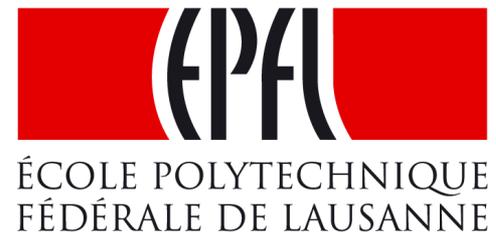


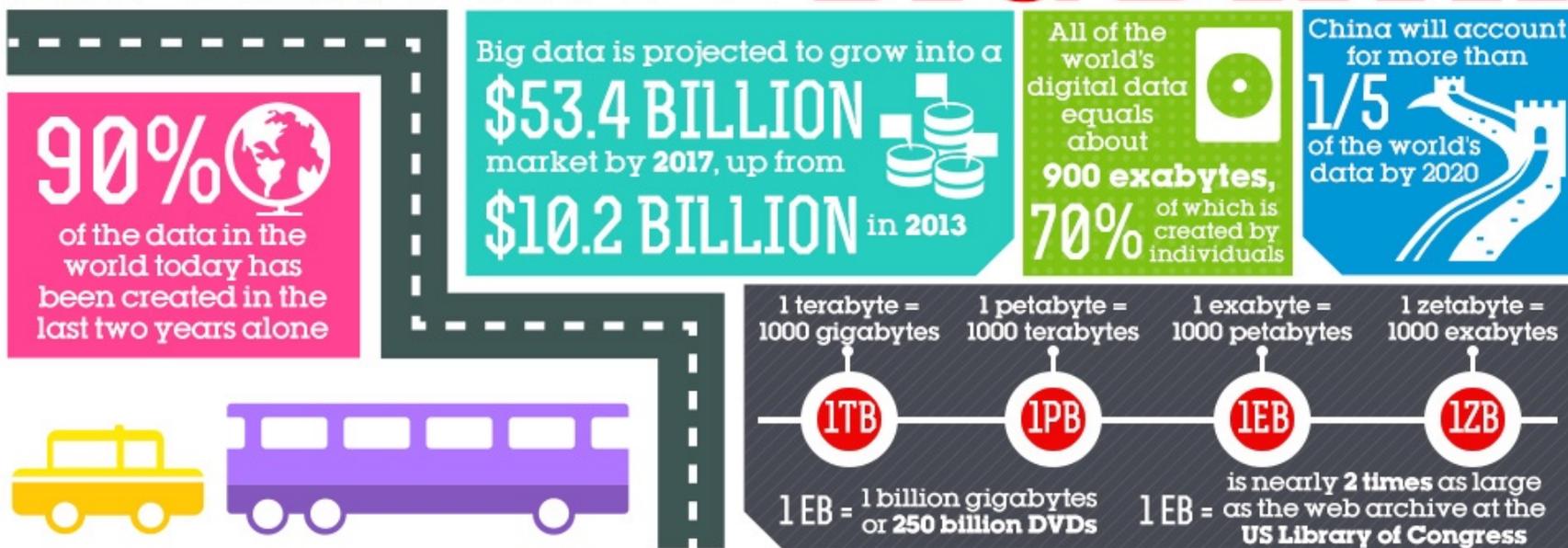
# Silicon Heterogeneity in the Cloud

Babak Falsafi  
ecocloud.ch





## THE LANDSCAPE OF **BIG DATA**



# Modern Datacenters are Warehouse-Scale Computers

- Millions of interconnected home-brewed servers
- Centralization helps exploit economies of scale
- Network fabric provides micro-second connectivity
- At physical limits
- Need sources for
  - Electricity
  - Network
  - Cooling



20MW, 20x Football Field  
\$3 billion

# Warning!

## Datacenters are not Supercomputers

- Run heterogeneous data services at massive scale
- Driven for commercial use
- Fundamentally different design, operation, reliability, TCO
  - Density 10-25KW/rack as compared to 25-90KW/rack
  - Tier 3 (~2 hrs/downtime) vs. Tier I (upto 1 day/downtime)
  - .....and lots more

Datacenters are the IT utility plants of the future

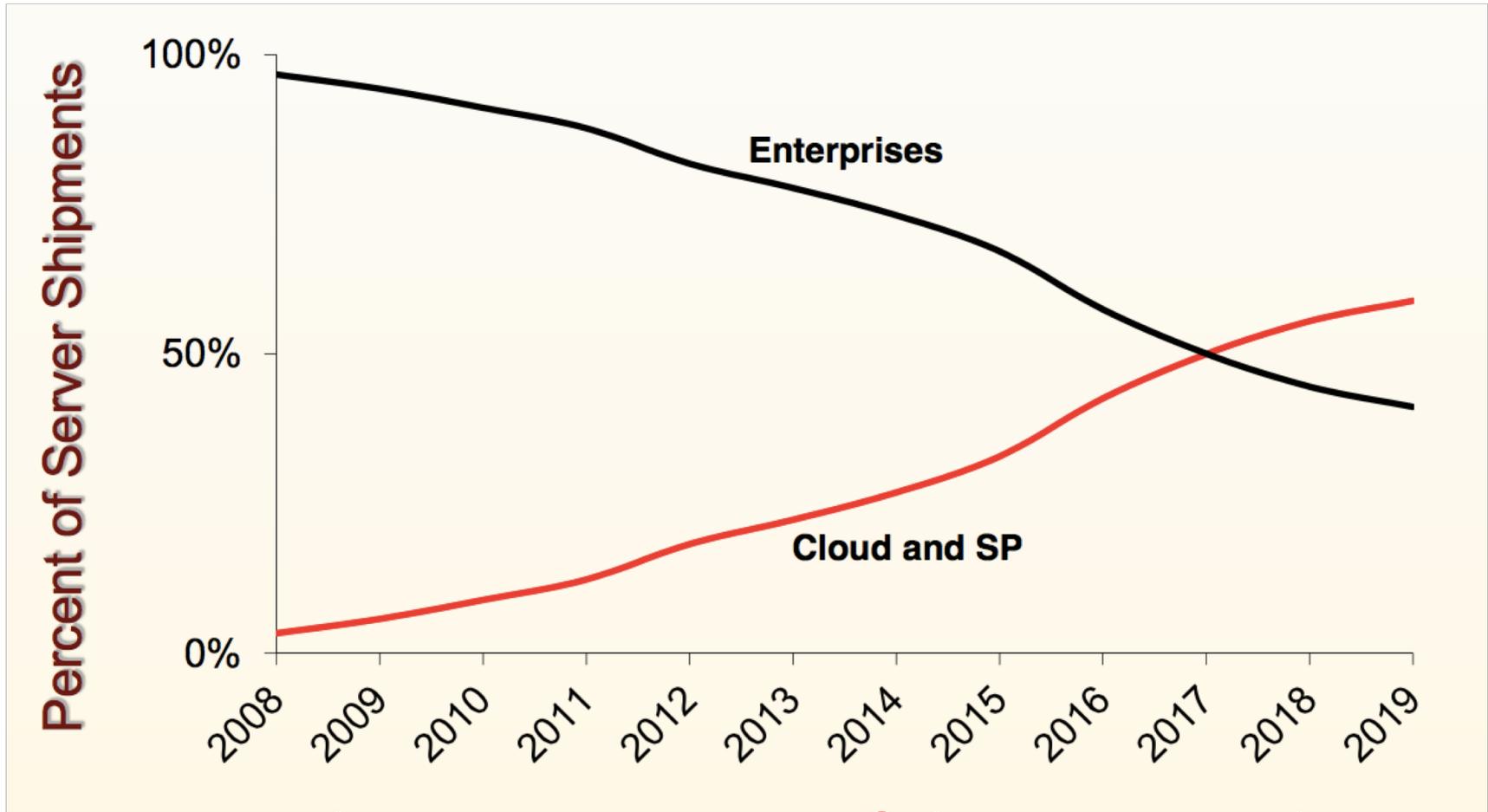


Supercomputing



Cloud Computing

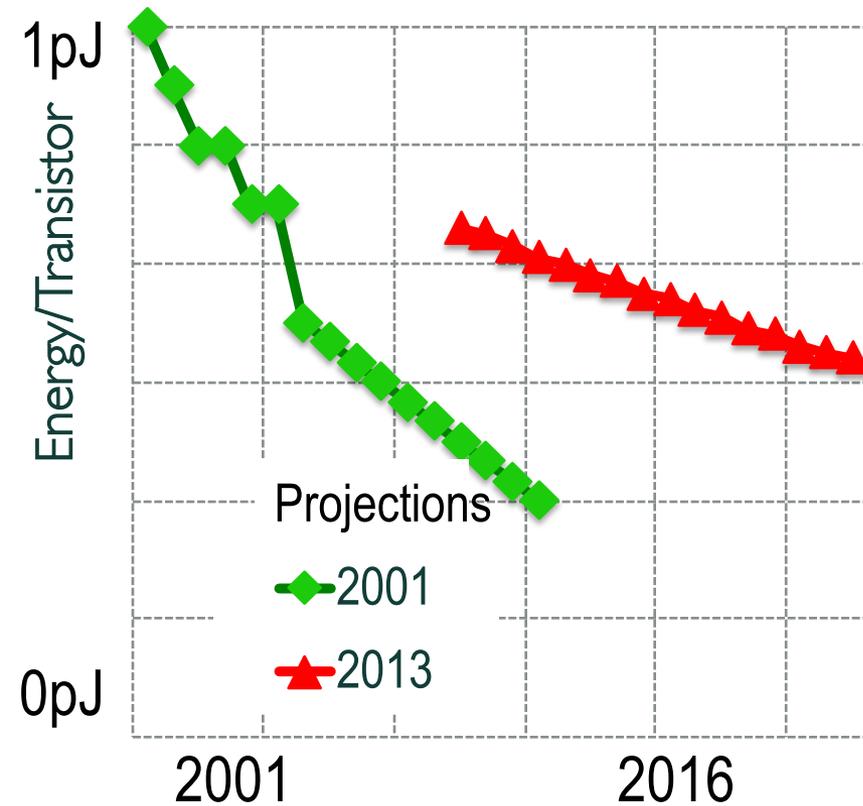
# Cloud Taking Over Enterprise



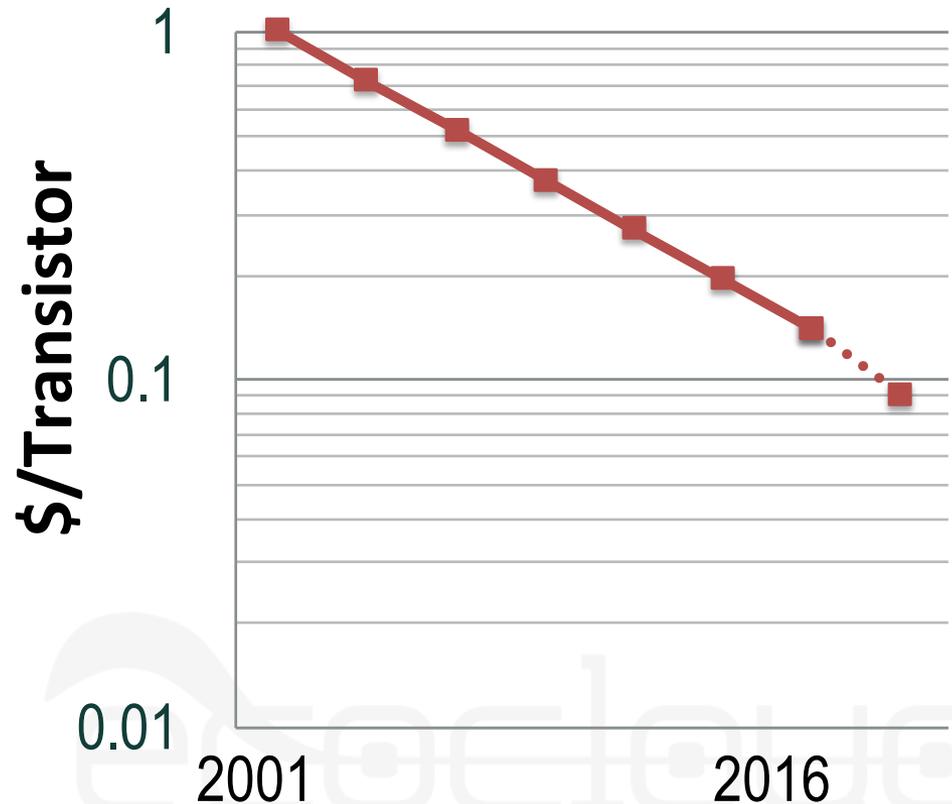
Source: Dell 'Oro 2Q15

# But, Silicon out of steam!

Silicon efficiency is dead  
(long live efficient silicon)



Moore's Law is Dead too!  
[Mark Bohr's Keynote, ISSCC'15]

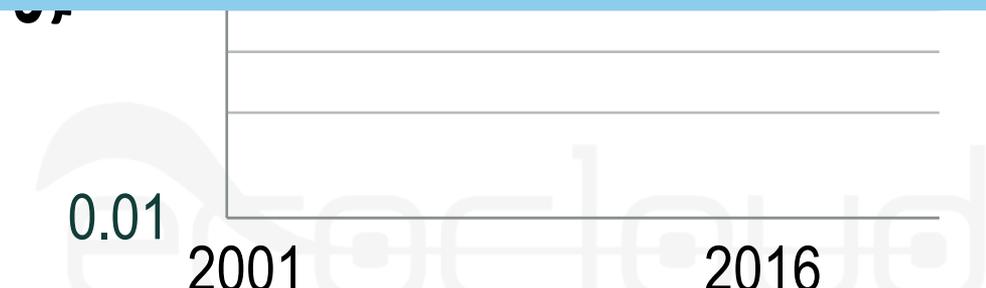


# But, Silicon out of steam!

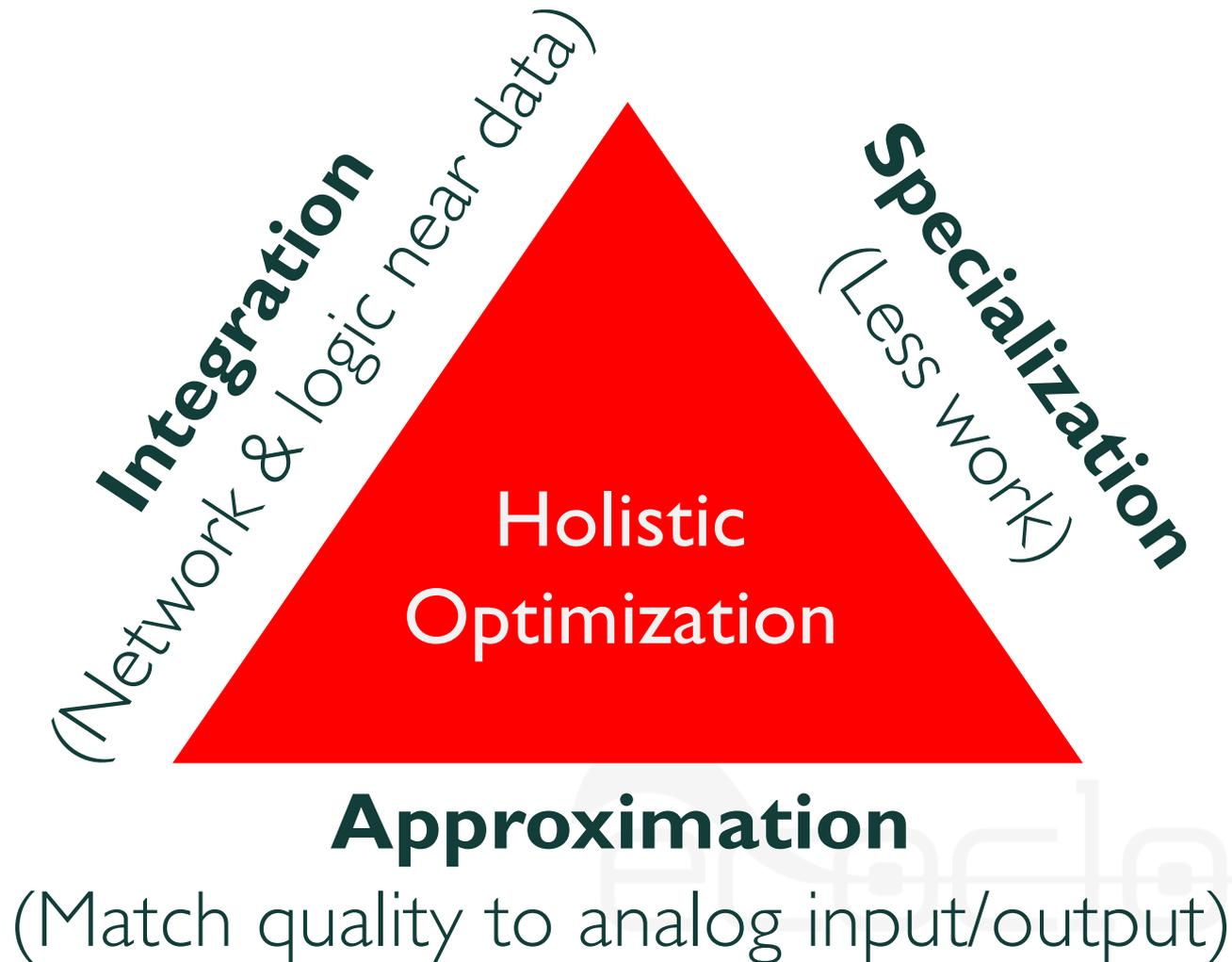
Silicon efficiency is dead  
(long live efficient silicon)

Moore's Law is Dead too!  
[Mark Bohr's Keynote, ISSCC'15]

## Global Foundries cancelled their 7nm on August 28, 2018!



# Optimization Opportunities: The ISA Triangle



# Scale-Out Datacenters

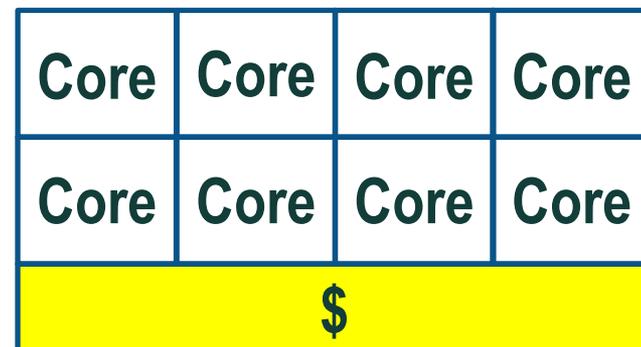
Vast data sharded across servers

Memory-resident workloads

- Necessary for performance
- Major TCO burden

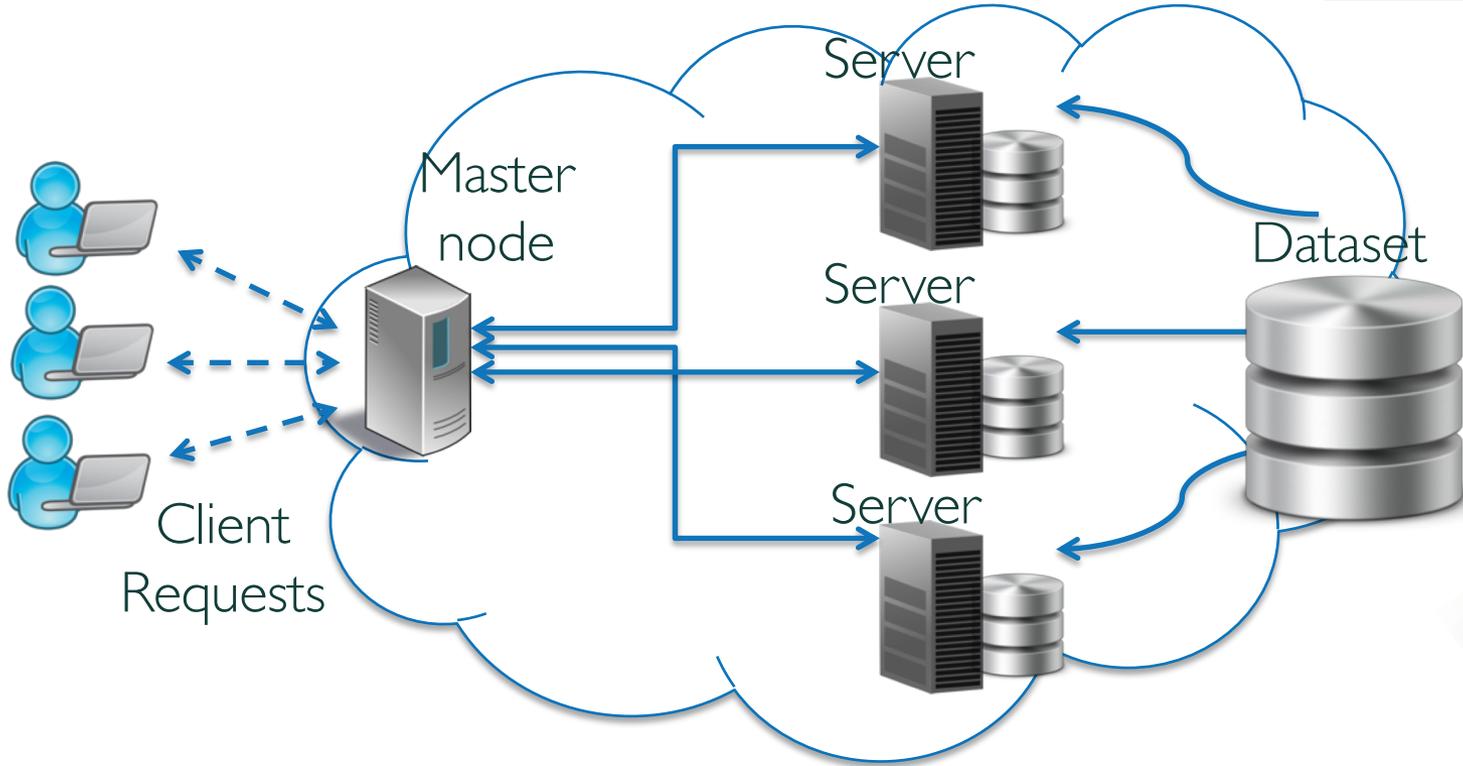
Put memory at the center

- Design system around memory
- Optimize for data services



Servers driven by the DRAM market!

# In-Memory Scale-Out Services



- Many independent requests/tasks
- Huge dataset split into shards
- Use aggregate memory over network



# Server Benchmarking with CloudSuite 3.0 (cloudsuite.ch)

Data Analytics  
Machine learning



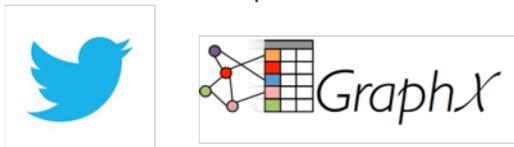
Data Caching  
Memcached



Data Serving  
Cassandra NoSQL



Graph Analytics  
GraphX



Media Streaming  
Nginx, HTTP Server



Web Serving  
Nginx, PHP server



Web Search  
Apache Solr & Nutch



In-Memory Analytics  
Recommendation System



**Building block for Google PerfKit, EEMBC Big Data!**

# Scaling CPU's: Manycores

- Parallelism has emerged as the only silver bullet
- Use simpler cores
  - Prius instead of Audi R8
- Restructure software
  
- Each core → fewer joules/op

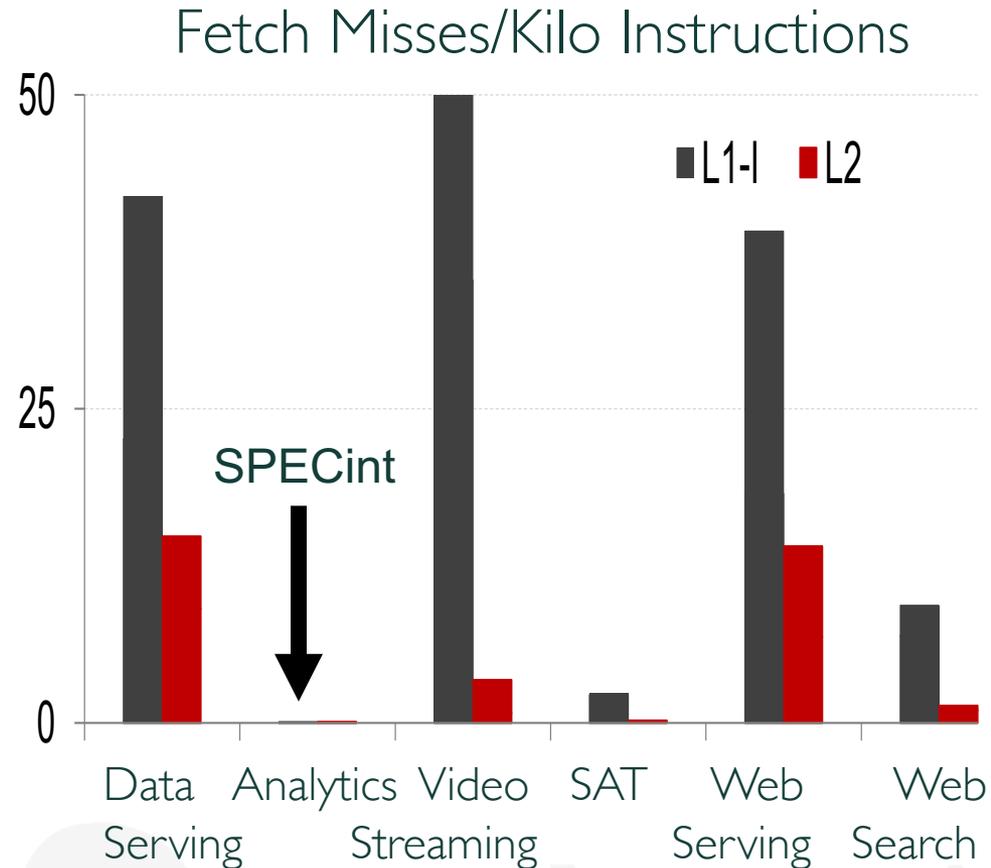
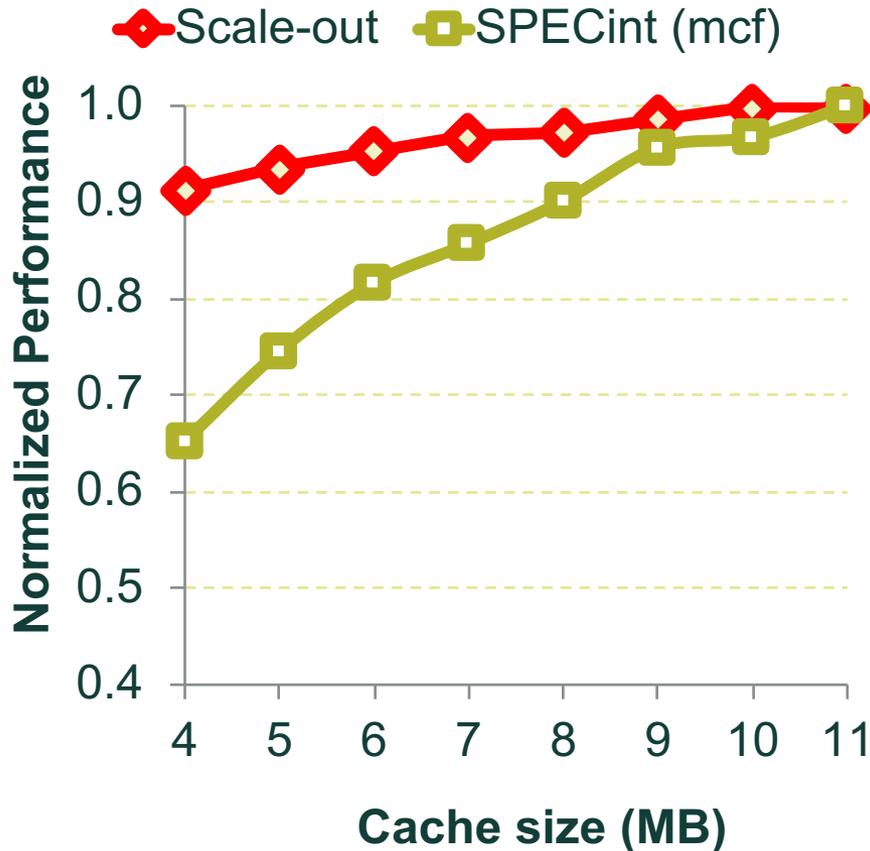
Conventional Server  
CPU (e.g., Xeon)



Modern Manycore  
CPU (e.g., Tilera)

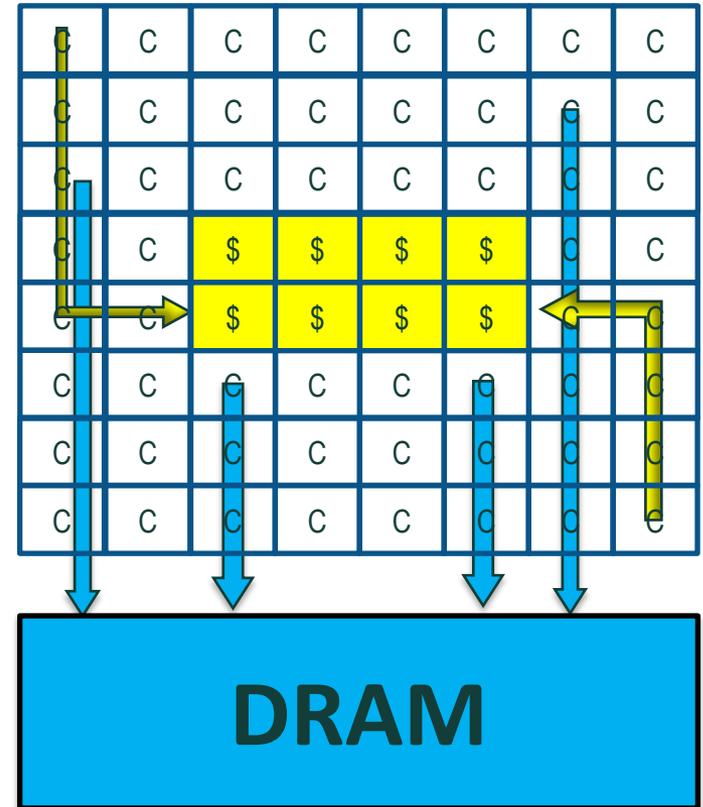
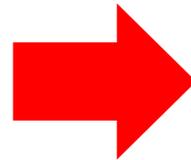
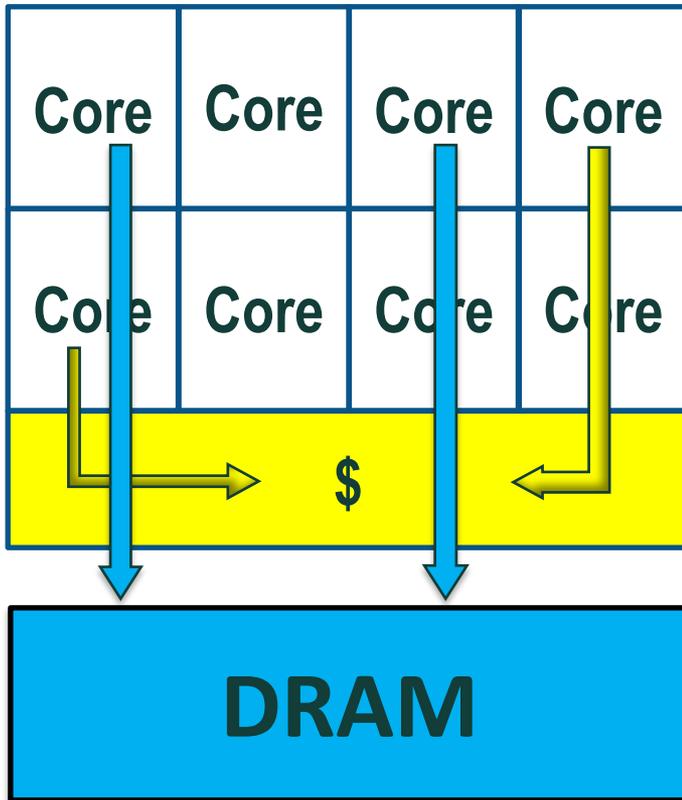


# But, Services Stuck in Memory (x86 servers) [ASPLOS'12]



- On-chip memory overprovisioned
- Instruction supply is bottlenecked

# Scale-Out Processors (SOP)



- x86 server CPU
- ✗ Logic 60% of silicon
- ✗ 6x bigger cores

- 3-way ARM manycore
- ✓ Logic 85% of silicon
- ✓ 7x more parallelism

# Innovation in SOP's

[ISCA, MICRO, ASPLOS, MemSys, IEEE Micro'12-18]

## Instruction supply:

- Core front-end (BP/BTB, Boomerang [Grot])

## On-chip networks:

- Core-to-\$ rather than core-to-core

## Off-chip connectivity:

- HBM, DRAM hierarchy, network



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

**Gopal Hegde**

VP/GM, Data Center Processing Group

## Cavium Thunder X

- Based on SOP @ EPFL
- Designed to serve data
- Optimized code supply
- Trade off SRAM for cores
- Runs stock software
- 10x faster than Xeon for CloudSuite

# Massively parallel cores

- Data parallelism
- Higher memory b/w

Super simple cores

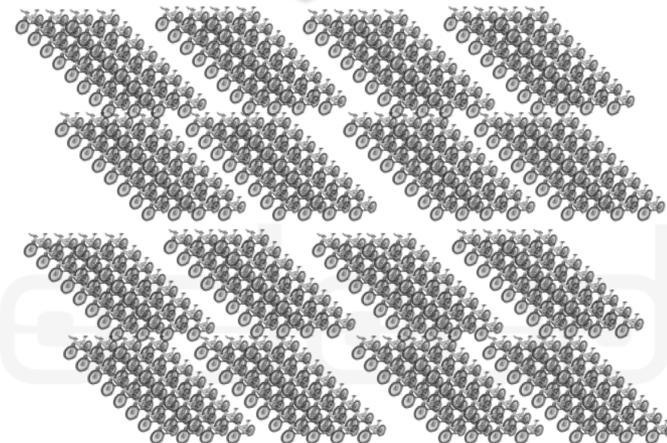
- Shared front end
- 10x slower clocks

Great for dense parallel computation

Conventional Server  
CPU (e.g., Xeon)



Modern GPU  
(e.g., Volta)

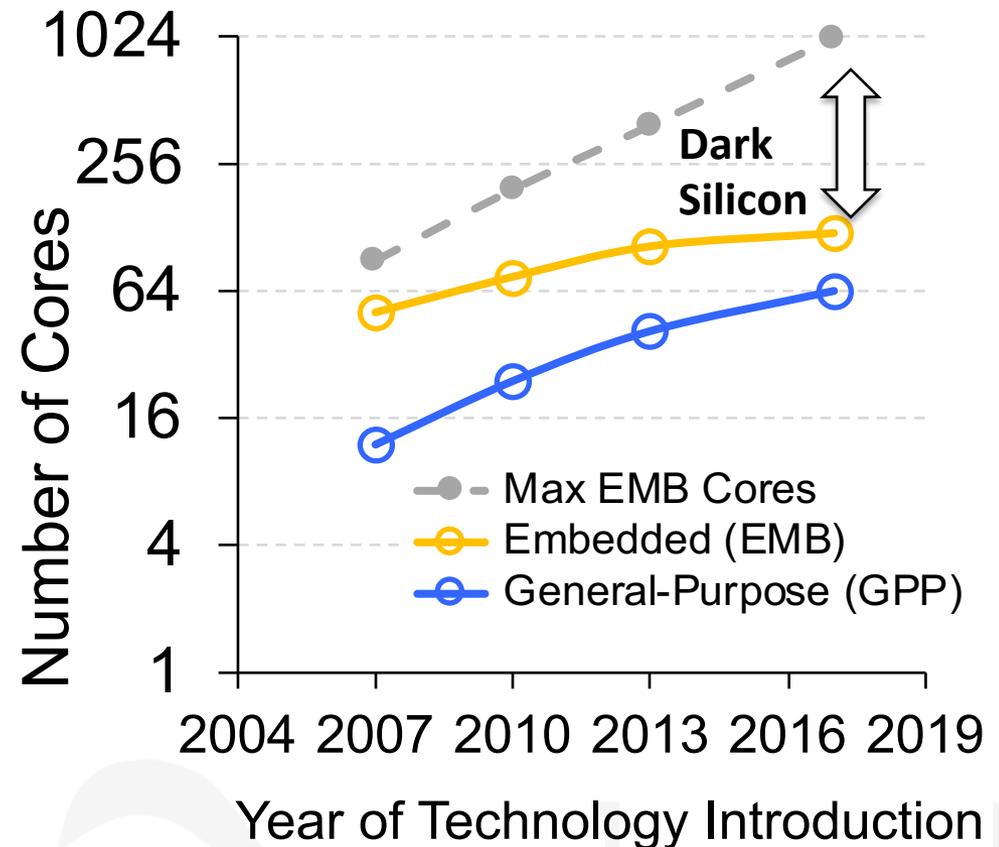


# Parallelism Alone Can't Help

Can populate chips  
But, can not operate all  
Today's chips are already  
“dark” (memory)

All future platforms will be  
heterogeneous

- Selectively activate parts



[source: Hardavellas et. al., “Toward Dark Silicon in Servers”, IEEE Micro, 2011]

# Custom Computing

[FPGA's vs. GPU's in Data centers, IEEE Micro'17]

Reconfigurable

- Best for spatial computing
- Not caching/reuse

Parallel, spatial computing

- 10x slower clocks
- Better for sparse arithmetic

Microsoft, Amazon & Intel

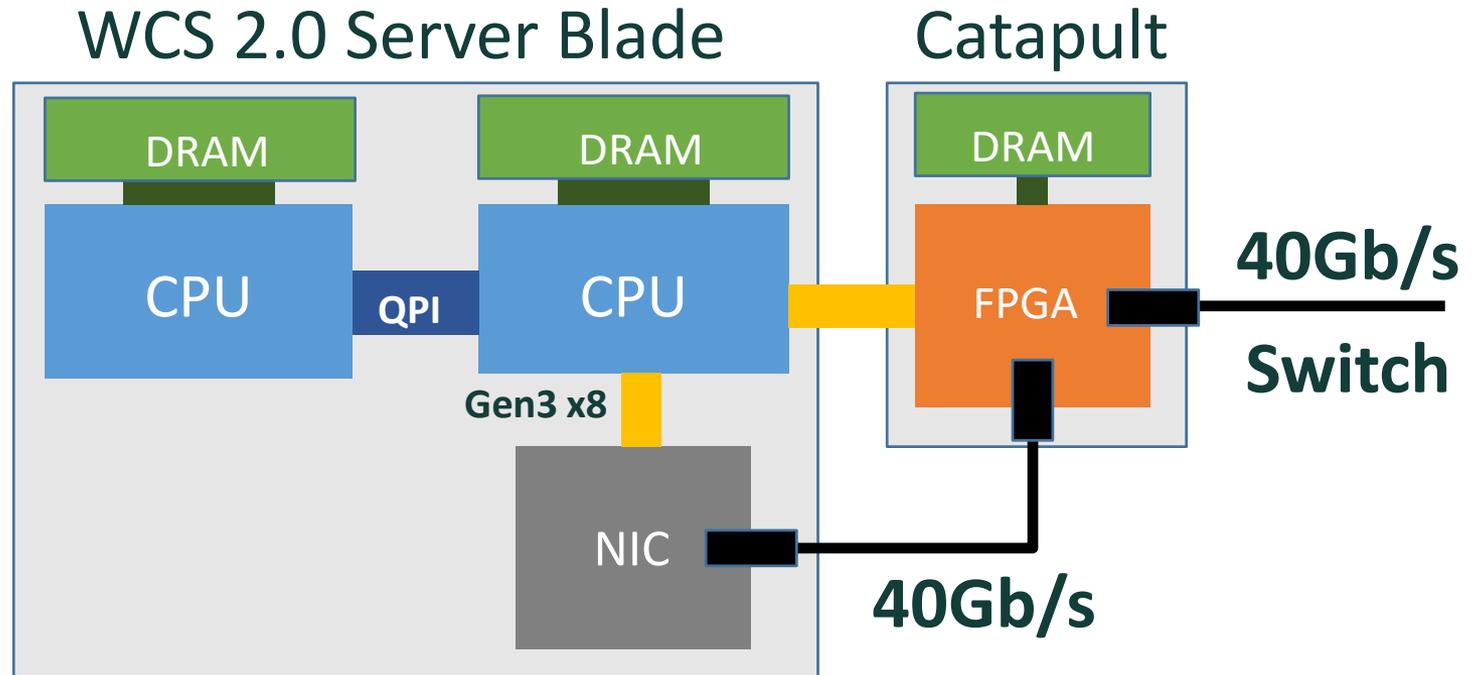
Conventional Server  
CPU (e.g., Xeon)



FPGA  
(e.g., Catapult)



# Microsoft's Catapult



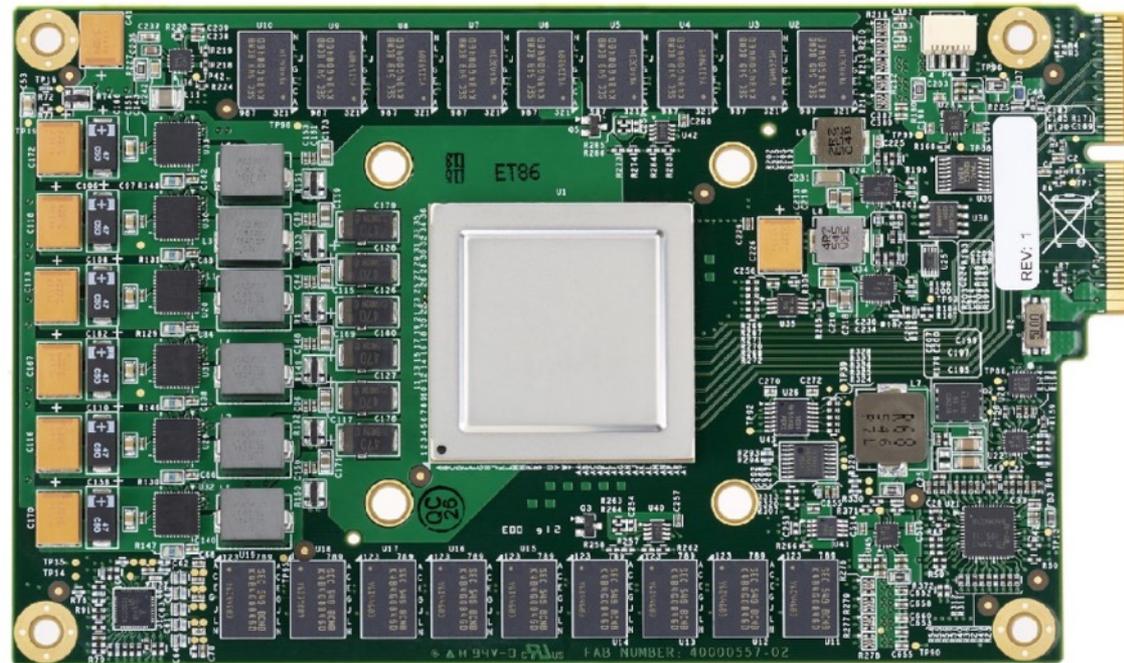
FPGA economies of scale:

1. Local/remote compute accelerator
2. Network/storage accelerator
3. Configurable cloud

# Google's TPU

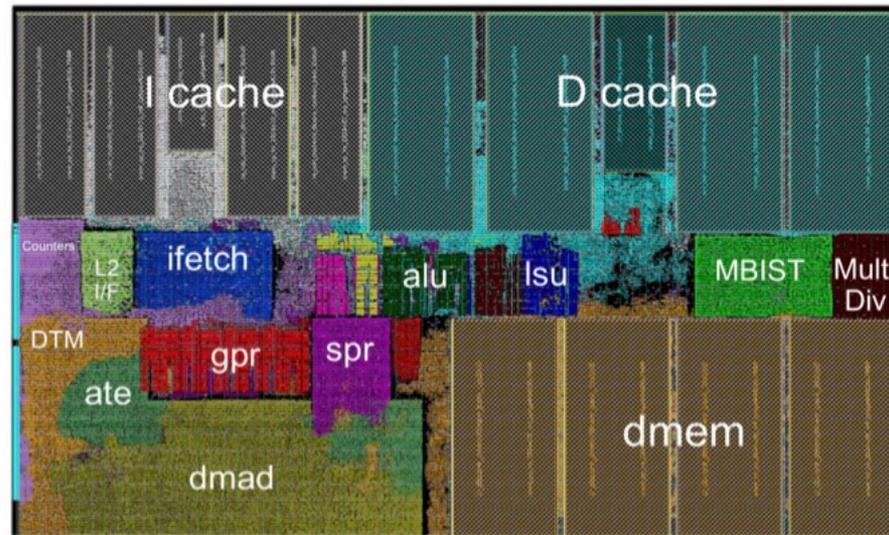
Custom array of arithmetic units:

- Linear algebra for ML/NN
- Currently memory bound
- 10x over GPU
- ML as a service



# Oracle's RAPID

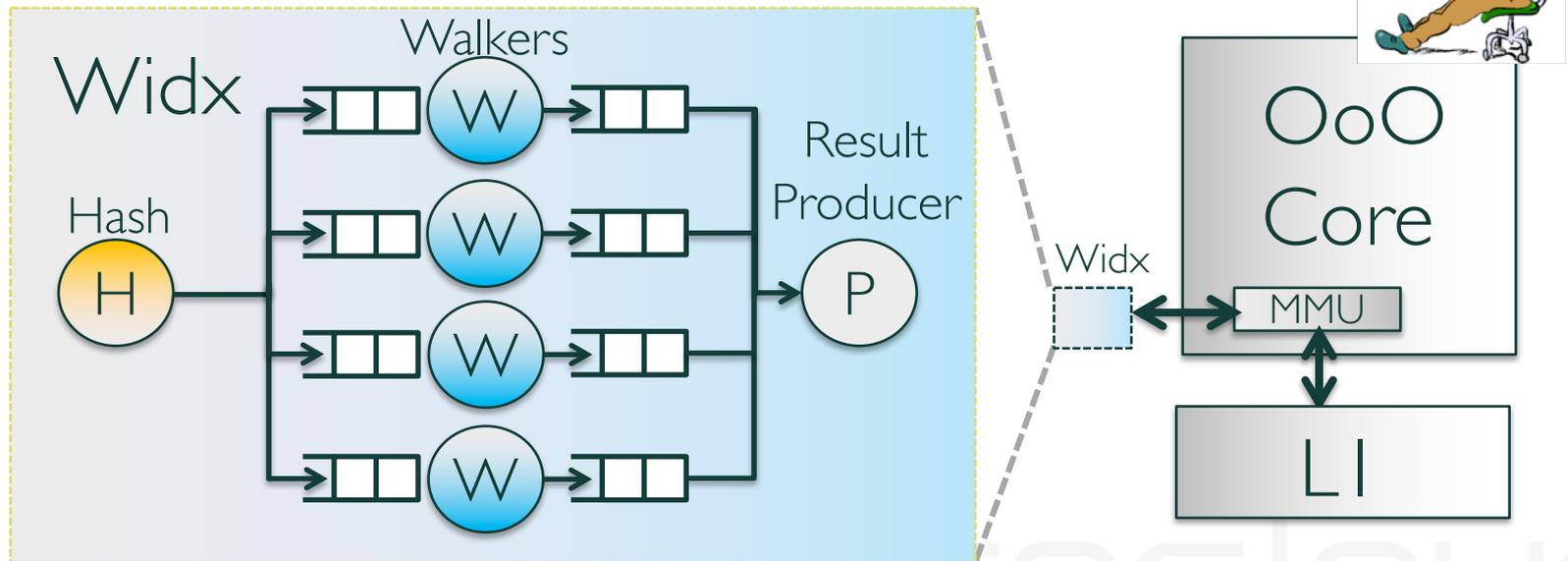
- Accelerator for analytics in SQL
- Data movement engine in hardware
- Custom message passing cores
- Up to 15x better perf/Watt over Xeon



ecocloud

# Walkers: CPU-Side Database Accelerators

- Pointer-based data structures (e.g., hash table, B-tree)
- Parallel lookups require traversing chains
- Decouple chains in co-designed hw/sw



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

# Walkers in Software [VLDB'16]

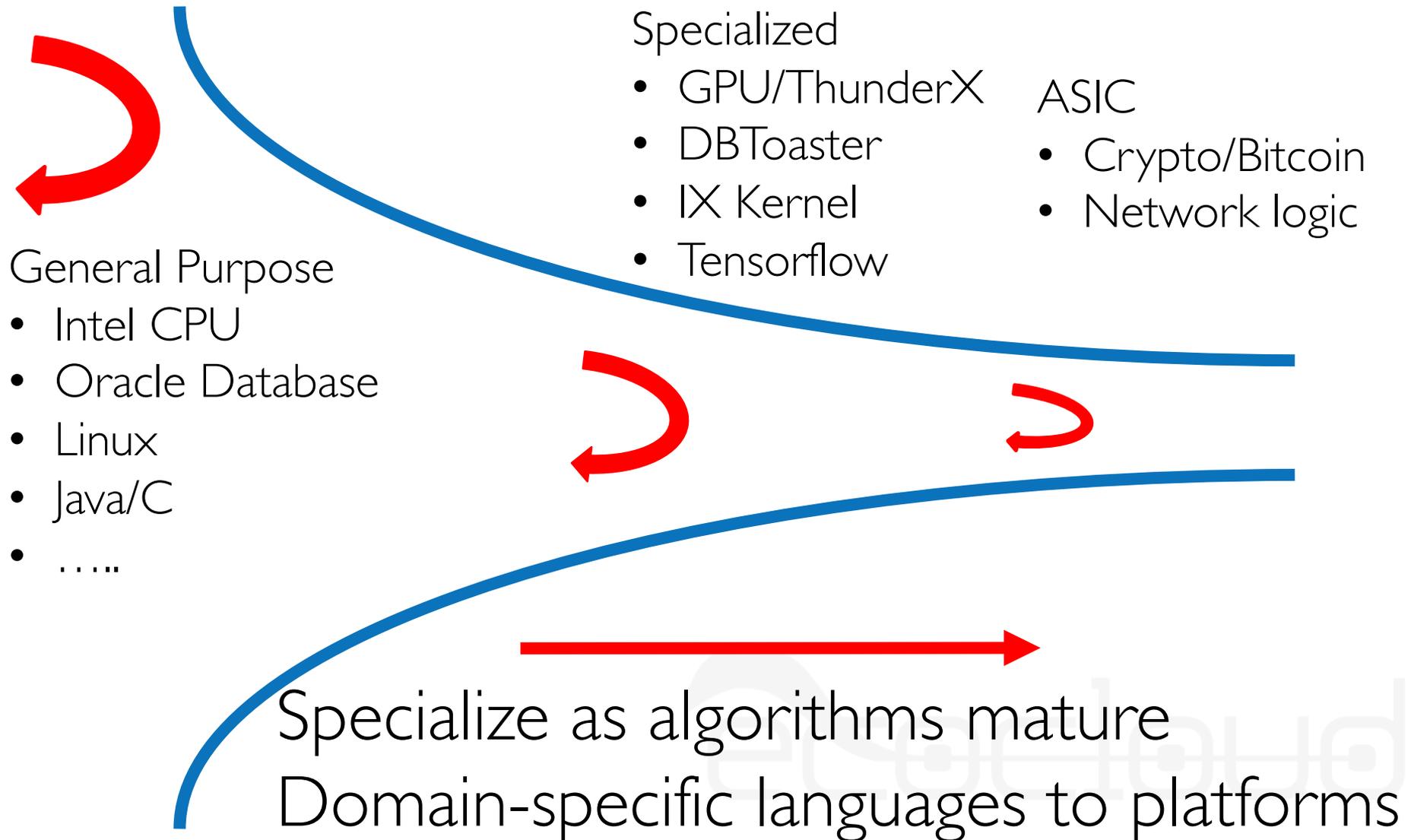
Use insights to help Xeon

- Decouple hash & walk 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
- To be integrated in SAP HANA [VLDB'18]

# Moving Forward: The Specialization Funnel



# Approximation

Modern apps/services are statistical

- Analog input, analog output

Key:

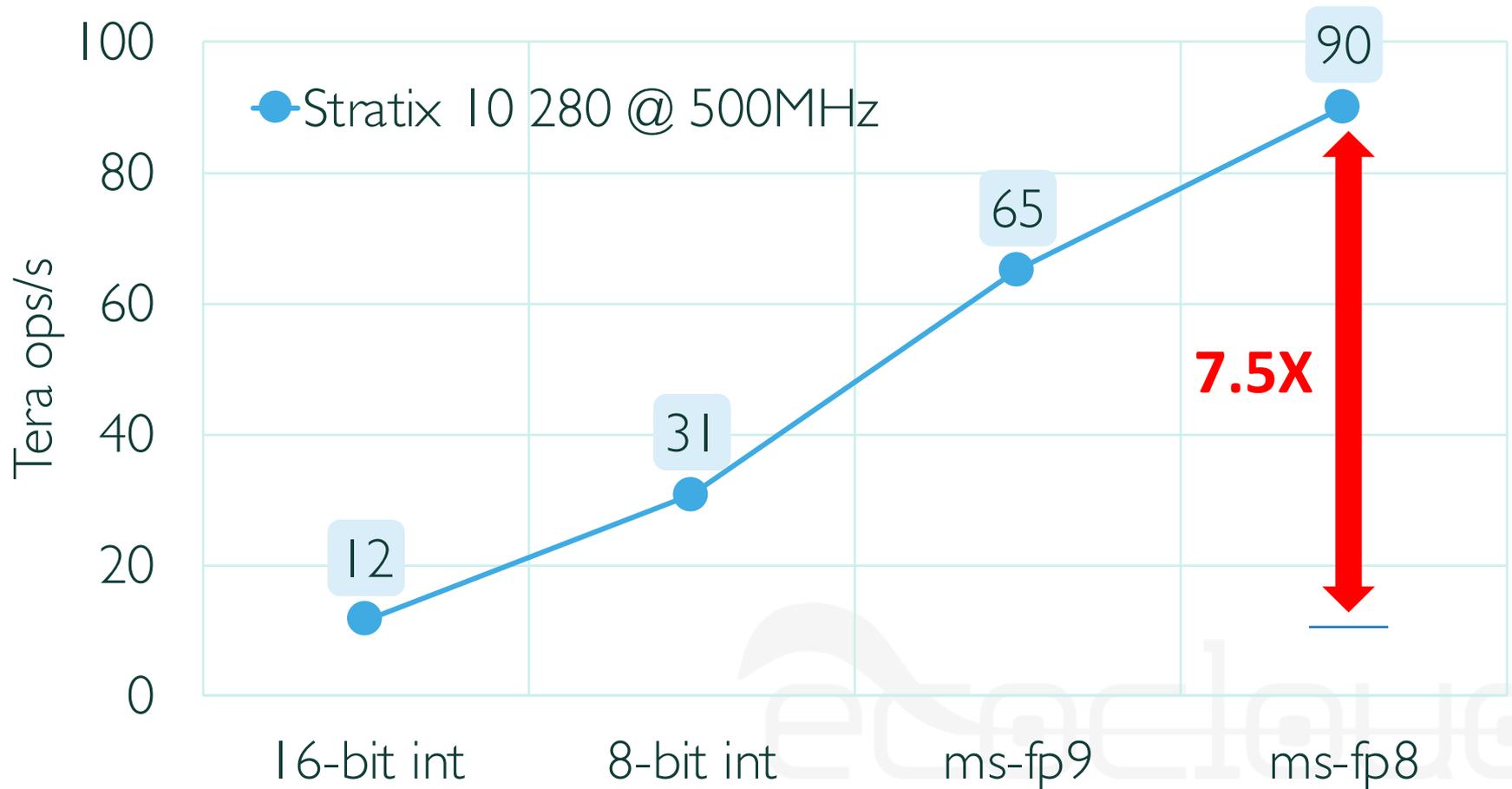
- Much redundancy in data/arithmetic
- Output quality not accuracy or error

Exploit in

- Processing, communication, storage

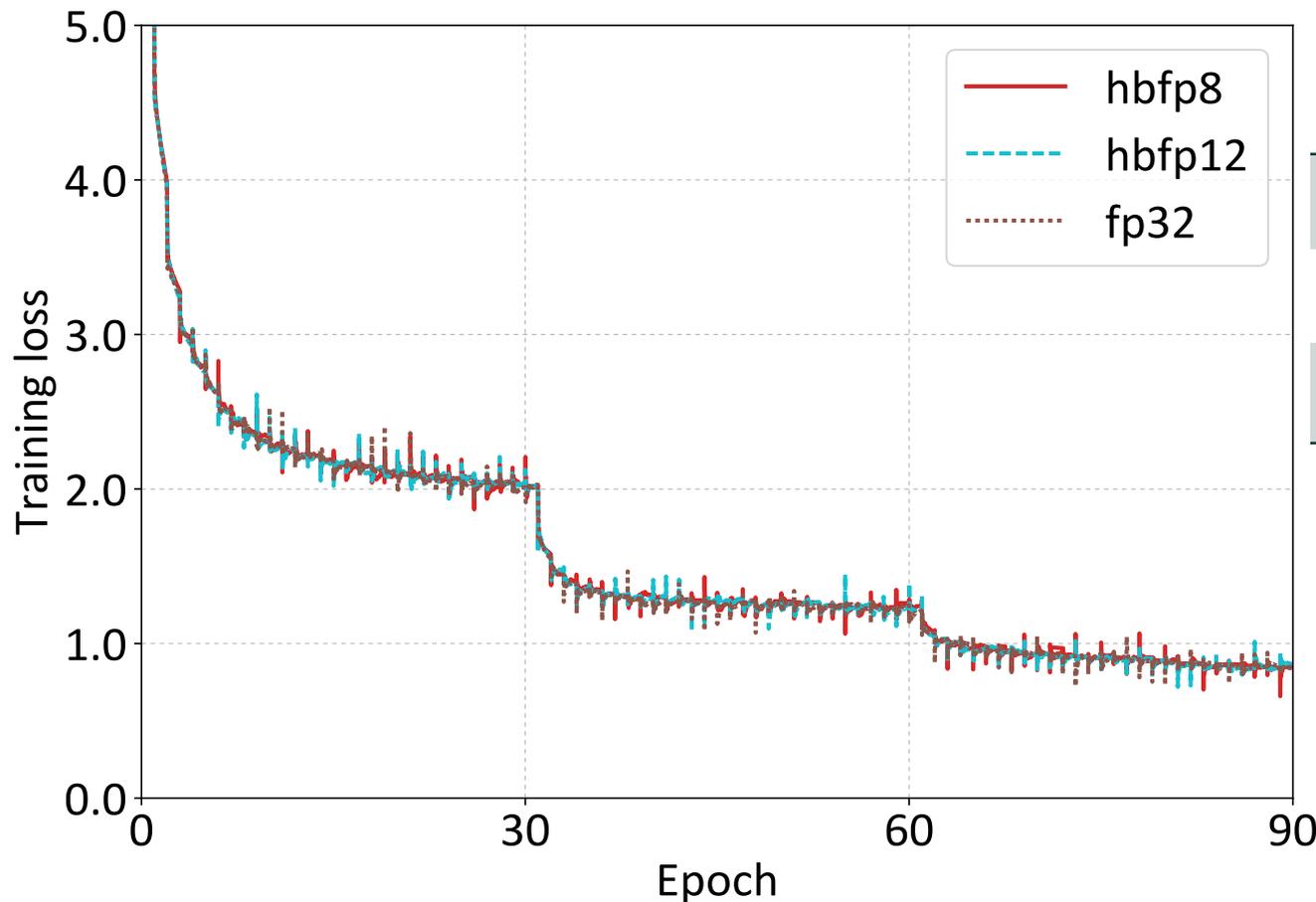
# Arithmetic in Deep Learning (Microsoft Brainwave)

## FPGA Performance vs. Data Type



# HBFP (Block FP) vs. FP32

Resnet-50 on ImageNet



Config.	Top-1 Error (%)
HBFP8	23.88
HBFP12	23.58
FP32	23.64



FP32 performance with 8-bit logic [NeurIPS'18]

# Memory Hierarchy

Faster



Bigger



**Regs**

**Caches  
(SRAM)**

**Main memory  
(DRAM)**

**SSD**

**Hard disk**

Today

**Regs**

**Caches  
(SRAM)**

**3D Caches  
(DRAM)**

**Main memory  
(DRAM)**

**Storage-Class  
Memory**

**SSD**

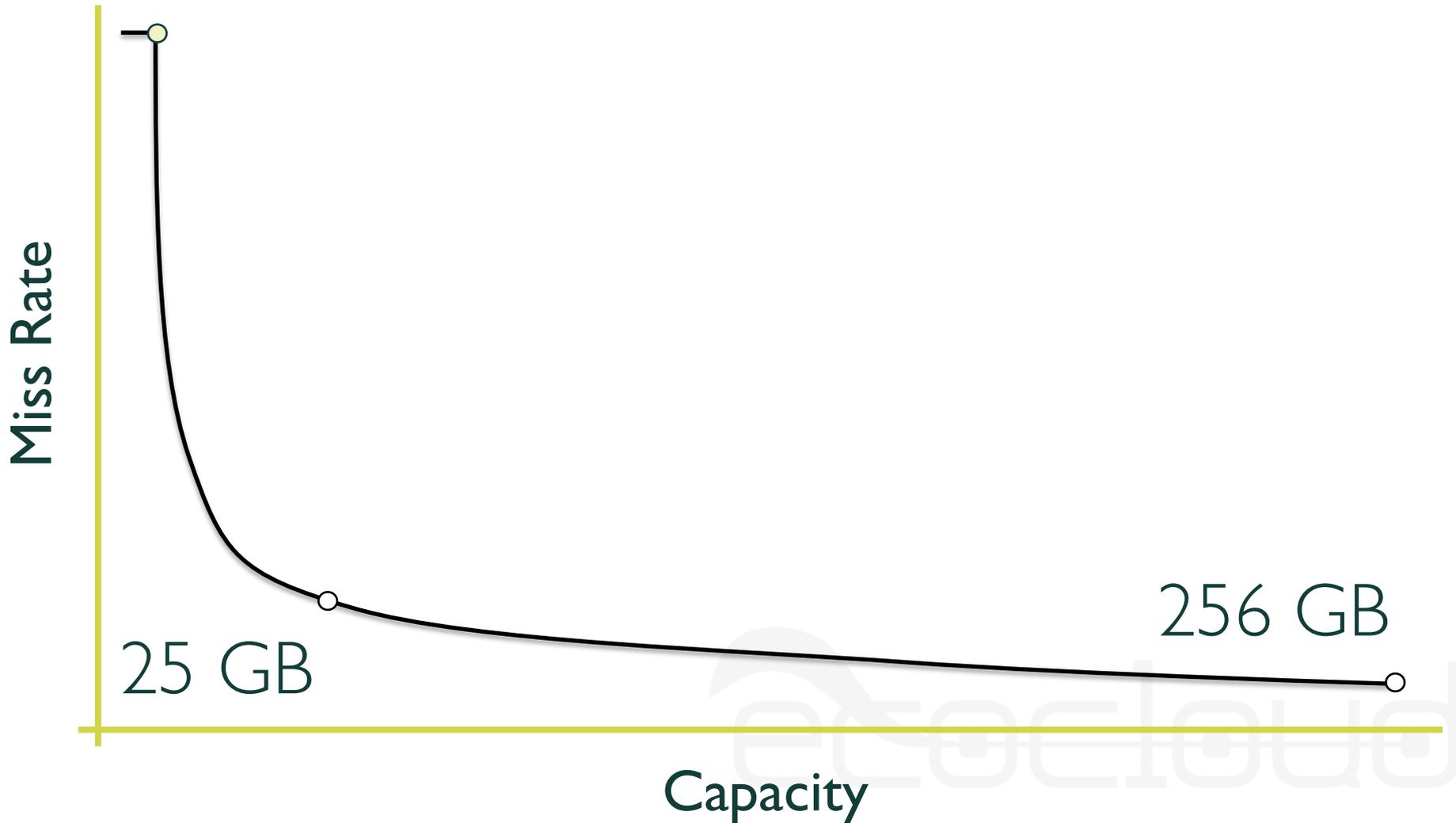
**Hard disk**

Coming Soon

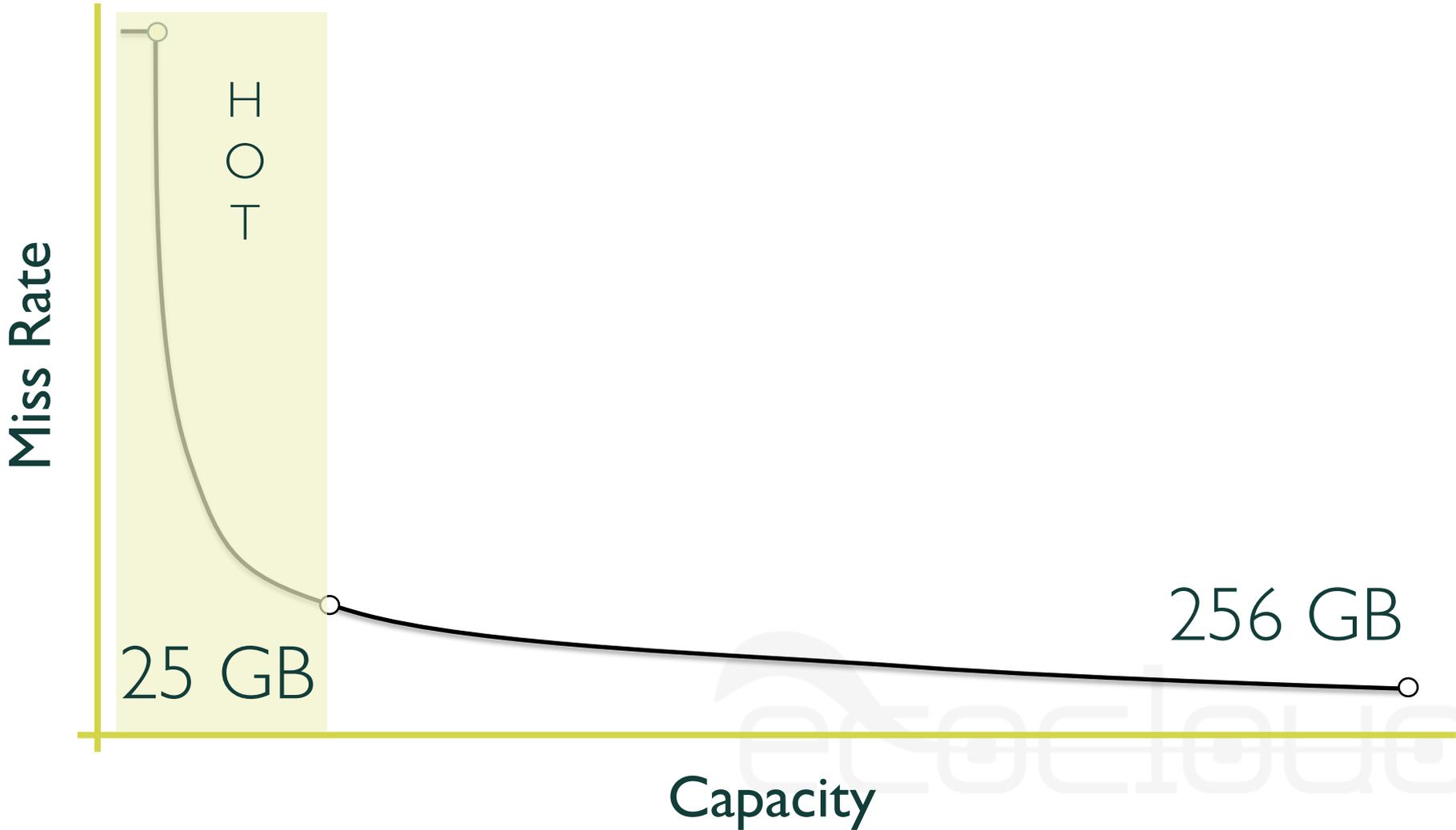
# Capacity/Miss Rate 101



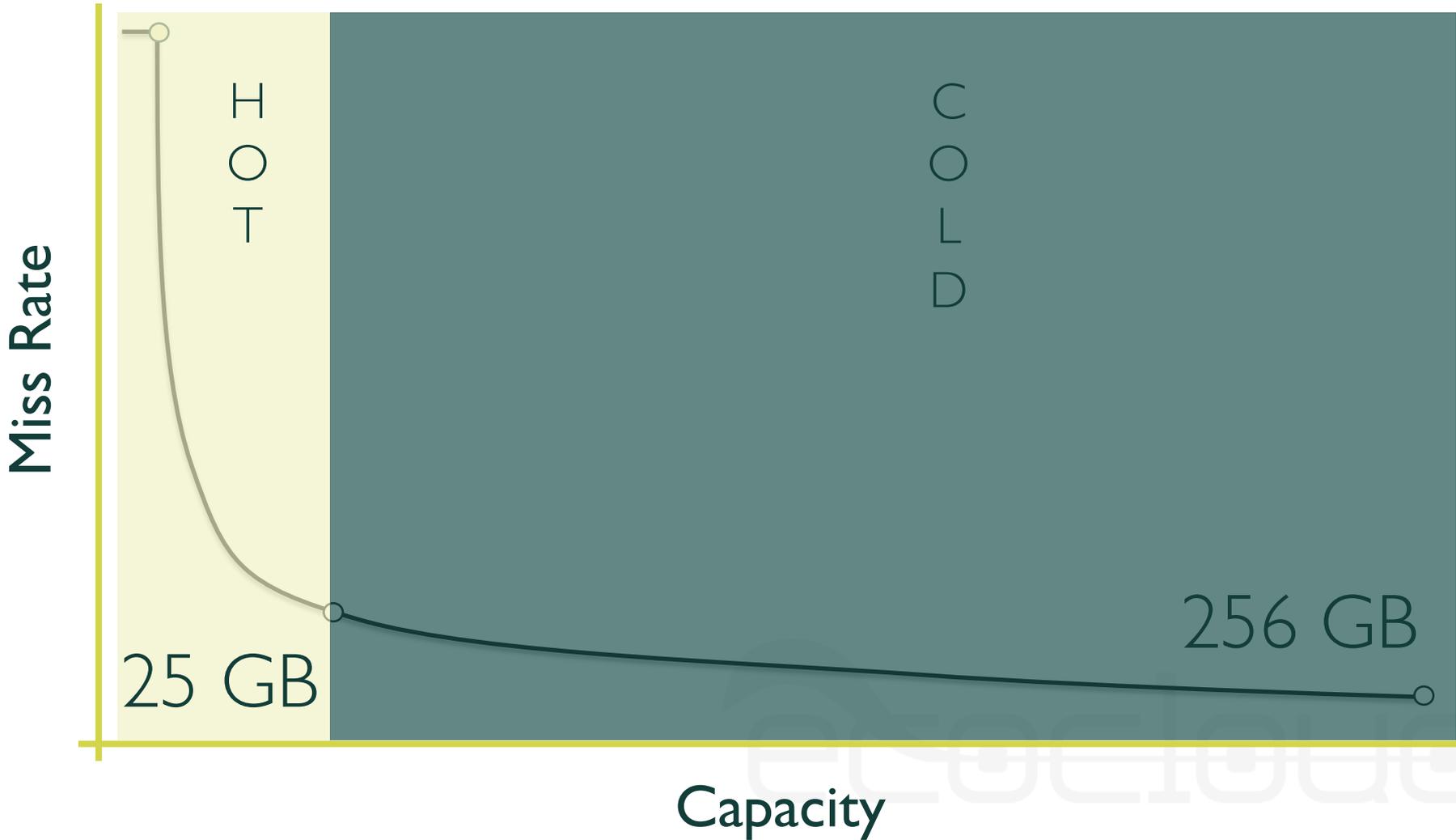
# Capacity/Miss Rate | 0 |



# Capacity/Miss Rate 101

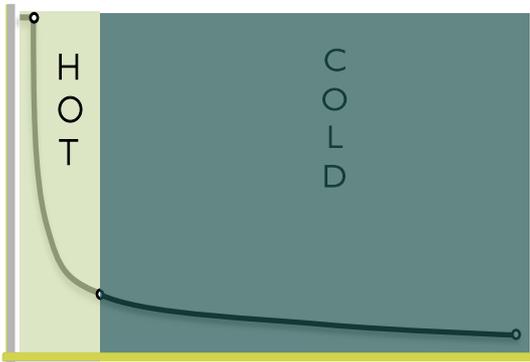
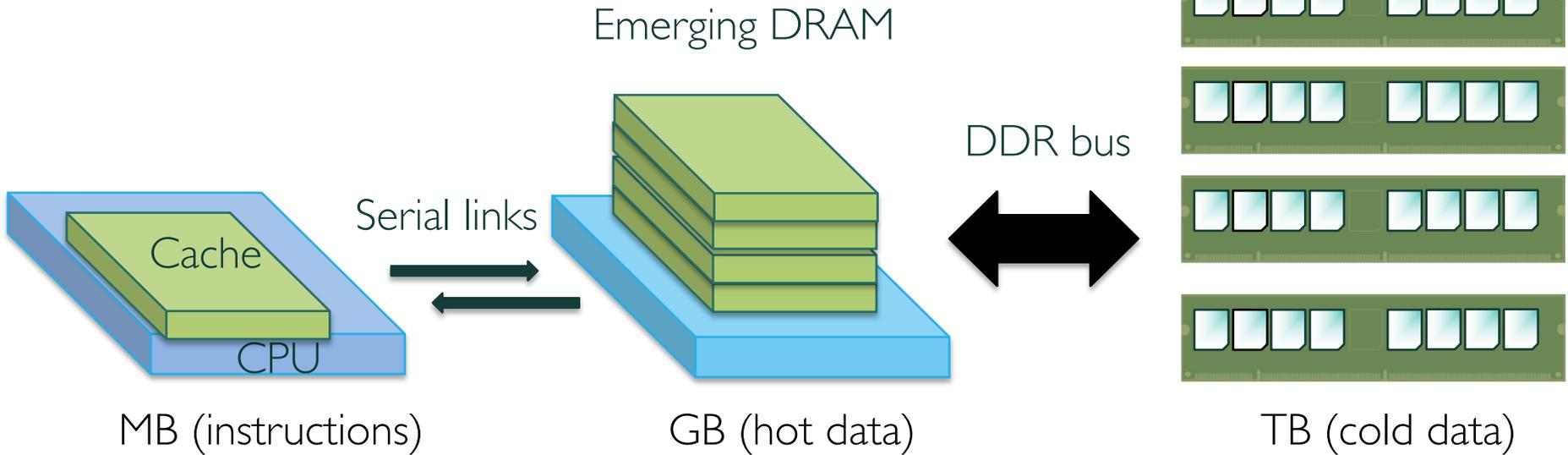


# Capacity/Miss Rate 101



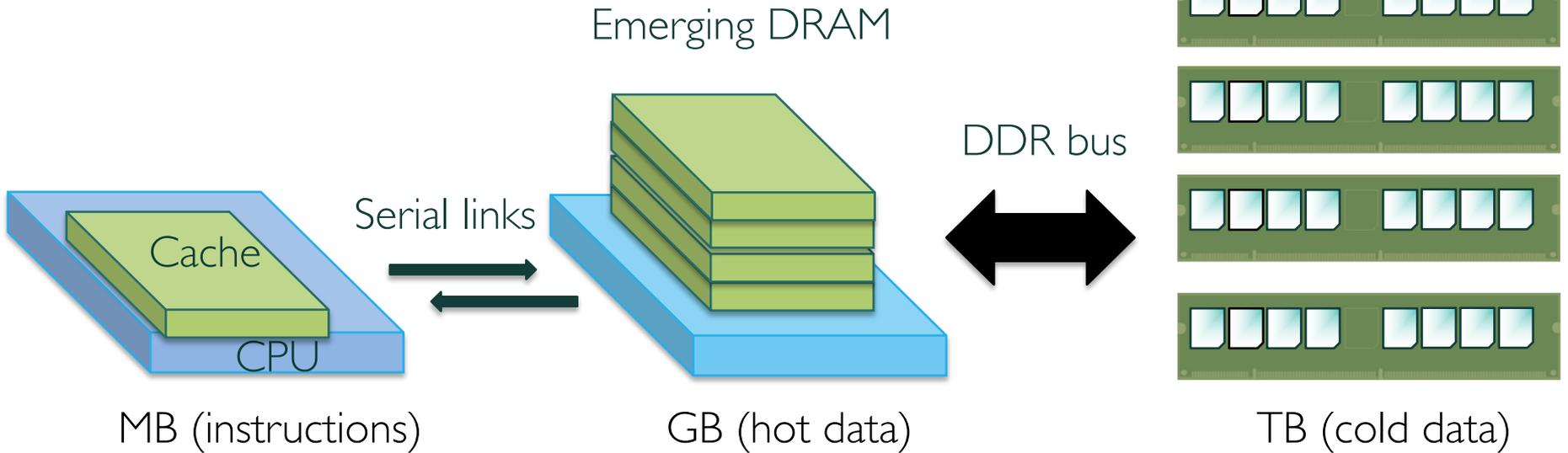
# Emerging Hierarchies

Storage-Class Memory



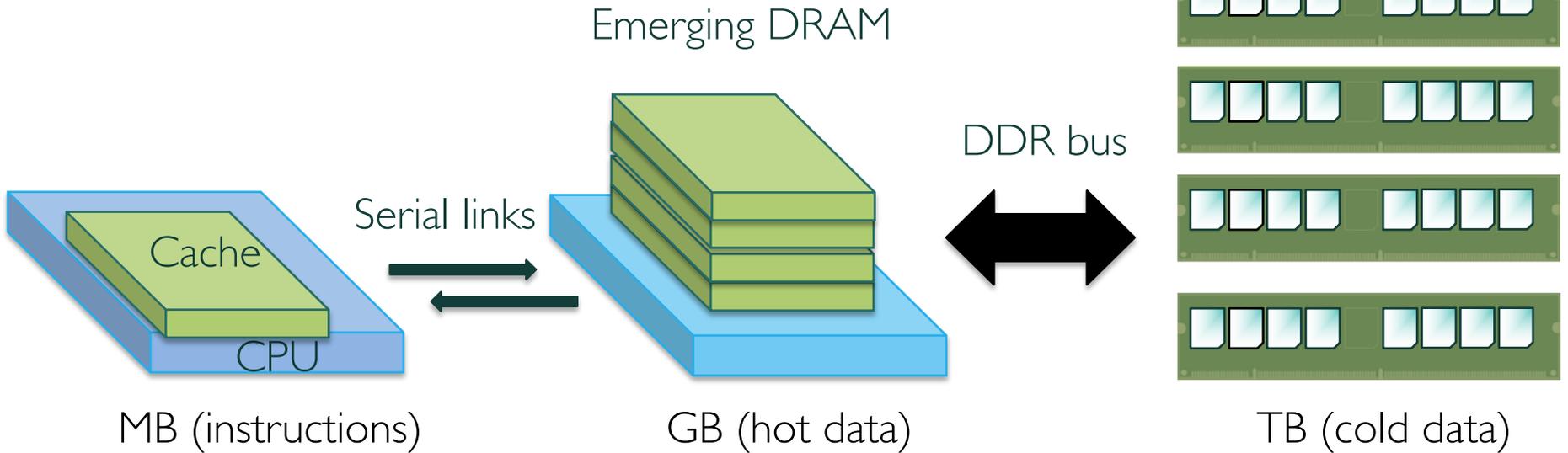
# Emerging Hierarchies

Storage-Class Memory



# Emerging Hierarchies

Storage-Class Memory



# Near-Memory Processing

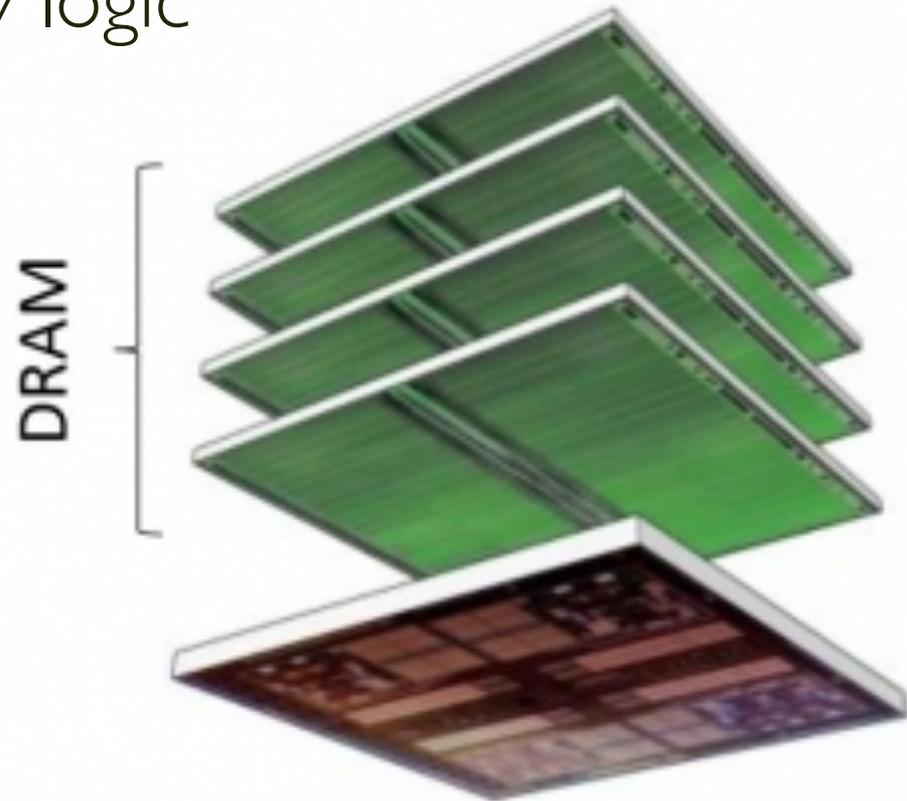
A stack of DRAM w/ nearby logic

- Minimize data movement
- Massive internal bandwidth

Limitations:

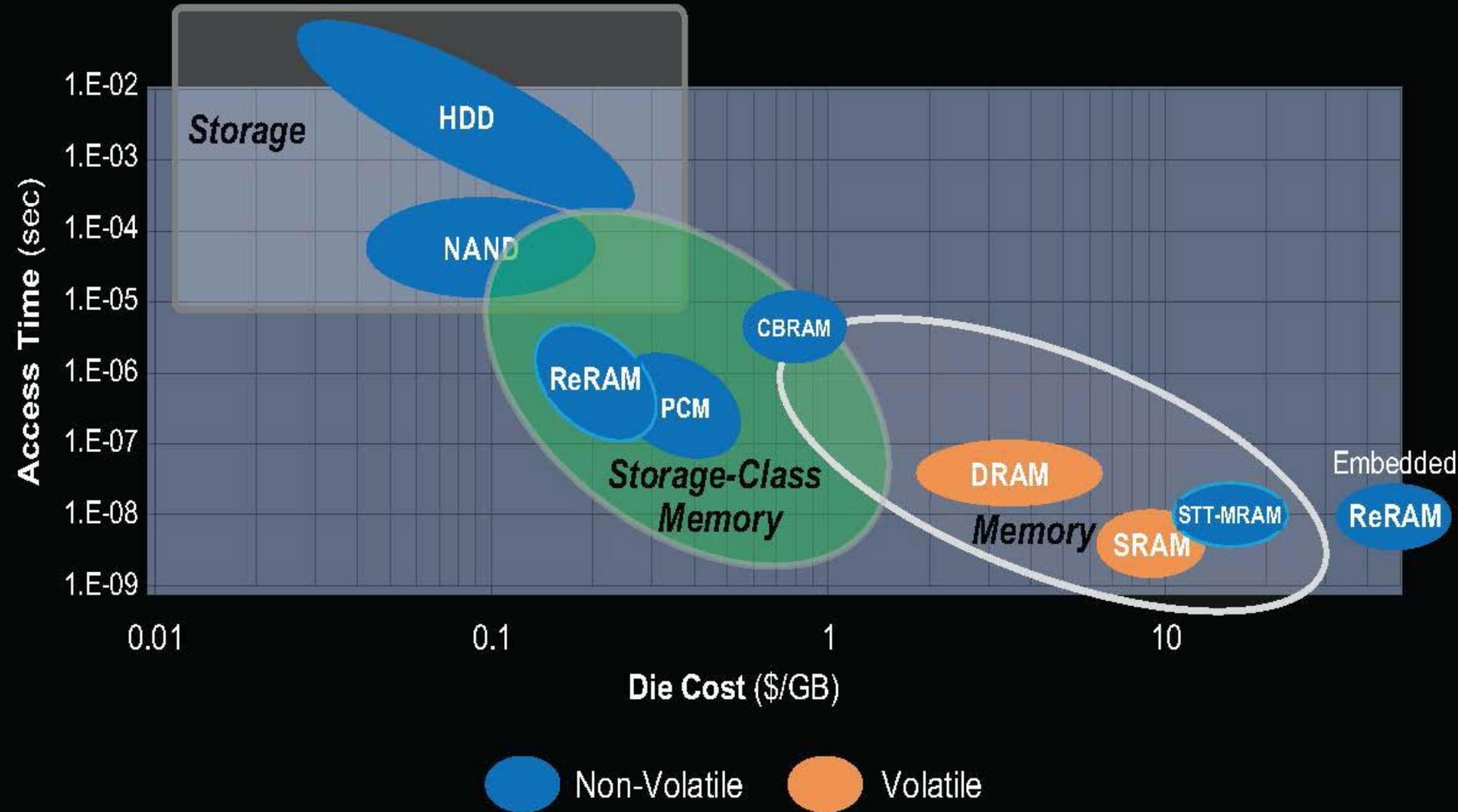
- A few layers of DRAM
- Logic power/thermals (3D)
- Thermals ok for HBM (2.5D)

[source: AMD]



Opportunities for algorithm/hardware co-design

# Memory & Storage Hierarchy



## Persistence

- 100's of nanosecond vs. microsecond
- Implications for logging & networks

## Disparity between reads/writes

- Can read at memory speed
- Writes must be batched/are slow
- Writes consume more power
- DRAM cache can help [MemSys'18]

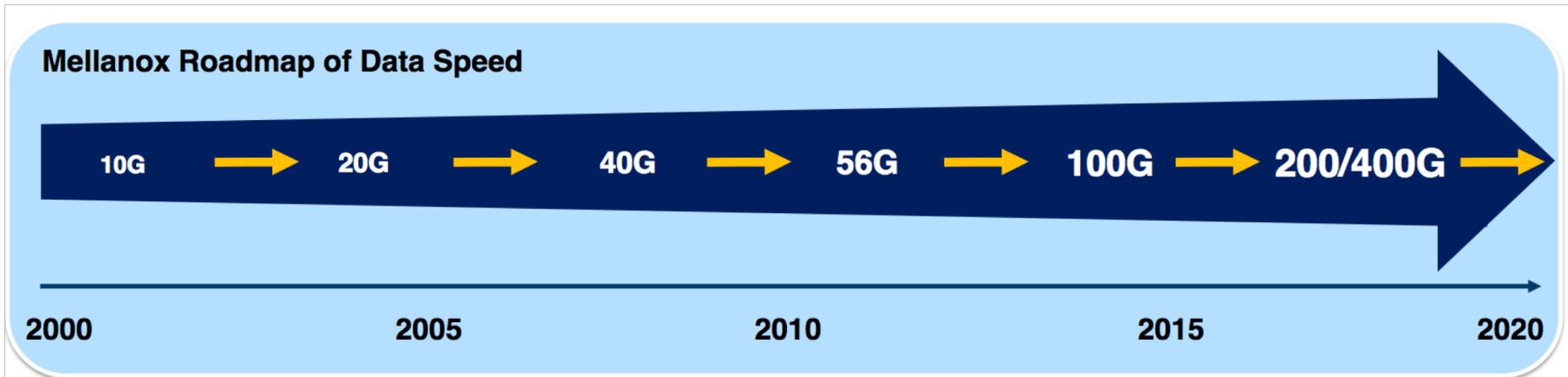
SSD is treated as storage

- Online data in DRAM
- But, DRAM costs dominate, slow scalability

Online services:

- Roundtrip tail latency 100's ms
- SSD access is in 50 us (1000x faster)
- SSD is 50x cheaper

Technology to bring SSD online!



Network stacks/interfaces are a bottleneck:

- Logic growing at 17%/year, network at 20%/year!
- $\mu$ Services emerging
- RPC stacks, scheduling/dispatch, data transformation, ....

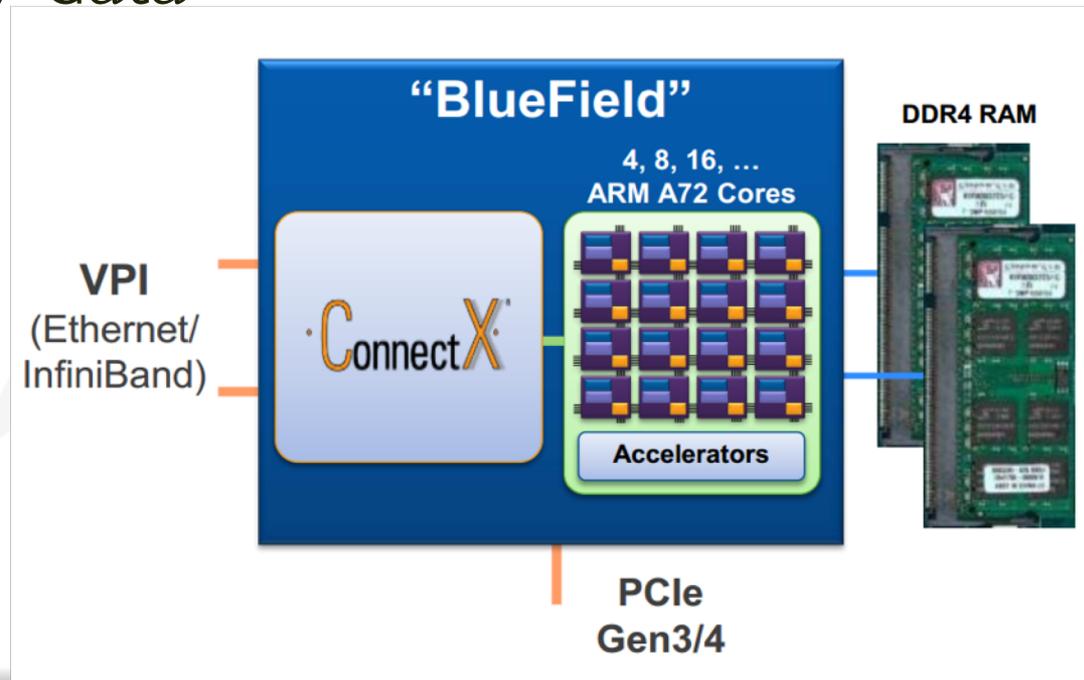
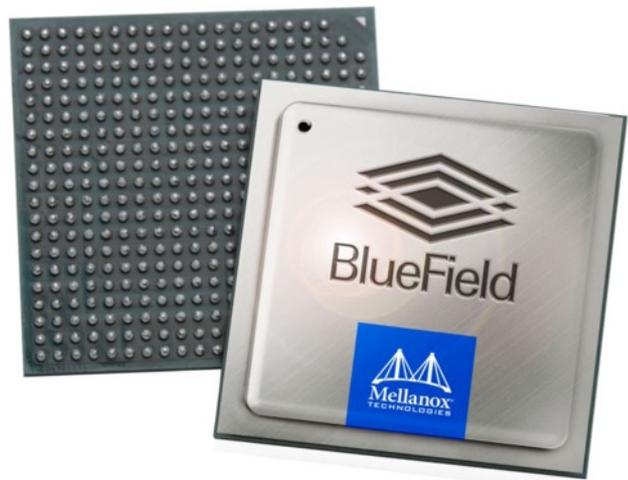
Key challenges:

- Abstractions for control/data planes
- Co-design of network stacks

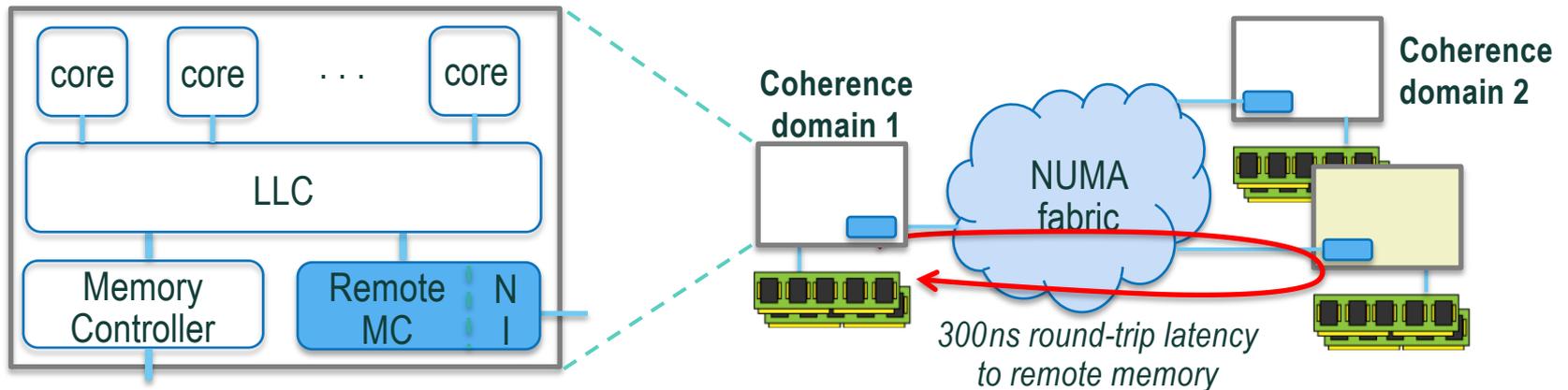
# Near-Network Processing: Nvidia BlueField

Network-interface integrated manycore

- Up to 32 cores w/ 0.5 TB of DRAM
- Can host an in-memory object store
- RPC over in-memory data



# Scale-Out NUMA [ASPLOS'14, ISCA'15, MICRO'16]



soNUMA:



- Socket-integrated network interface
- Protected global memory read/write + synch
- Fine-grain (~64B) & bulk objects (~1 MB)
- Remote memory ~ 2x local memory latency
- Extensions for messaging & RPC [Daglis' thesis]

## Server design centered around data

- In-memory services offered over the network

## Witnessing end of Moore's Law

- Emerging heterogeneous logic + memory

## Future servers nodes:

- Logic & memory with multiple network access points
- Tool chains to go from DSL's → accelerators

**Integrate + Specialize + Approximate**

Thank You!

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



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

