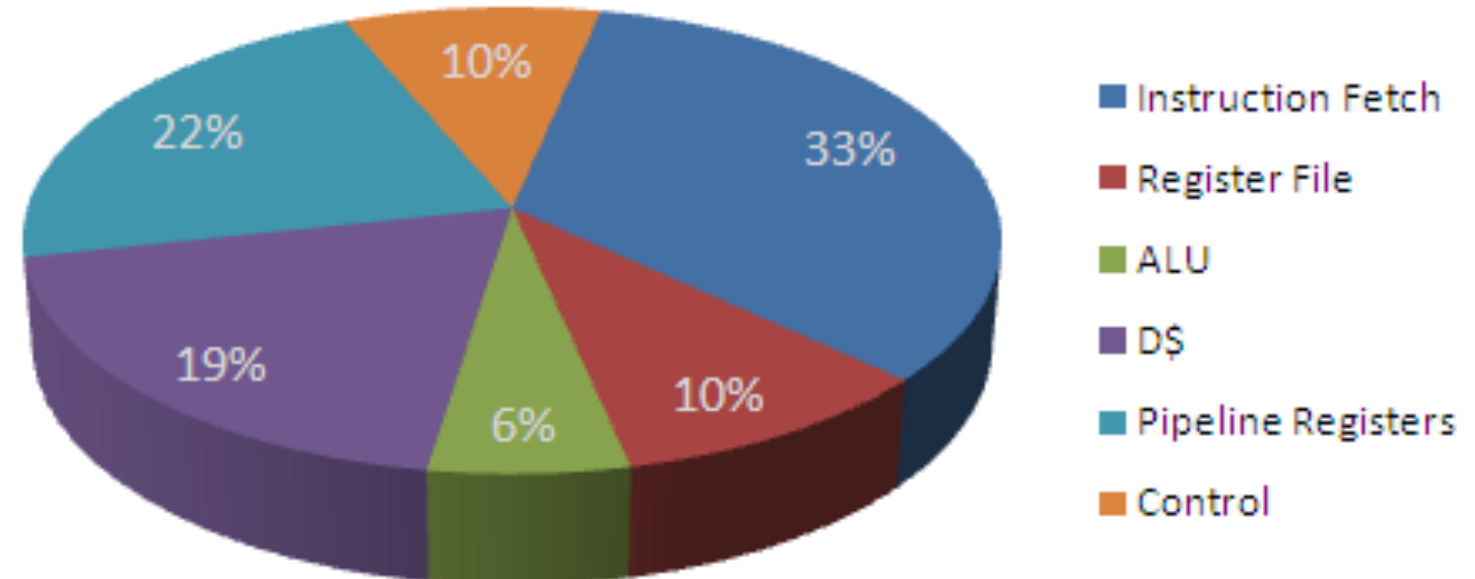


# THE LIMITS OF CPUS

CPUs follow the von Neumann machine organization

- Machine instructions fetched from memory
- Operands fetched/written to memory
- Referred to as von Neumann bottleneck

Only **6%** power in Pentium 4  
spent in arithmetic (ALU)



[src: Chen, et. al., IEEE Transactions, 2006]

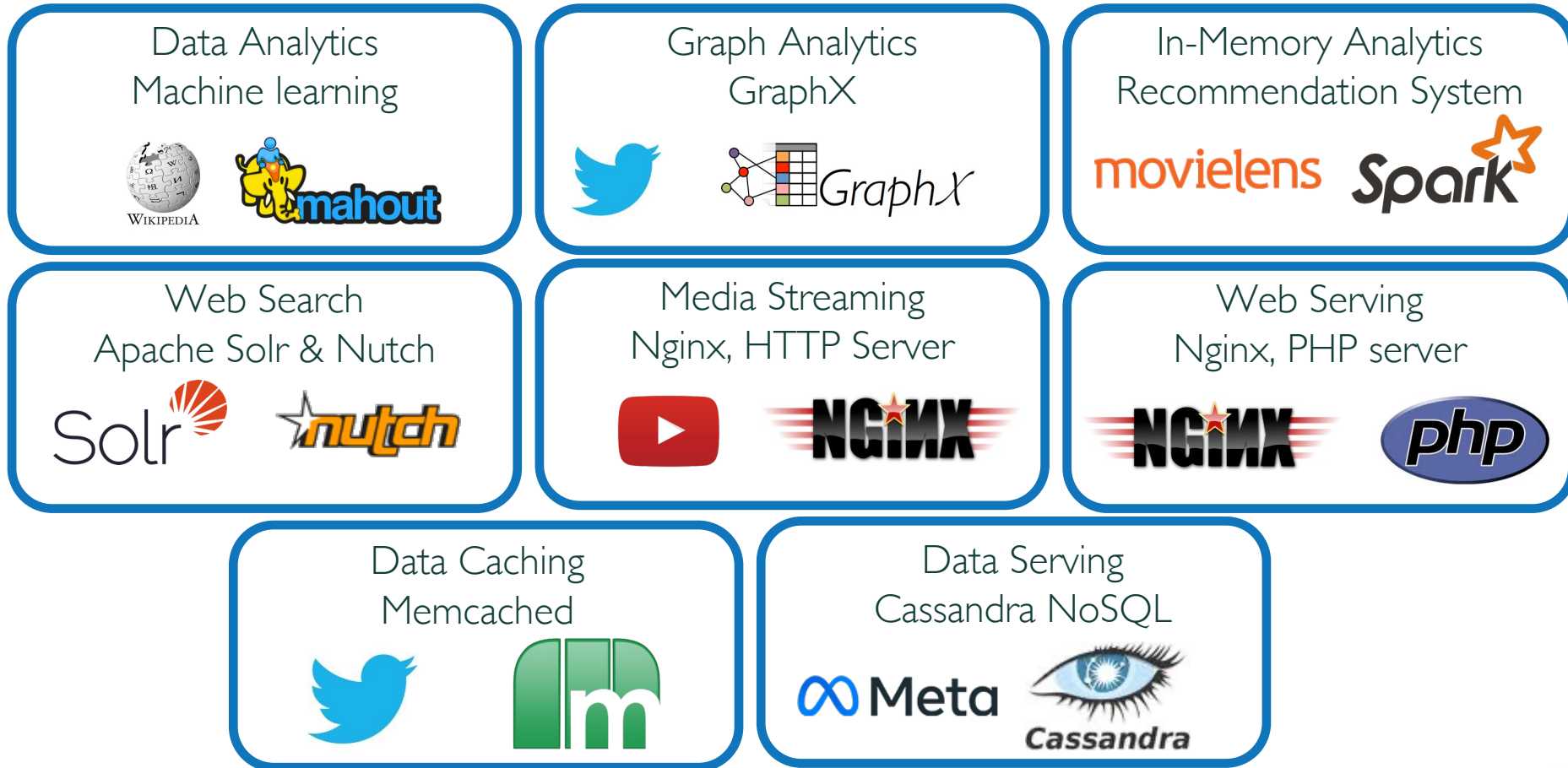


- First-party workloads (e.g., search, retail, media)
  - Data management
  - Analytics
  - Monoliths to microservices
- Third-party workloads (cloud)
  - Containerized
  - Serverless





# CloudSuite (4.0 release @ [cloudsuite.ch](http://cloudsuite.ch))

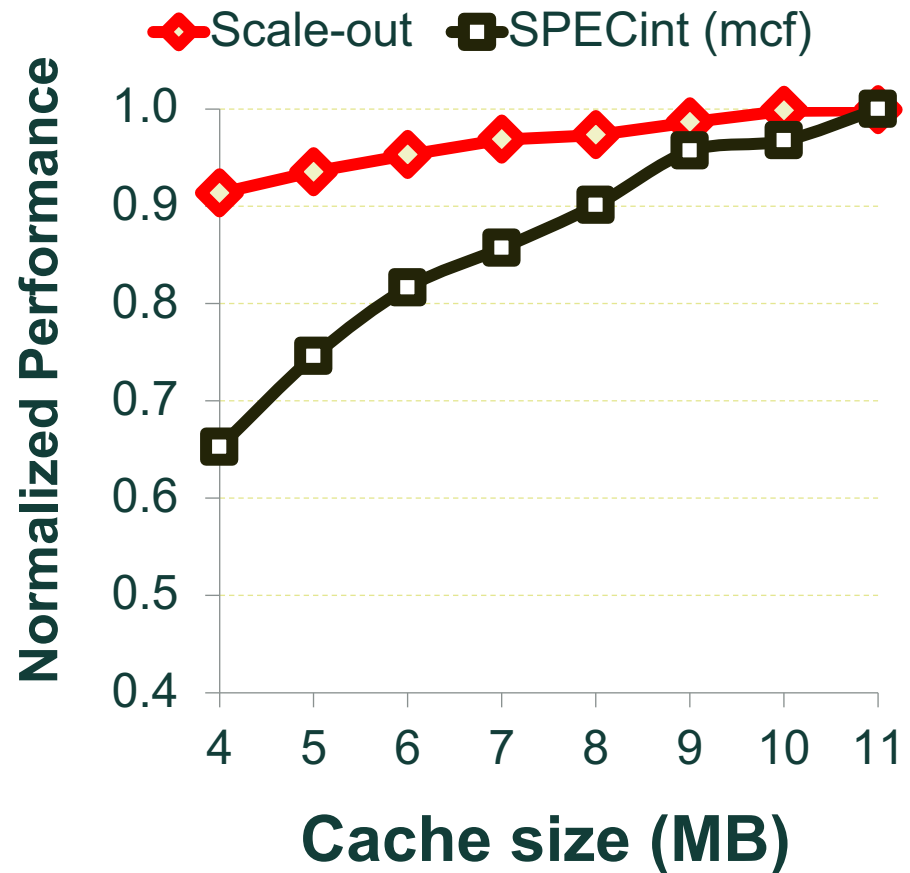


Supports x86, ARM64, RISC-V (coming)

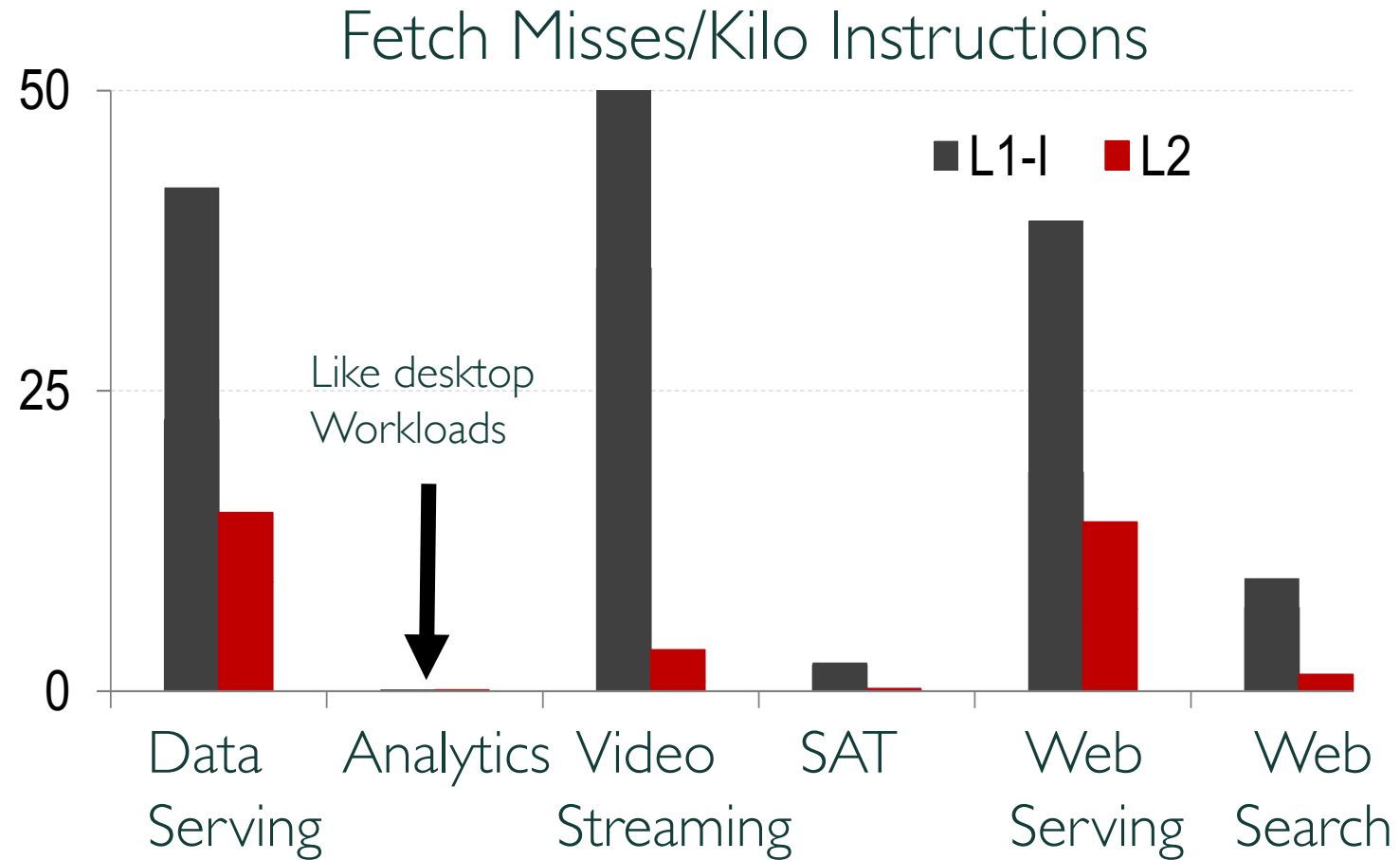




# SERVICES STUCK IN MEMORY [ASPLOS'12]



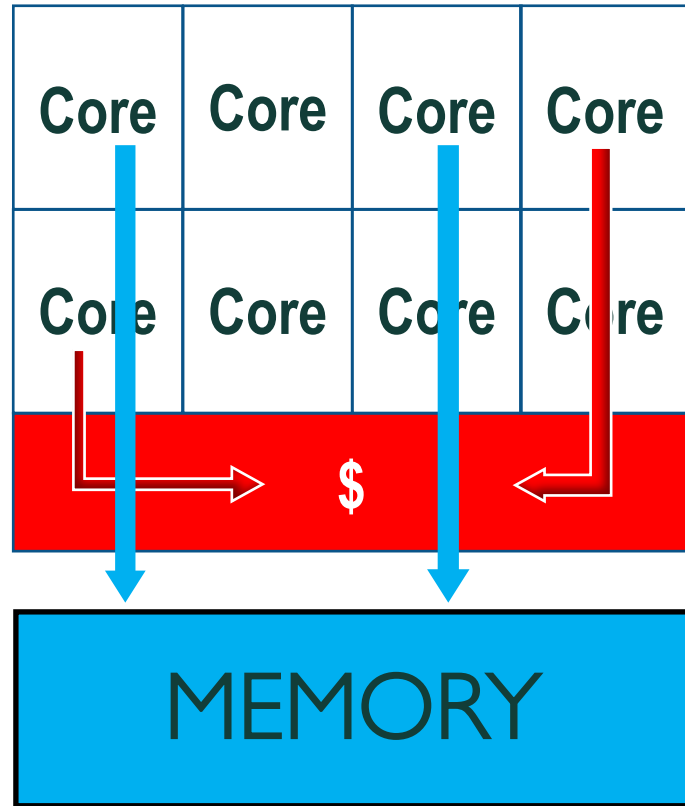
Cache overprovisioned



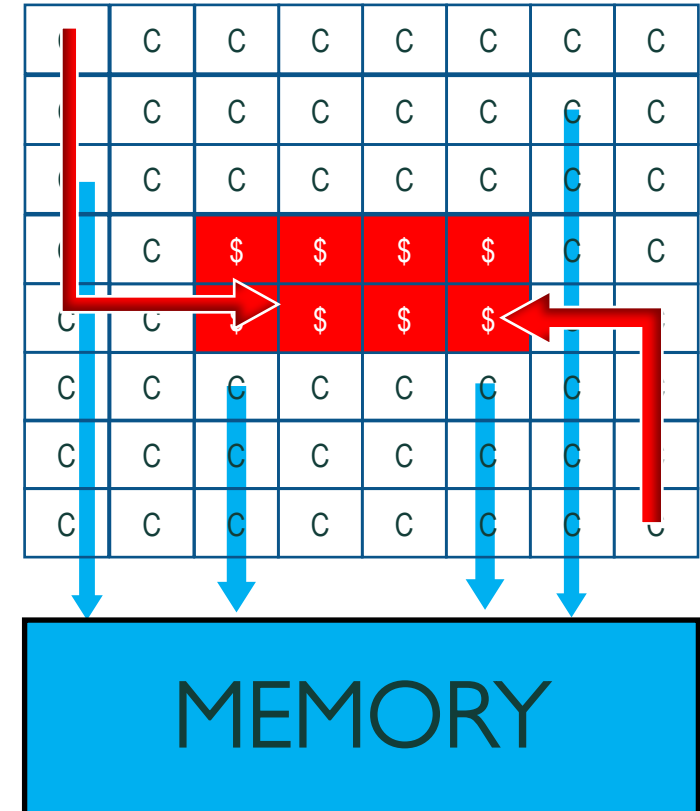
Instruction supply bottlenecked



# SCALE-OUT PROCESSOR (SOP)



- General-purpose CPU
- ✗ Logic 60% of silicon
- ✗ 6x bigger cores



- 3-way OoO ARM
- ✓ 85% logic, 7x more cores
- ✓ Faster instruction supply



# CLOUD-NATIVE CPU [c.a. 2014]



Case for Workload  
Optimized Processors  
For Next Generation  
Data Center & Cloud

**Gopal Hegde**

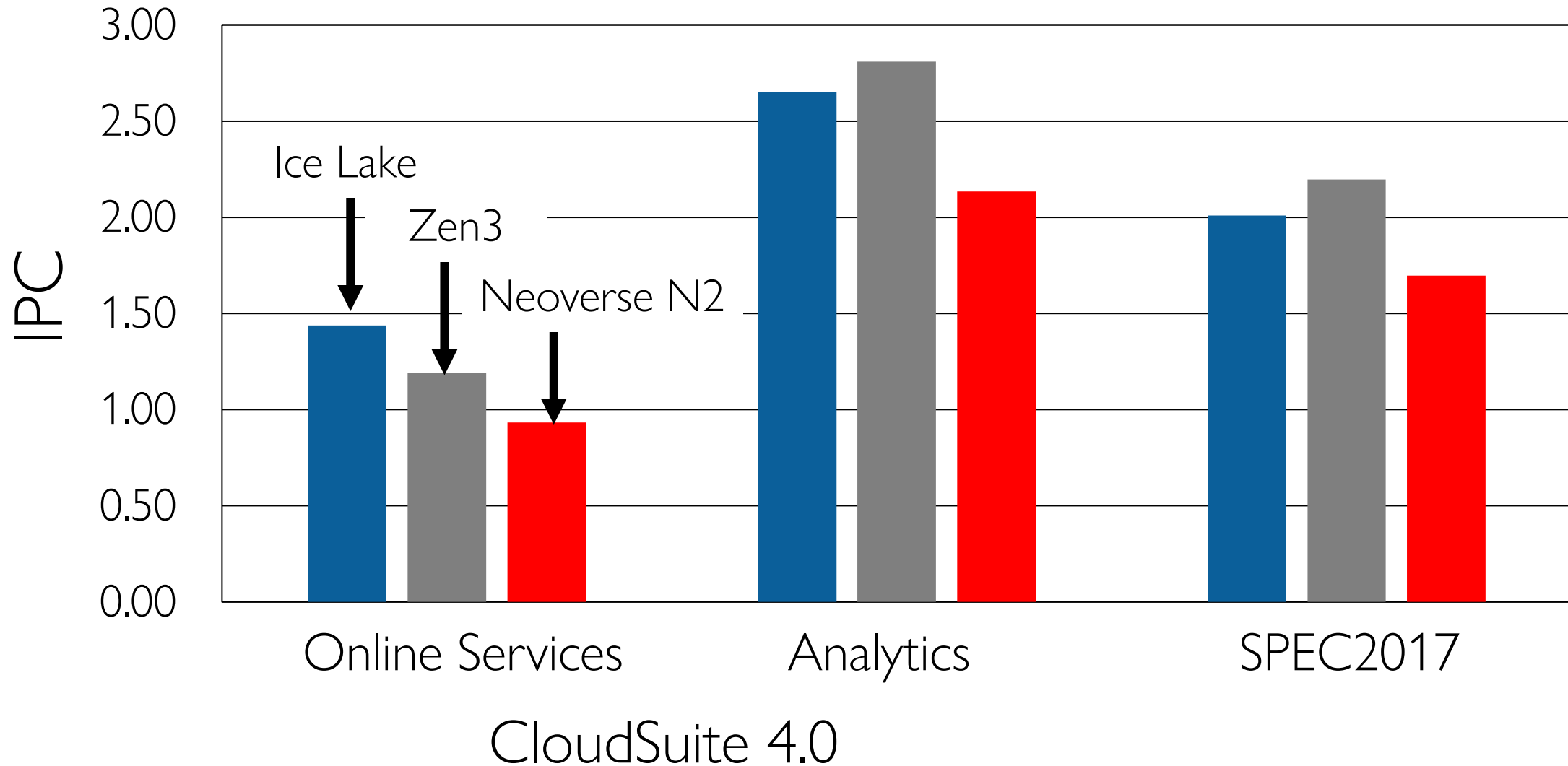
VP/GM, Data Center Processing Group

## Thunder X

- Based on SOP blueprint
- Designed to serve data
- 7x more core than cache
- Optimizes instruction supply
- Ran stock software
- 10x throughput over Xeon

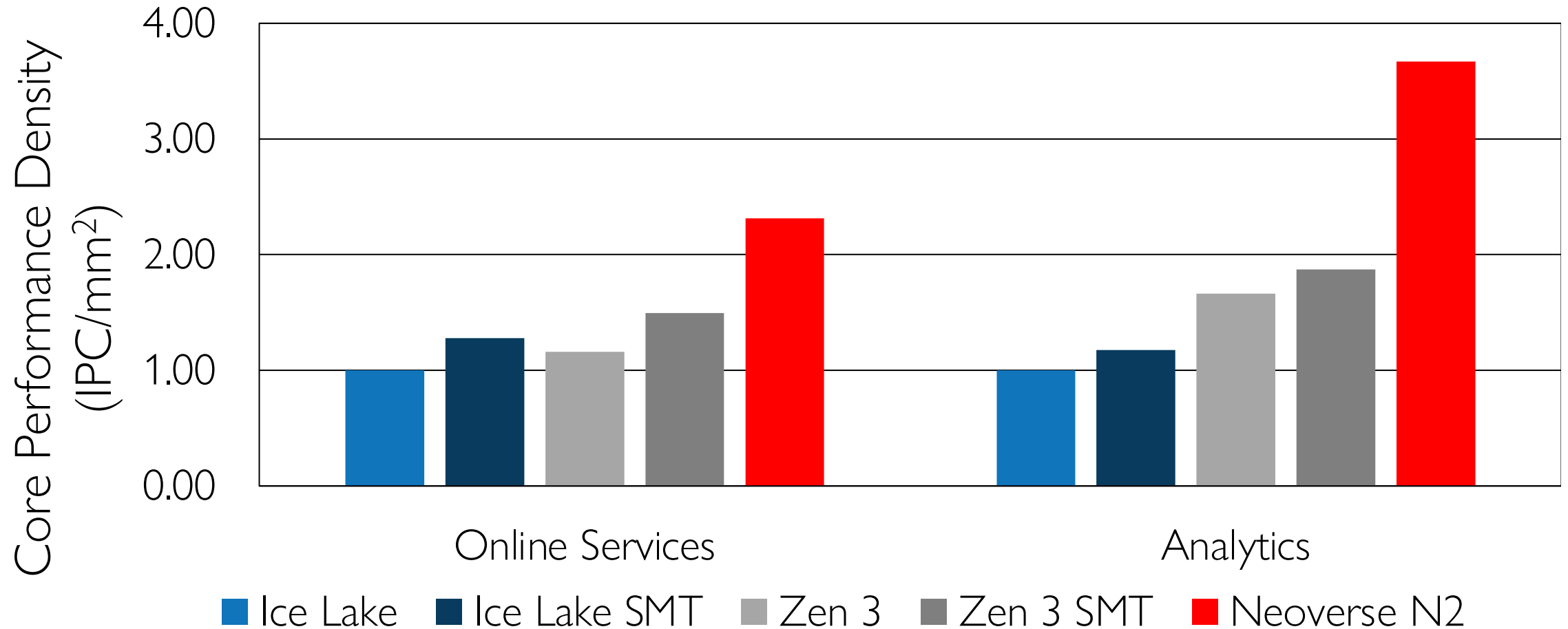


# X86 VS. ARM SINGLE THREAD





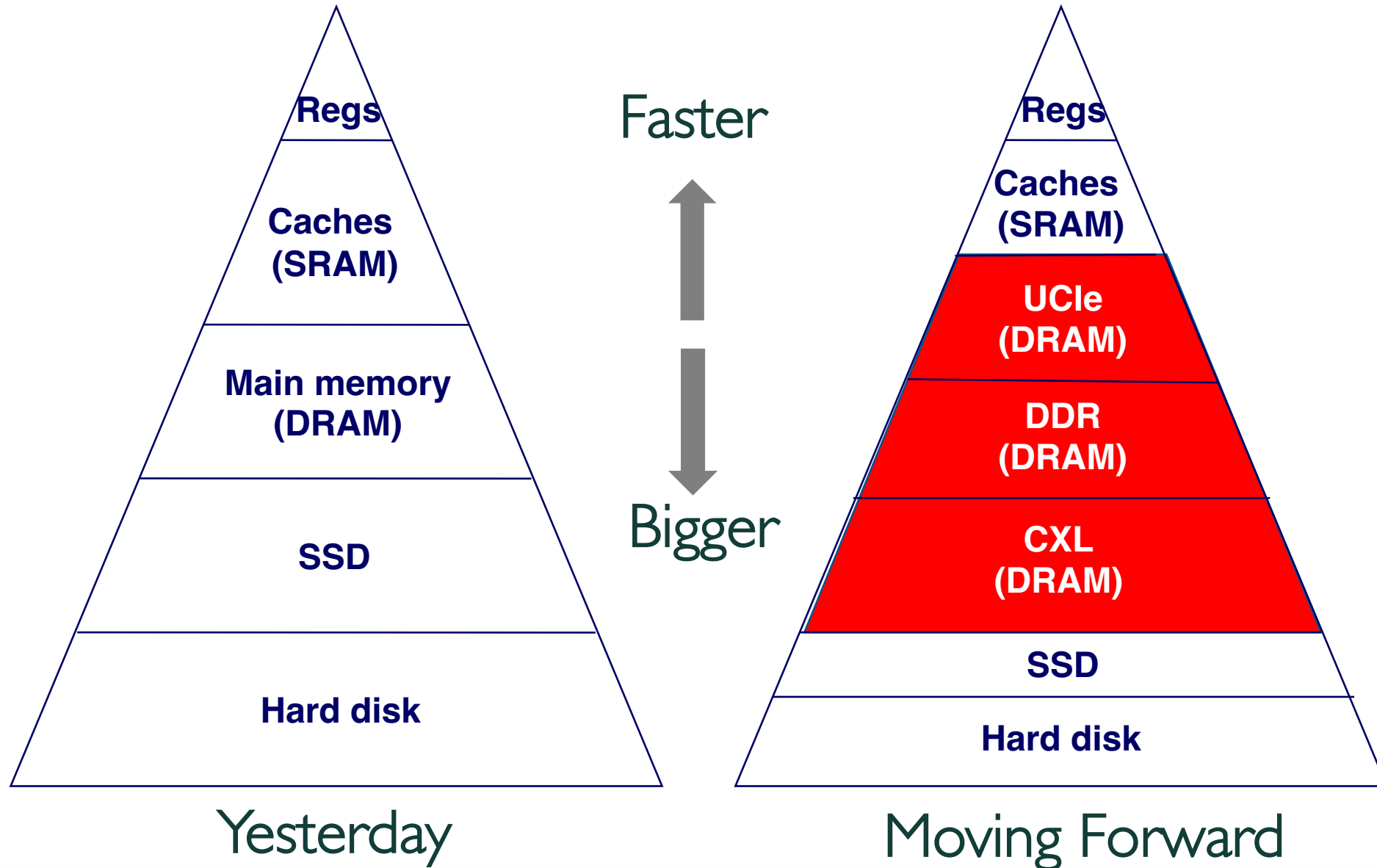
# SILICON EFFICIENCY IN X86 VS. ARM



The online services meet end-to-end tail latency requirements

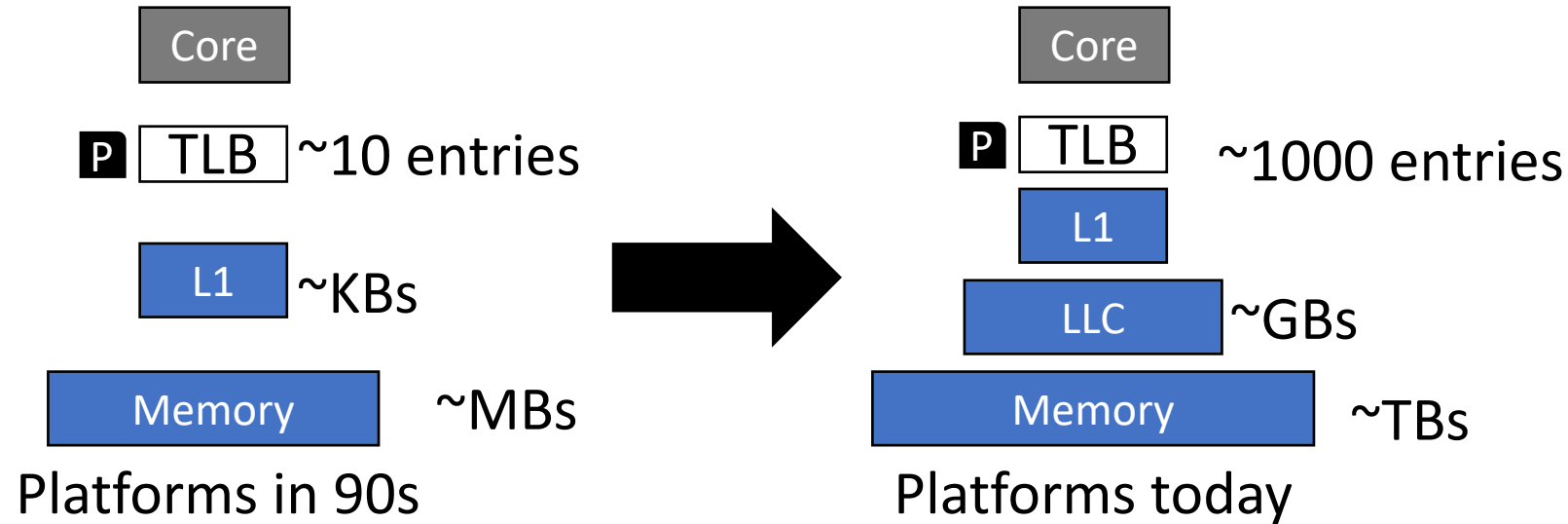


# TB-SCALE HIERARCHIES





# THE VM BOTTLENECK: TLBs



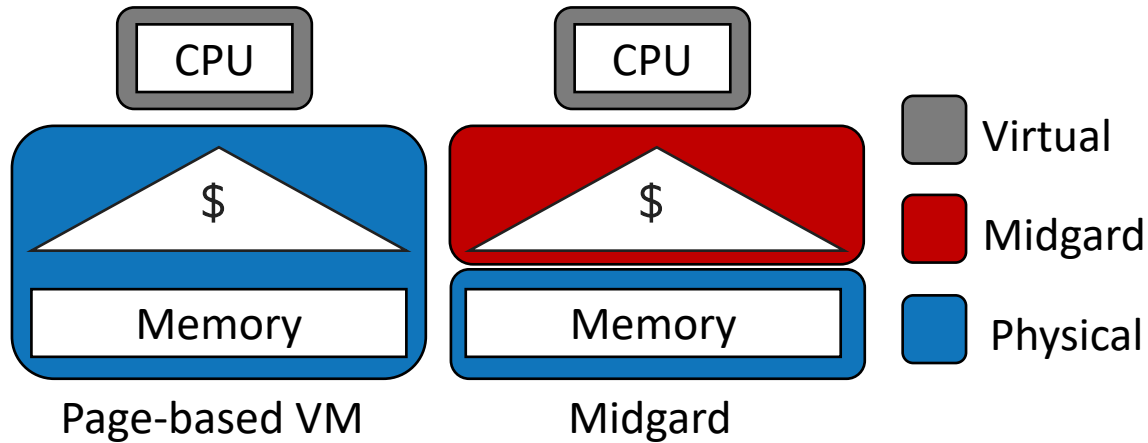
Product	Year	Cores	Cache capacity	TLB entries	Coverage (4KB)
Intel P4	2000	1	256KB SRAM	64	256KB
Intel KabyLake	2016	4	128MB eDRAM	1536	6MB
Apple M1	2020	8 (4+4)	16MB SRAM	3096	12MB (16KB)
AMD Zen3	2021	64 (8x8)	256MB SRAM	2048	8MB
Intel Sapphire Rapids	2022	56 (14x4)	64GB HBM2	?	?



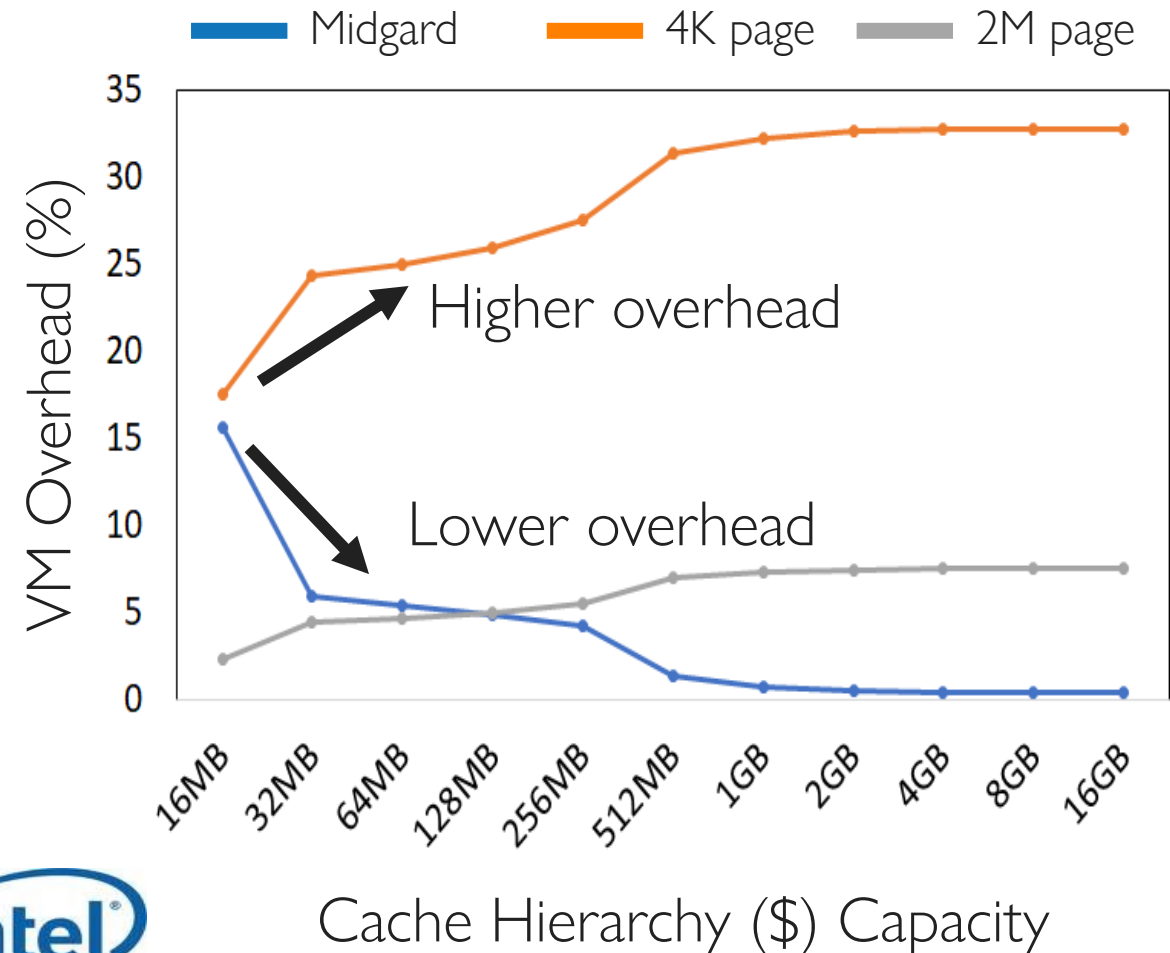
# TB-SCALE MEMORY WITHOUT TLB



[midgard.epfl.ch](https://midgard.epfl.ch)



- Keeps POSIX (VMA) interface to apps
  - Linux, MacOS/iOS, Android
- Eliminates page-based translation in \$
- ✓ Unclogs virtual memory for security, virtualization, accelerators





# TB-SCALE MEMORY WITHOUT TLB



## Midgard Roadmap:

CPU microarchitecture/OS [ISCA'21'23]

Compartmentalization [IEEE S&P'23]

Virtualization/Containerization

Accelerator ecosystem/IO

Monolith/ $\mu$ services/serverless

....

**EPFL**



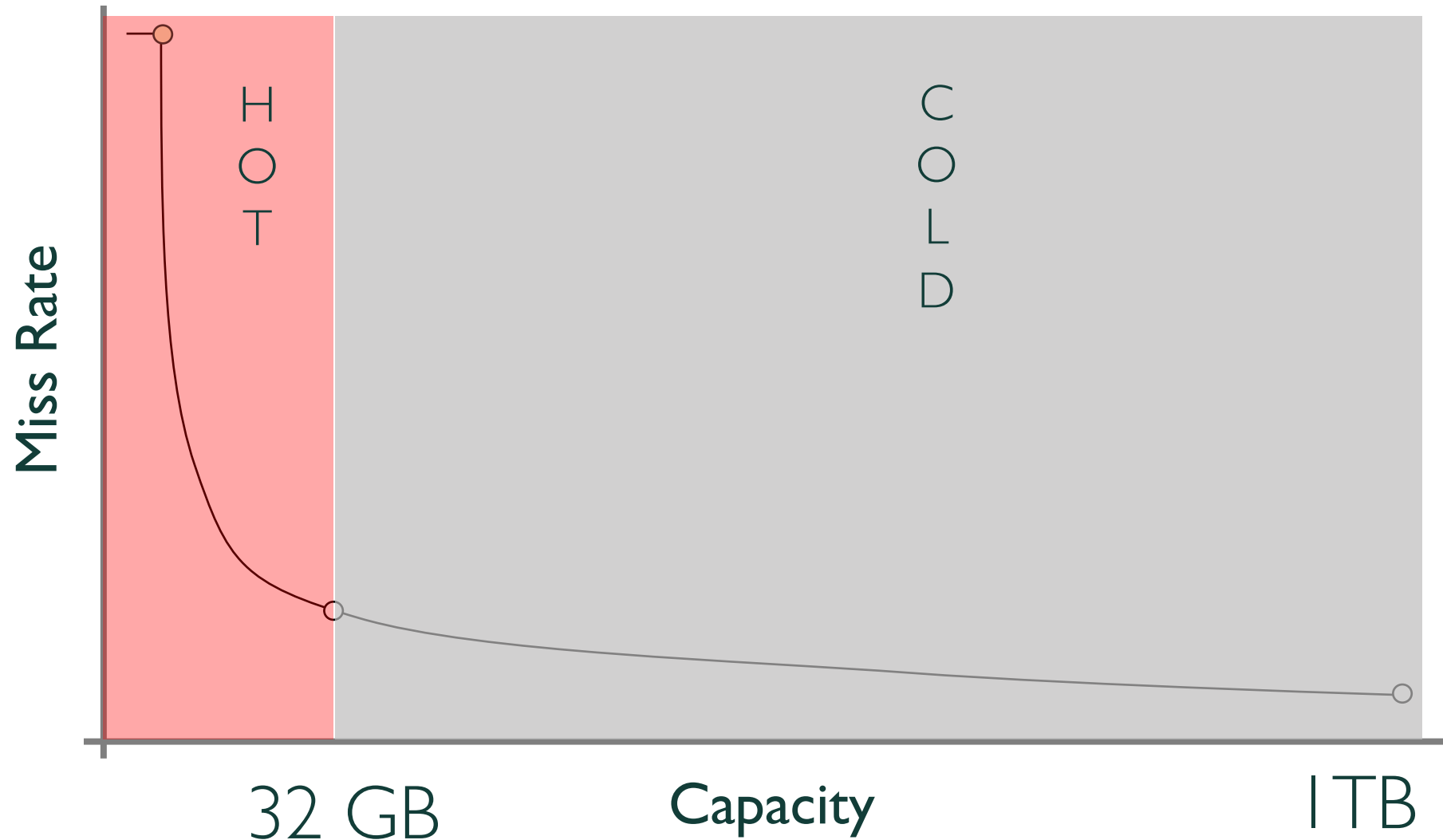
**Yale University**

Intel Transformative Server  
Architecture Center



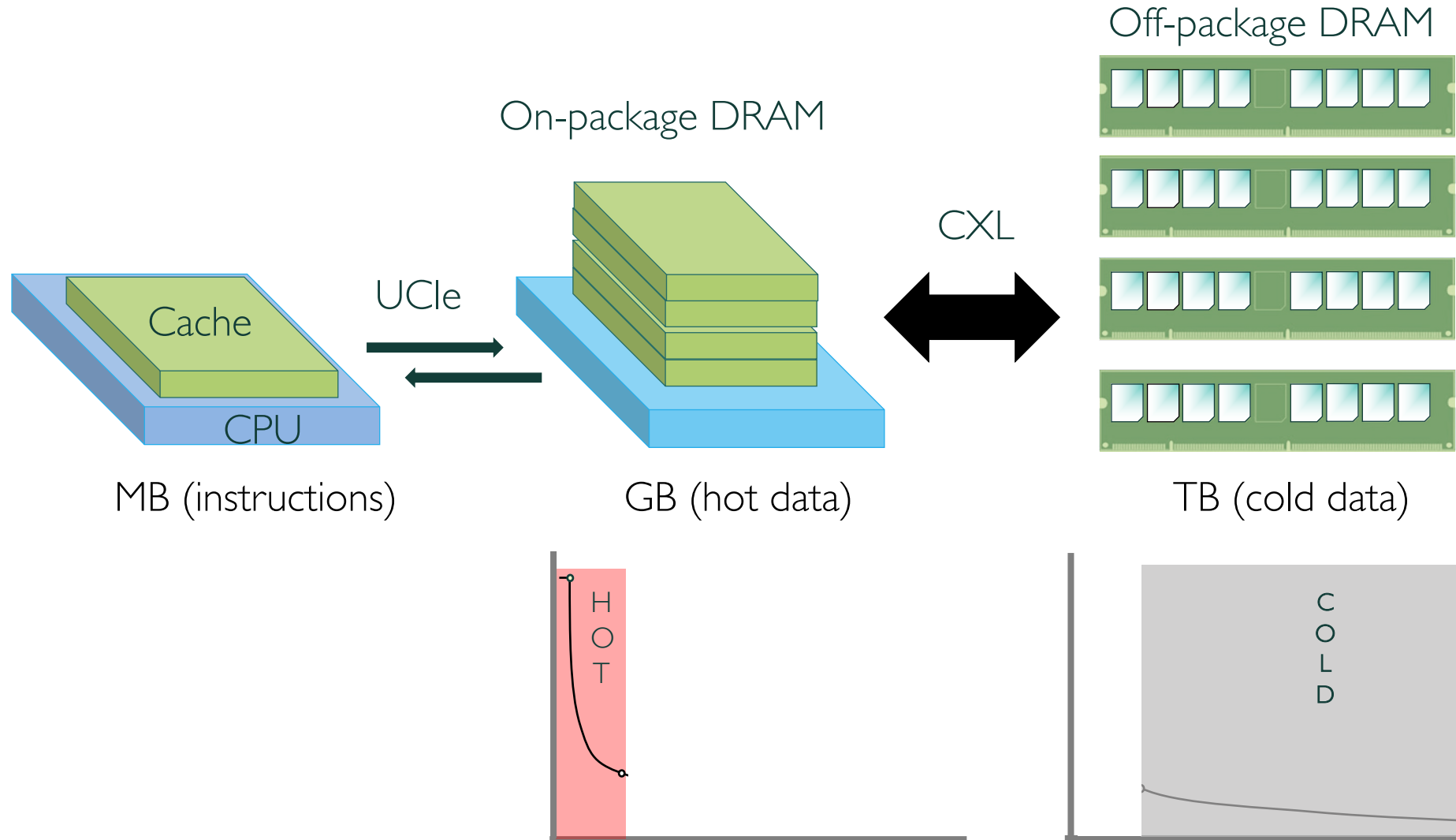


# CLOUD-NATIVE MEMORY HIERARCHIES



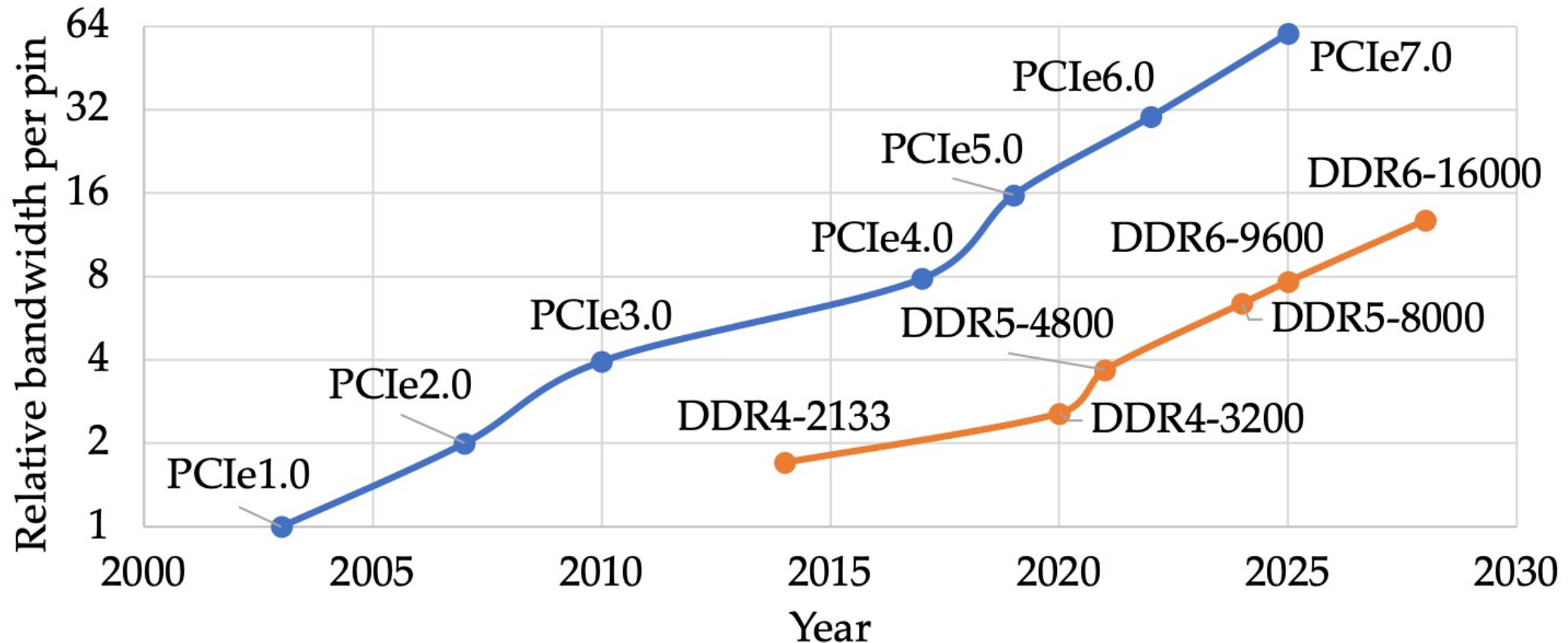


# EMERGING MEMORY HIERARCHIES





# PCIe BANDWIDTH PROJECTIONS



[A Case for CXL-Centric Servers, Cho, et. al.]

- CXL will eventually subsume DDR



- Memory is currently wasted (50% in containers)
- Much memory is stranded (Pond, ASPLOS'23)
- Pool memory in proximity (e.g., Scale-Out NUMA)
  - Both requests and data partitions are skewed
  - Helps with load balancing
- Memory accounts for much of fabrication emissions
  - Use CXL to keep (old not new) DDR memory



# ON-PACKAGE OPPORTUNITIES: ACCELERATORS

## Data management

- Data copy/move
- Compression/decompression
- Serialization/deserialization
- Encrypt/decrypt
- Scatter/gather

## Analytics

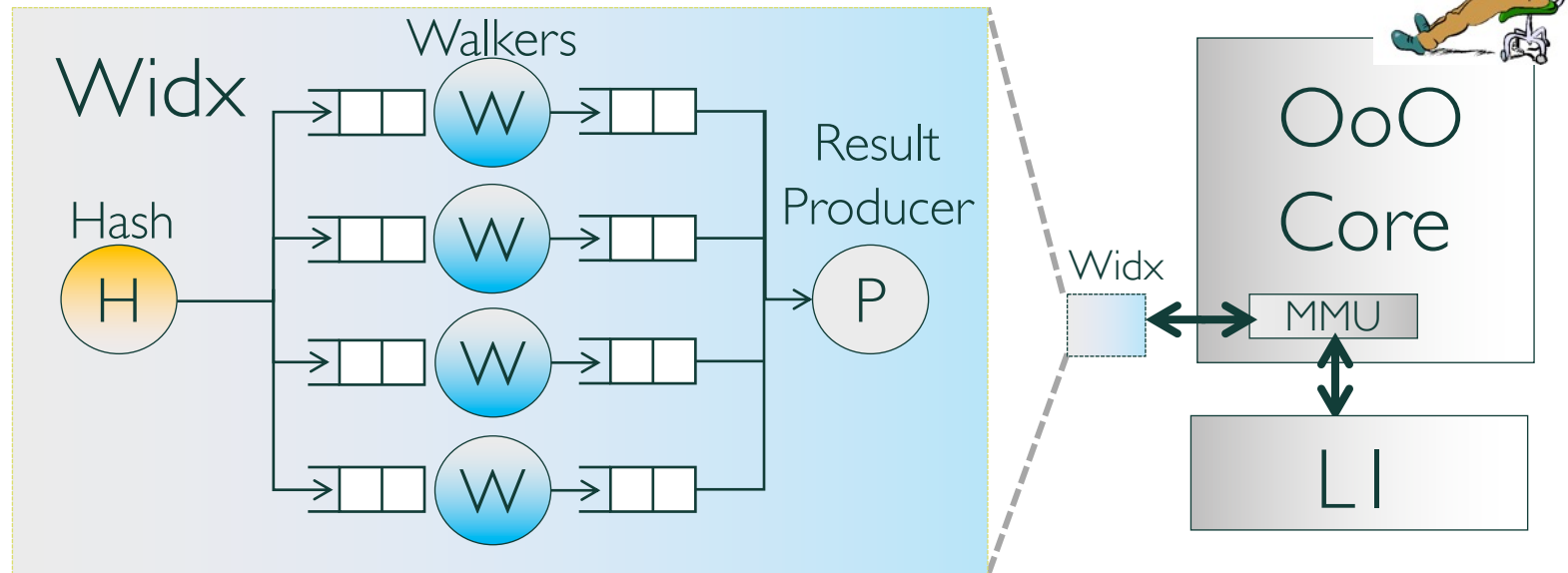
- Spark operations  
(e.g., scan, filter, groupby)
- SQL operations  
(e.g., intersect, union, join)
- Matrix operations

Memory views: rows for data management , columns for analytics



# CHASING POINTERS W/ WALKERS

- Traverse data structures (e.g., hash table, B-tree)
- Parallelize pointer chains
- Overlap pointer access across chains



**15x better performance/Watt over Xeon**



# WALKERS IN SOFTWARE [VLDB'16]

Use insights to help CPUs

- Decouple hash & walk(s) in software
- Schedule off-chip pointer access with co-routines

2.3x speedup on Xeon

- Unclogs dependences in microarchitecture
- Maximizes memory level parallelism
- DSL w/ co-routines
- Integrated in SAP HANA [VLDB'18]



# NEAR-MEMORY PROCESSING [ISCA'17]

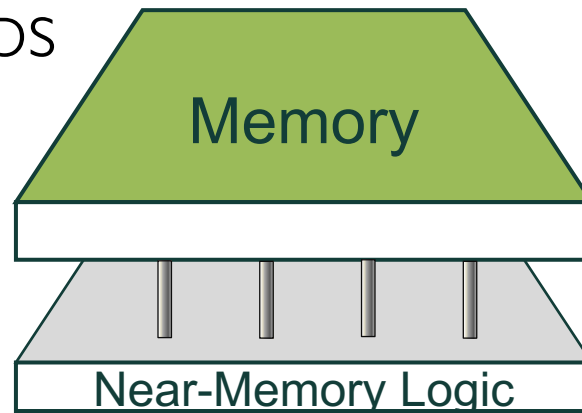
Stream data out of row buffers

SIMD cores + data streaming

- Saturates b/w with parallel SIMD streams
- 1024-bit SIMD @ 1 GHz
- No caches

Runs Spark Analytic Ops

50x over Xeon



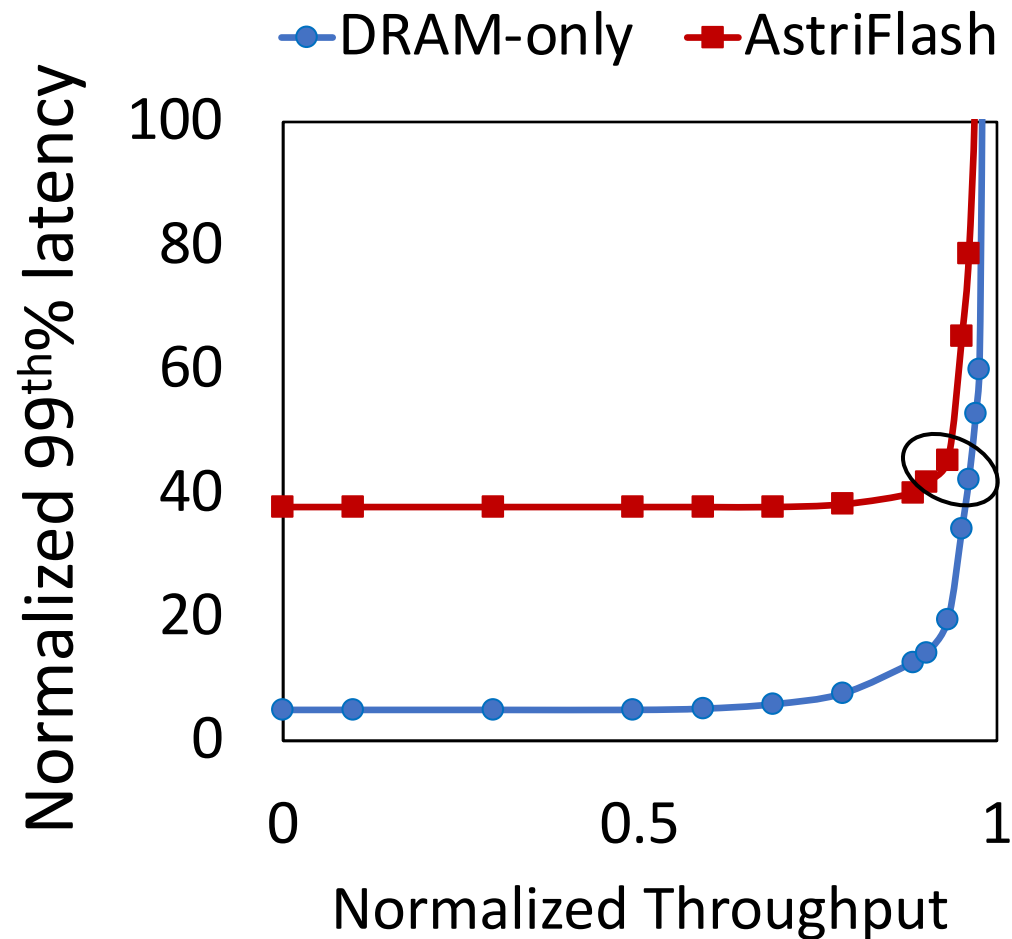
Algorithm/hardware co-design maximize near-memory performance



# ONLINE SERVICES IN FLASH [HPCA'23]

	Cost	Latency
DRAM	1x	~100 ns
SCM	1/5x	1-10 $\mu$ s
SSD	1/30x-1/50x	> 50 $\mu$ s (OS)

- Host & serve data from SSD
- Map SSD as memory, DRAM as cache
- Co-design CPU & OS for
  - paging
  - $\mu$ s-level thread switch



Maintains tail latency with only 5% lower throughput

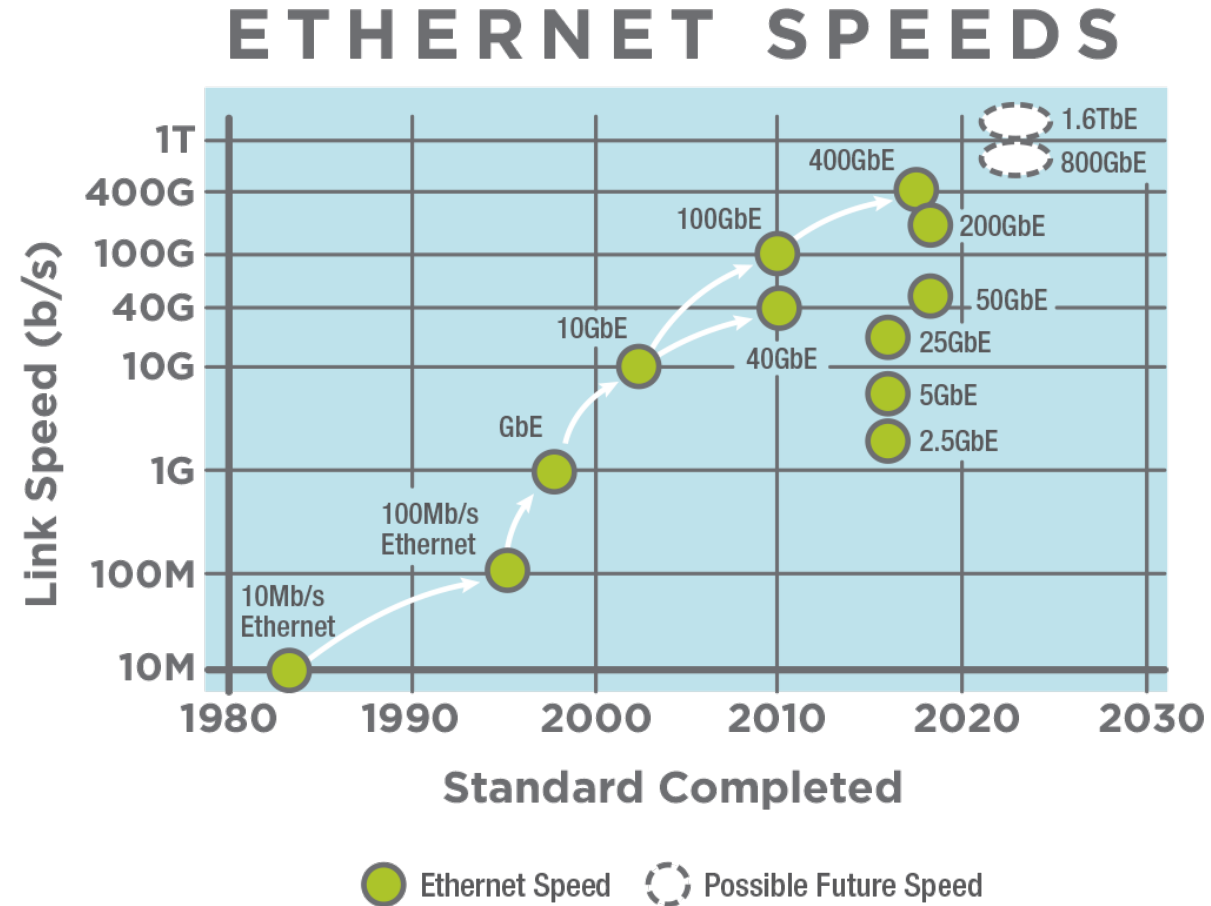


Network stack bottleneck:

- B/W growing faster than silicon
- Emerging  $\mu$ Services + serverless
- RPC, orchestration, ....

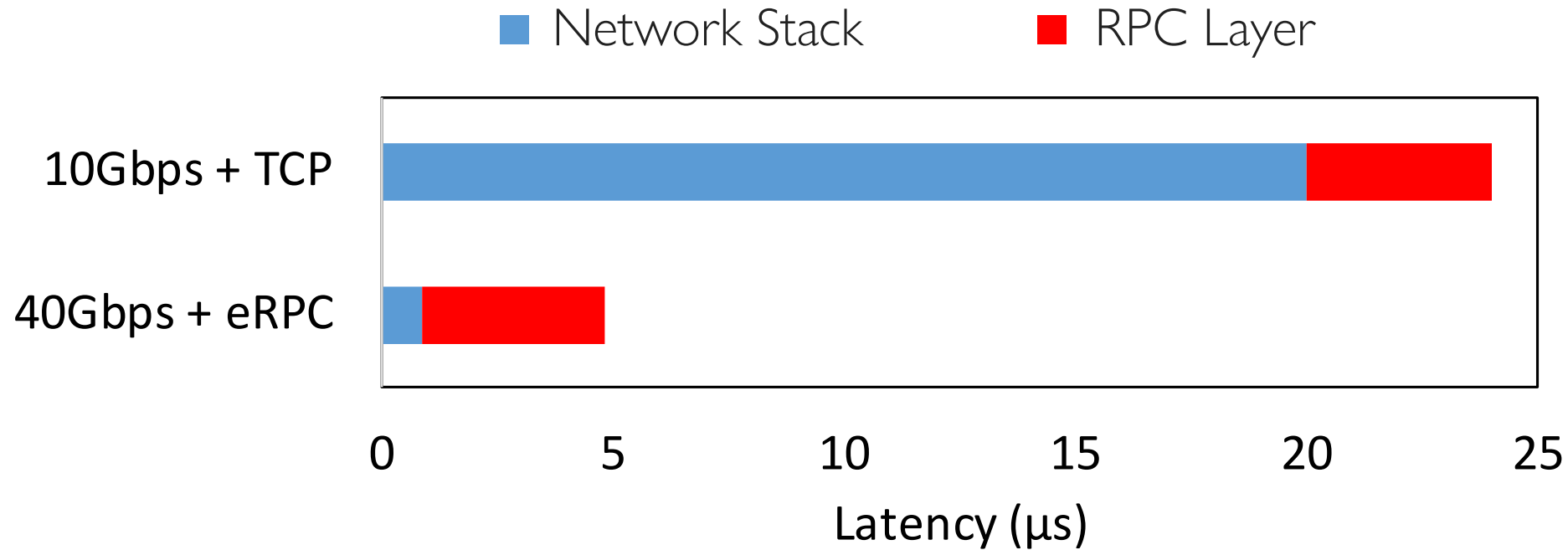
Key challenges:

- New abstractions
- Co-design of network stacks





# RPC ACCELERATORS



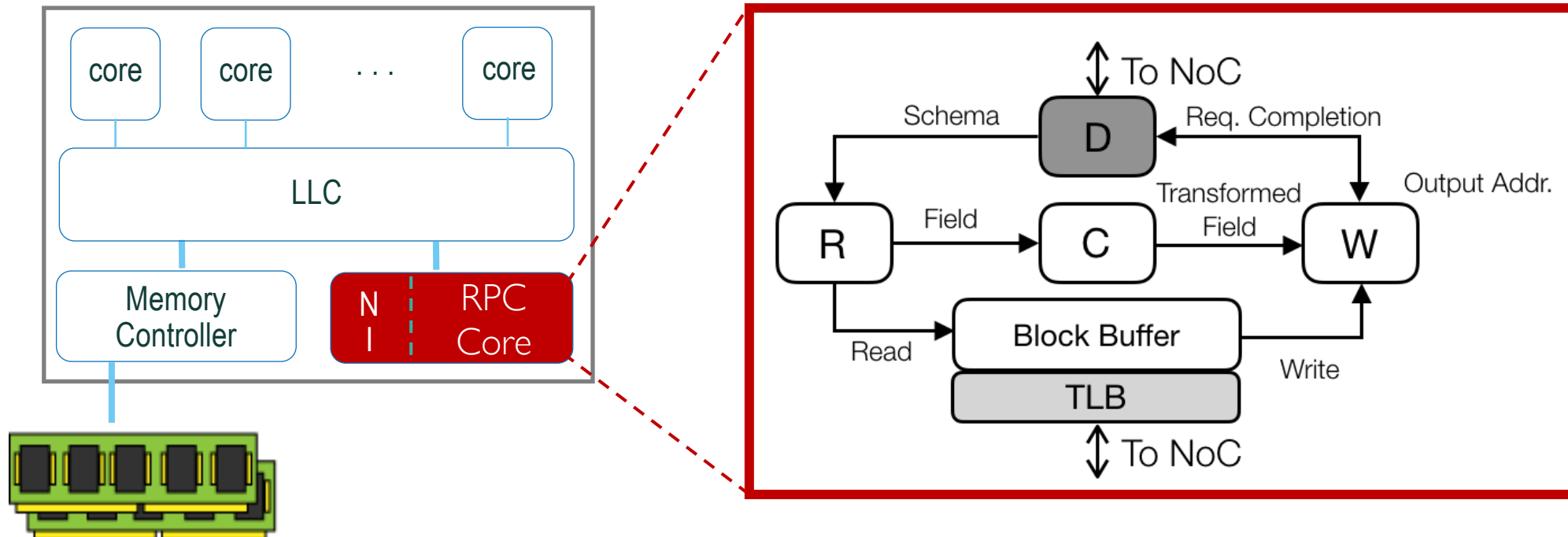
- Wire time and protocol stacks have shrunk
- RPC dominates CPU cycles in  $\mu\text{Services}$
- E.g., data transformation @  $\sim 2.4\text{Gbps}$  w/ Thrift on Xeon



# CEREBROS & NEBULA [ASPLOS'20, ISCA'20, MICRO'21]

RPC processing at line rate:

- A "schema" (not instructions) interface to an RPC core
- Implements load balancing/affinity scheduling for  $\mu$ Services



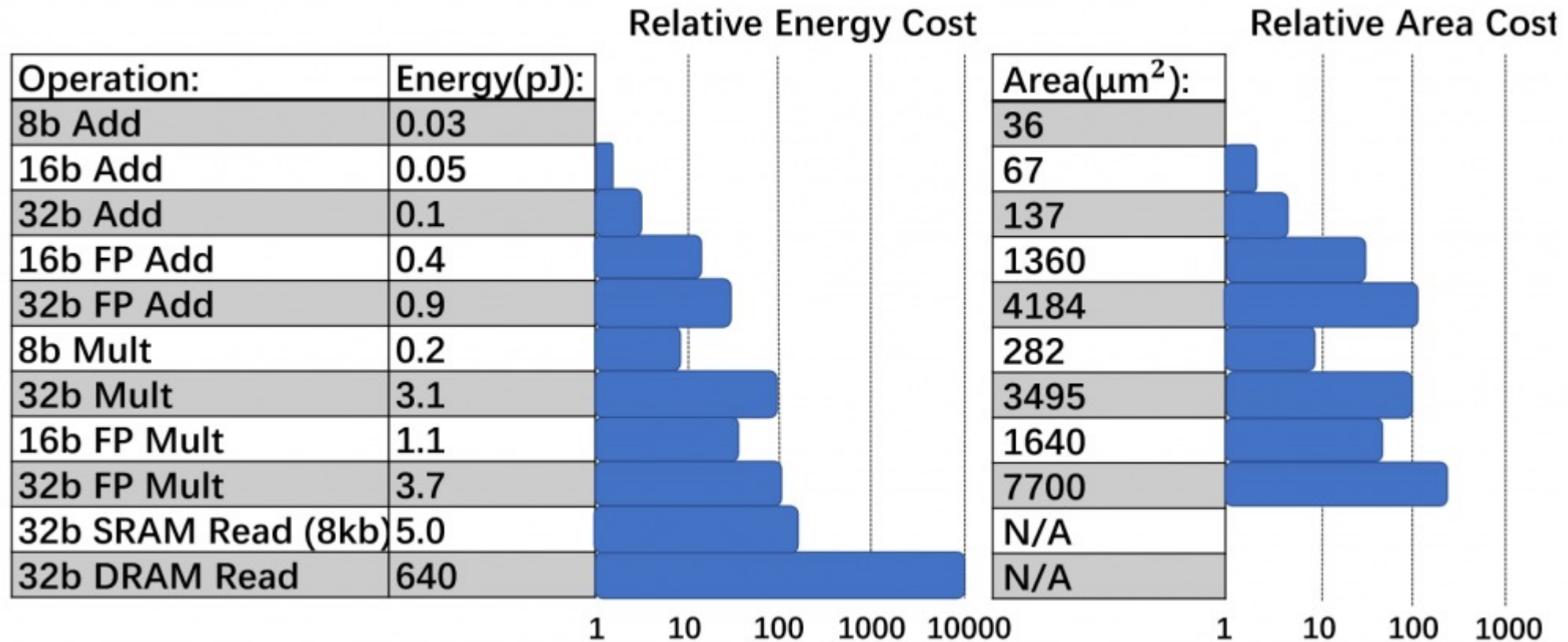


# OUTLINE

- ~~Overview~~
- Post-Moore servers
  - ~~Today's servers~~
  - ISA opportunities
    - ~~CPU/Memory/Storage/Network~~
    - AI
- Datacenter sustainability
- Summary



# COST OF LOGIC VS. MEMORY



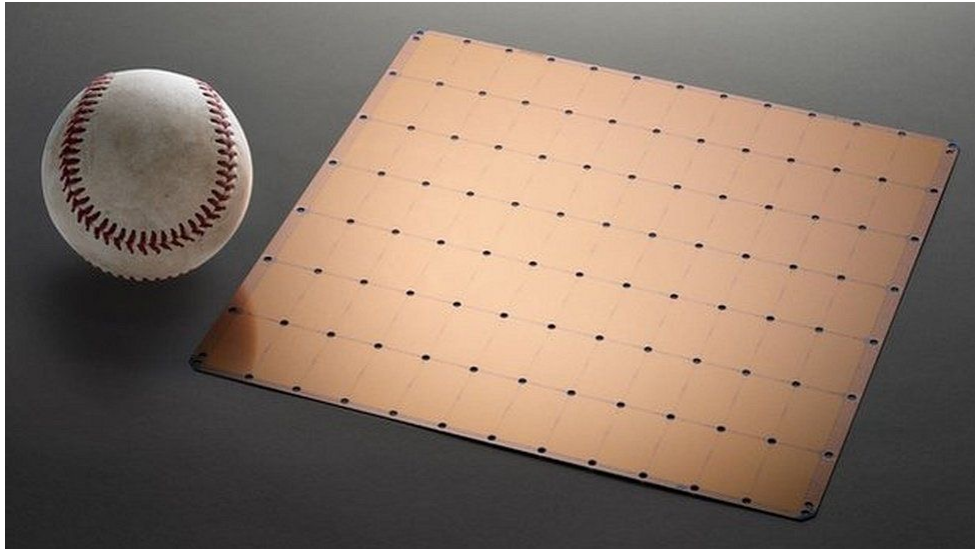
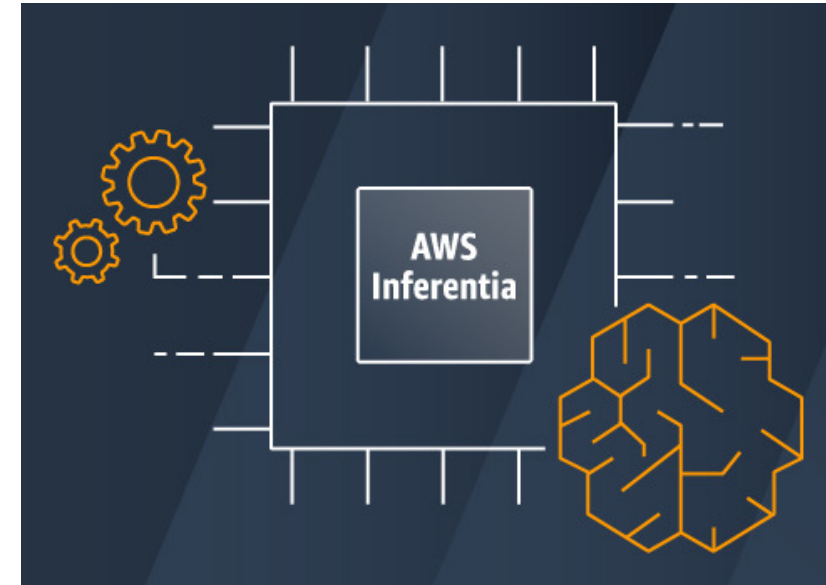
[src: Gholami, et. al.]



# DNN PLATFORM DIVERGENCE

Inference platforms:

- Tight latency constraints
- Ubiquitous deployment
- Relies on fixed-point arithmetic



Training platforms:

- Throughput optimized
- Server deployment
- Requires floating-point arithmetic

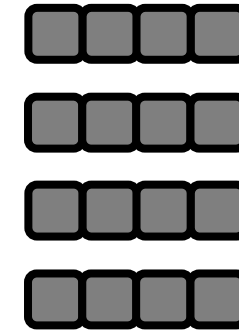


# HYBRID BLOCK FLOATING POINT (HBFP)

## 1. Block floating point (BFP): one exponent/tensor

- Low magnitude variation in tensor products
- > 90% of all arithmetic operations

Block of Mantissas



Exponent

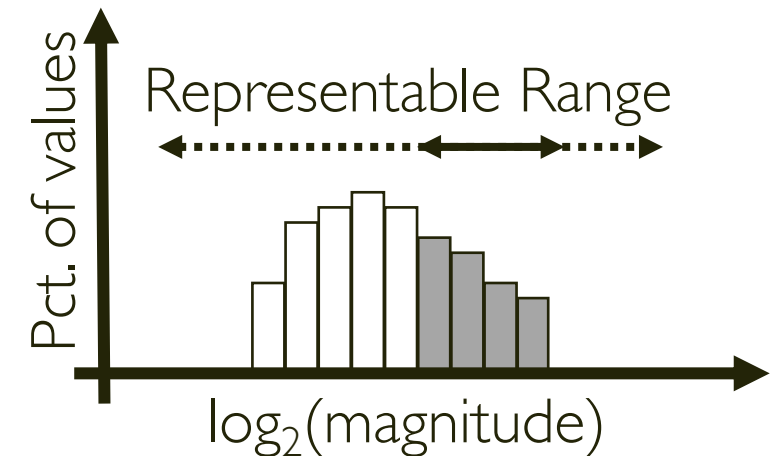


## 2. FP32 for all activations

- High magnitude variation in gradient updates

## Co-Located Training & Inference (ColTrain)

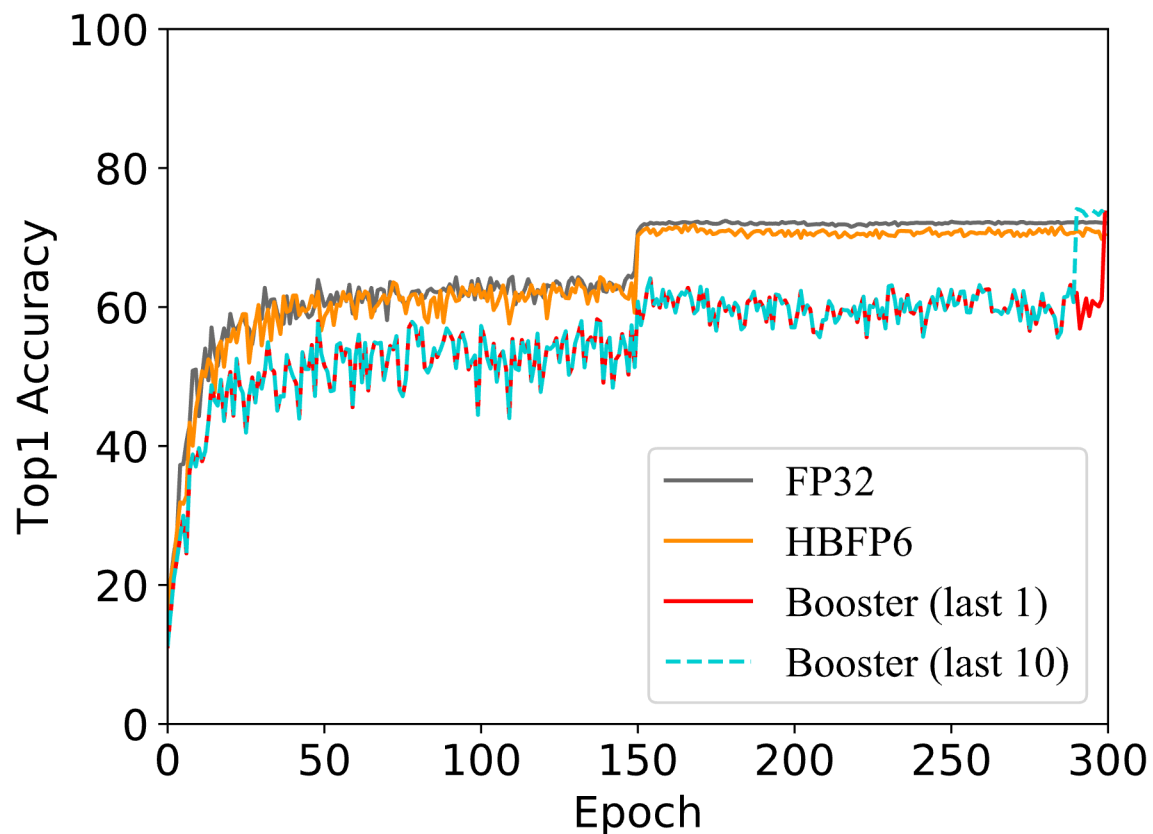
- ✓ One accelerator for training and inference
- ✓ Eliminates quantization
- ✓ Enables online learning





# MIXED-MANTISSA HBFP VS. FP32

DenseNet40 on CIFAR100



Transformers

Configuration	BLEU Score
FP32	34.77
HBFP6	34.47
HBFP4	32.64
<b>Booster</b>	<b>36.08</b>

FP32 level accuracy while using HBFP4 for majority of operations with 21.3x higher density



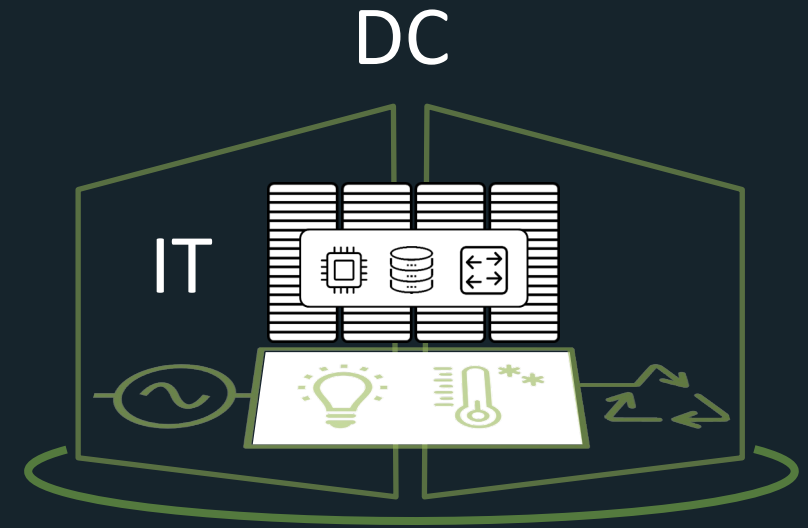
# OUTLINE

- ~~Overview~~
- ~~Post Moore servers~~
- Datacenter sustainability
- Summary



Today's efficiency metric  
power usage efficiency

$$\text{PUE} = \frac{\text{Total DC Power}}{\text{IT Power}}$$



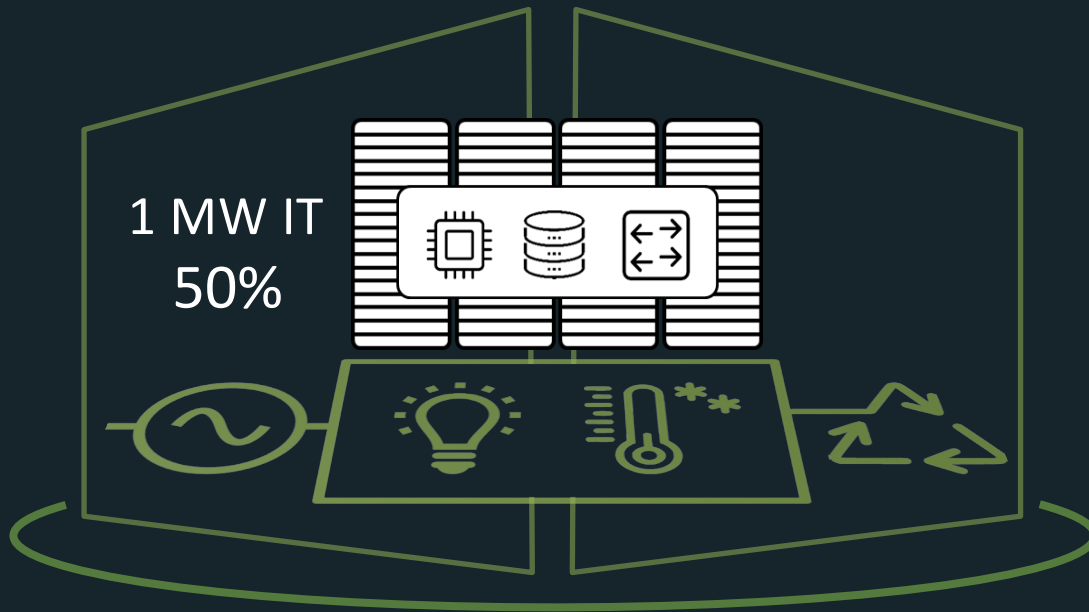
PUE has been around for two decades



# What is wrong with PUE?

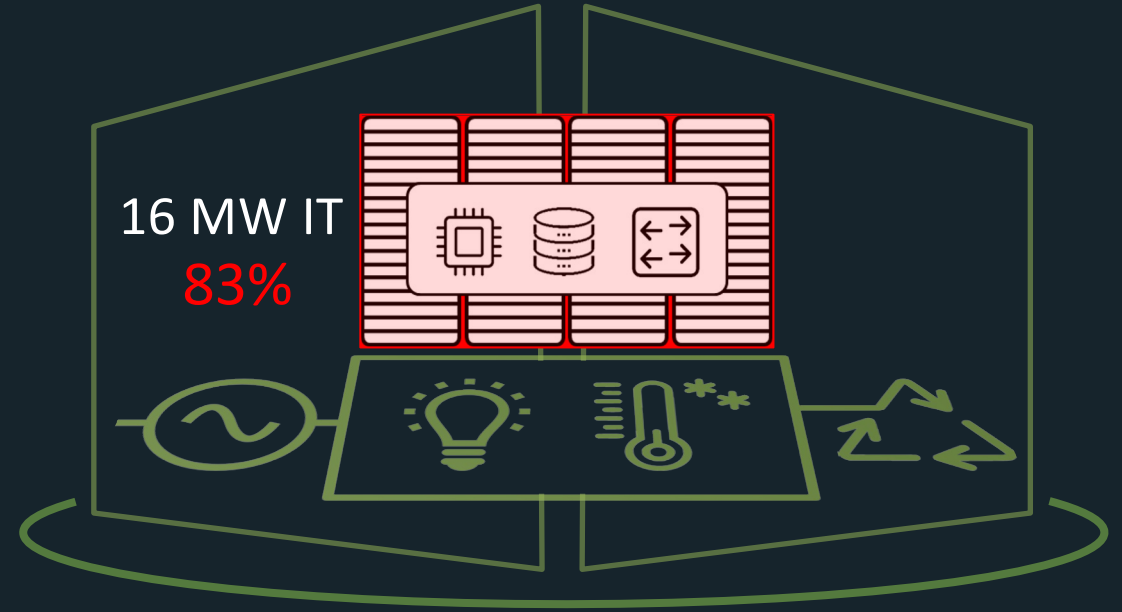
DC electricity goes to IT

2 MW DC in 2010



PUE = 2.0

20 MW DC in 2020



PUE = 1.2



# GOODBYE PUE! HELLO FULL-STACK EFFICIENCY

## DC INFRASTRUCTURE EFFICIENCY

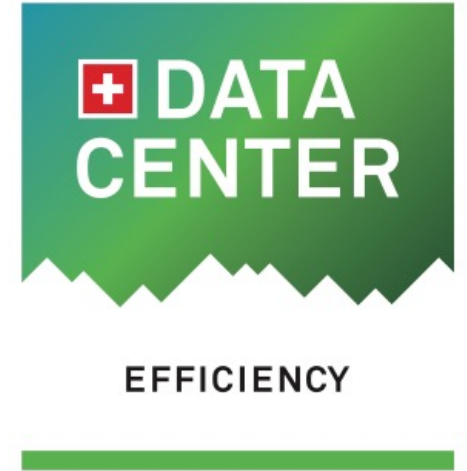
- electricity w/ renewables, cooling, heat recycling

## IT INFRASTRUCTURE EFFICIENCY

- + compute, storage, network and workloads

## DC CARBON FOOTPRINT

- + emissions from input electricity sources



Visit [sdea.ch](https://sdea.ch) to find out  
about our label





# DATACENTER BEST PRACTICES

Need:

- Metrics for datacenter output
  - E.g., matrix multiply in Python is 10x more work than in C
- Metrics for chip design
  - E.g., speedup not a great metric for accelerators
- Need life cycle from fabrication to recycling
  - Half of the emissions are from fabrication



# SUMMARY

## Post-Moore datacenters:

- Integration + Specialization + Approximation
- Revisit legacy abstractions, SW/HW interfaces
- Holistic algorithm/SW/HW co-design
- Division of control vs. data plane

## Datacenter sustainability:

- Best practices
- Metrics



THANK YOU!

For more information, please visit us at  
[parsa.epfl.ch](http://parsa.epfl.ch)

**EPFL**