# FireBox: Past, Present, and Future

NATHAN PEMBERTON

12/12/2017

**Compute moving to extremes: edge + cloud**

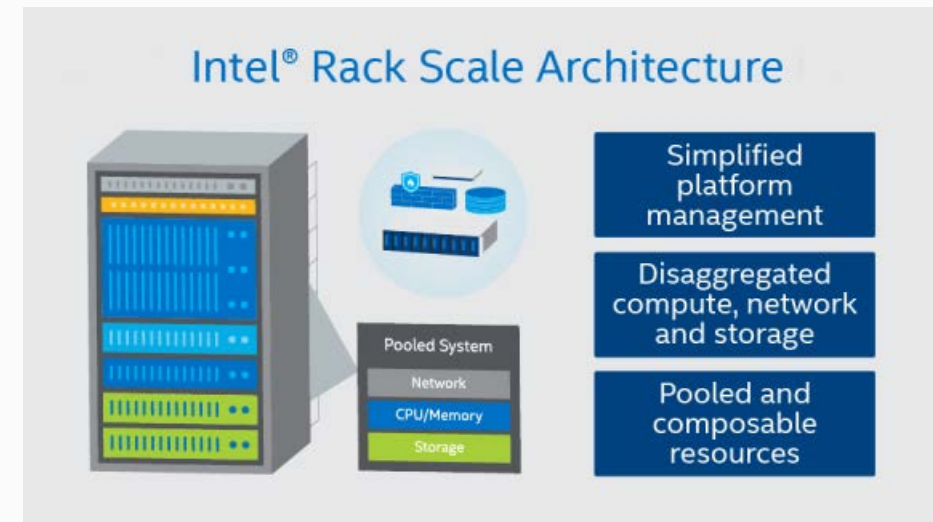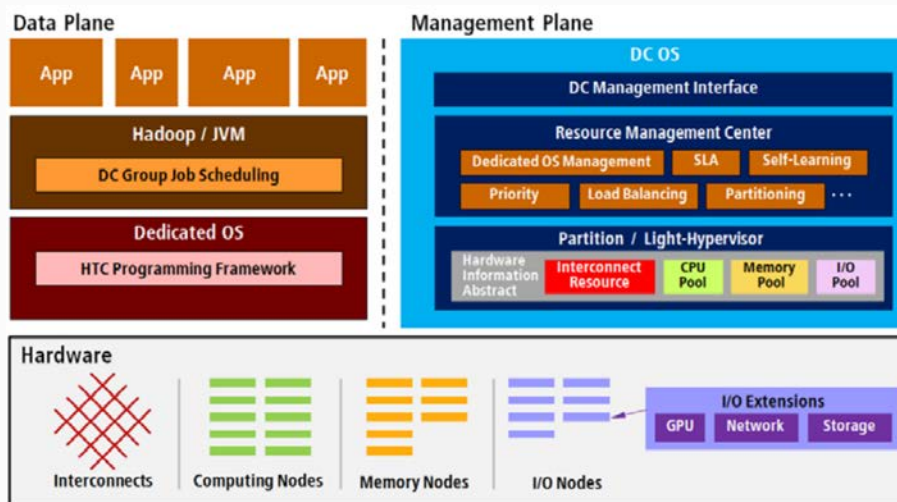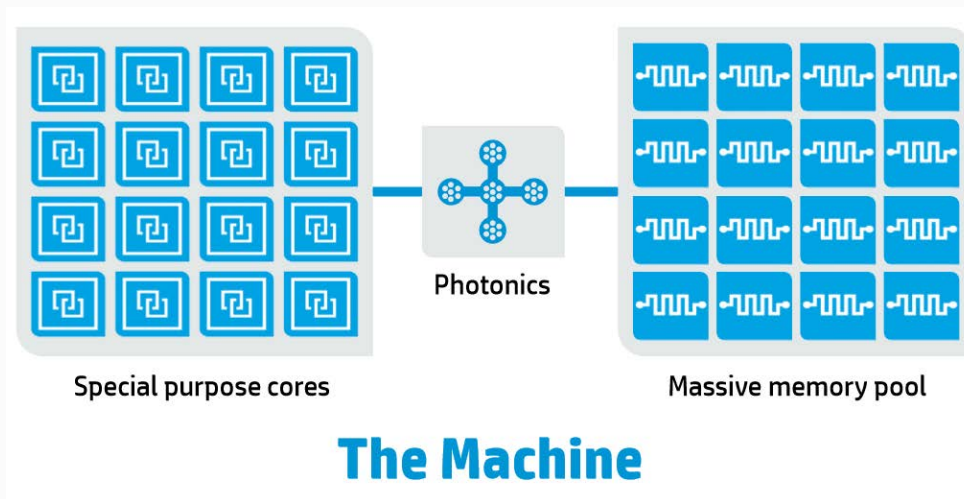**Specialization for the datacenter**
- → ~2000
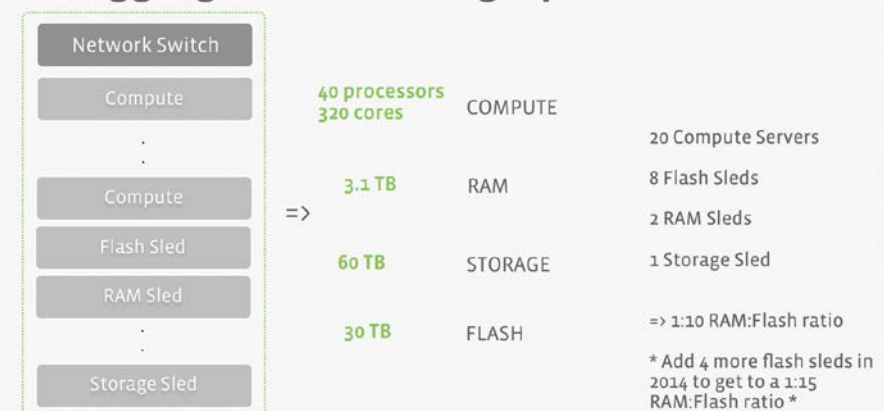  - Commercial Off-The-Shelf (COTS) computers, switches, & racks
- → ~2010
  - Custom computers, switches, & racks but build from COTS chips
- → ~2020
  - Custom computers, switches, & racks using custom chips

# Hardware Specialization

## Compute

**Moved faster than expected**

GPGPU

Google TPU

Microsoft Catapult

Amazon F1

## Networks

**Si Photonics not quite here yet**

50 μm    Transmitter

**But NWs have improved rapidly**

INFINIBAND

100G

GENZ    RoCE™

## Storage and Memory

**New physics moved slower than expected**

**Huge growth in flash**

## New Software Models

→ Service-Oriented Architecture (SOA) - All components are designed to be services

- Communicate only via the network (cluster-first mentality)
- Services reusable and interoperable

→ Serverless – Event-driven stateless functions

## Cluster-wide management

→ "Datacenter OS" (Mesos, Kubernetes, etc…)

→ Local OS less powerful (cluster-first mentality)

# Challenges

## Tail Latency
→At scale, slow == failed
→Need techniques to get good latency from unpredictable parts
→New Metric: 99% tail (average much less important)
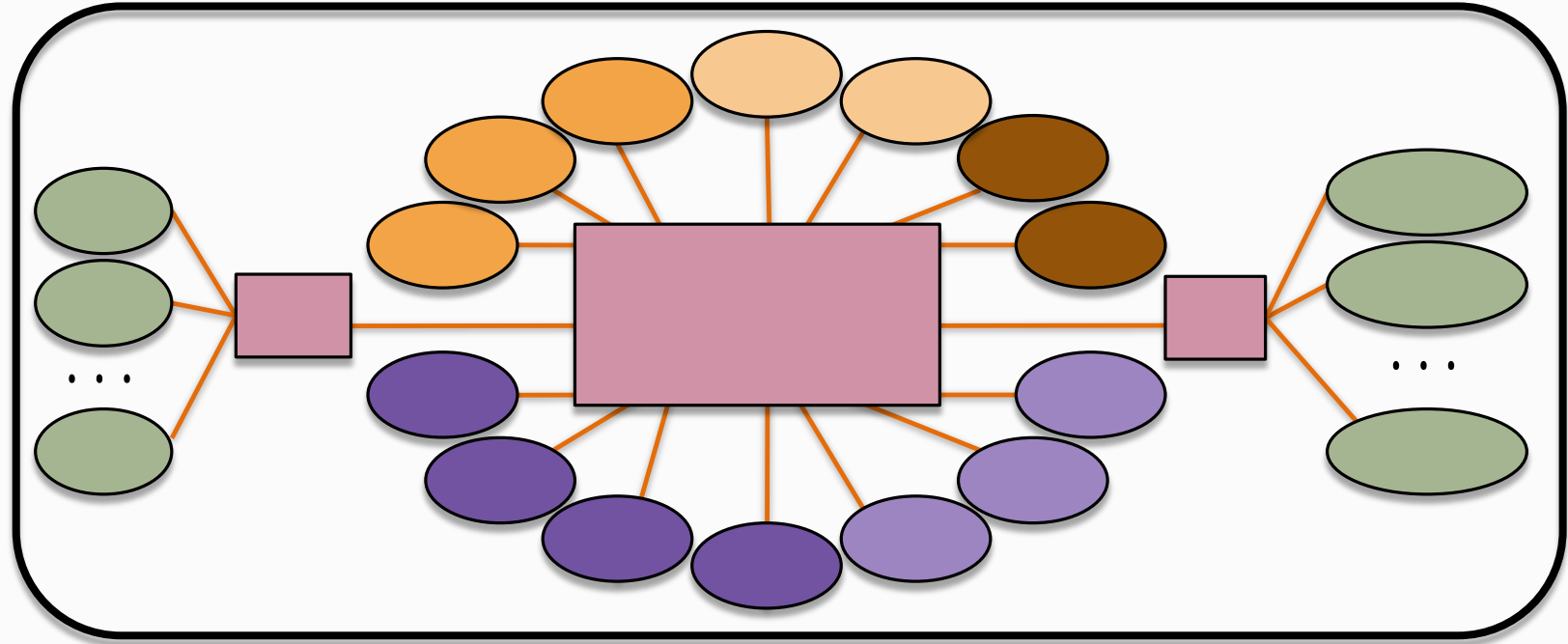
## Complex Memory Hierarchies
→Deeper: HBM, DRAM, NVM, Flash, Disk, Tape
→Wider: NUMA, RDMA, Disaggregation

## Security
→No longer considered safe within WSC
→Encryption demanded at all times

# Firebox at Berkeley

**What's Happened?**
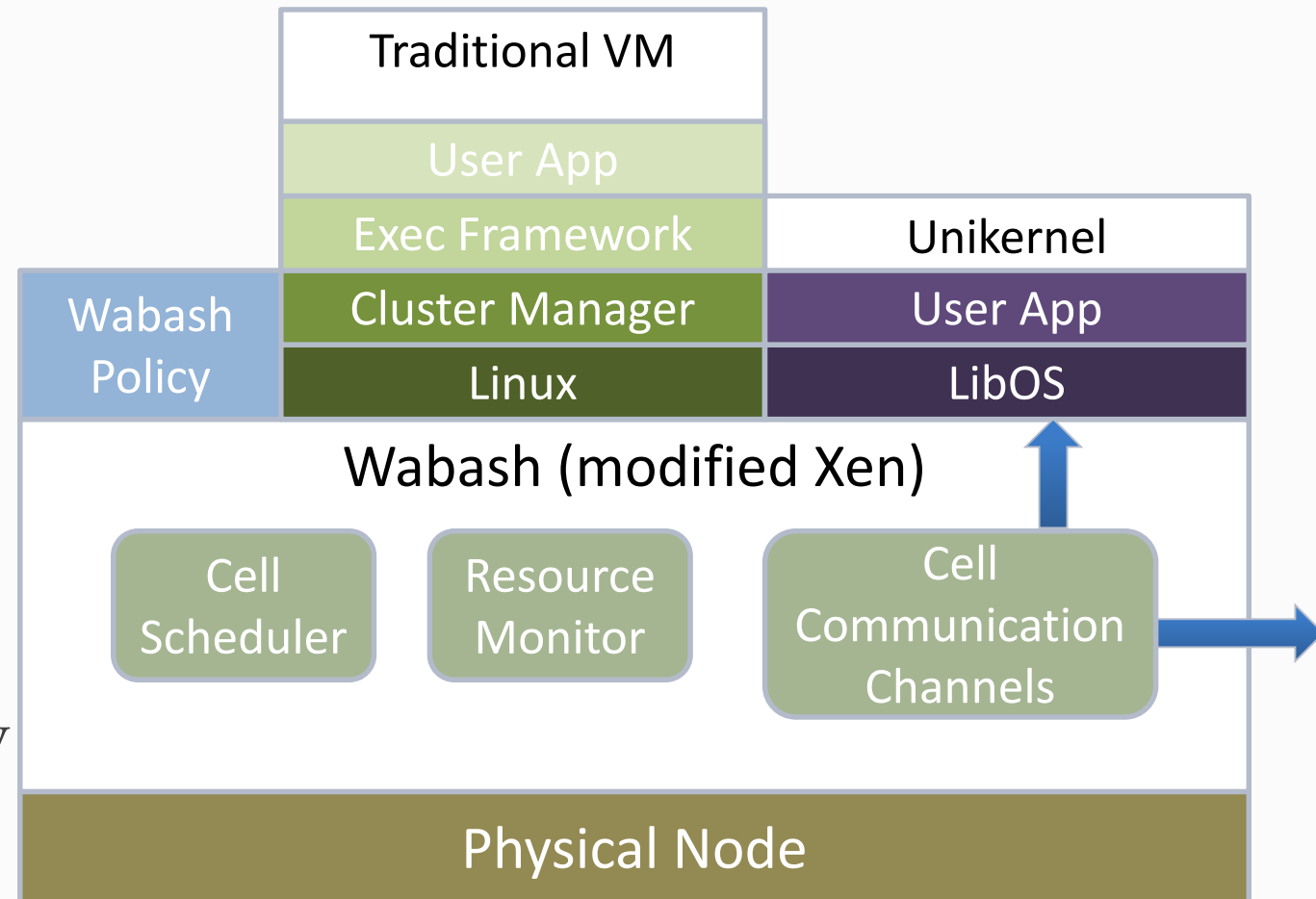
# Systems and Runtimes
## Wabash OS

**Project Goals and Approach:**

→ Push "DCOS" concept further

→ Based on concepts from Tessellation

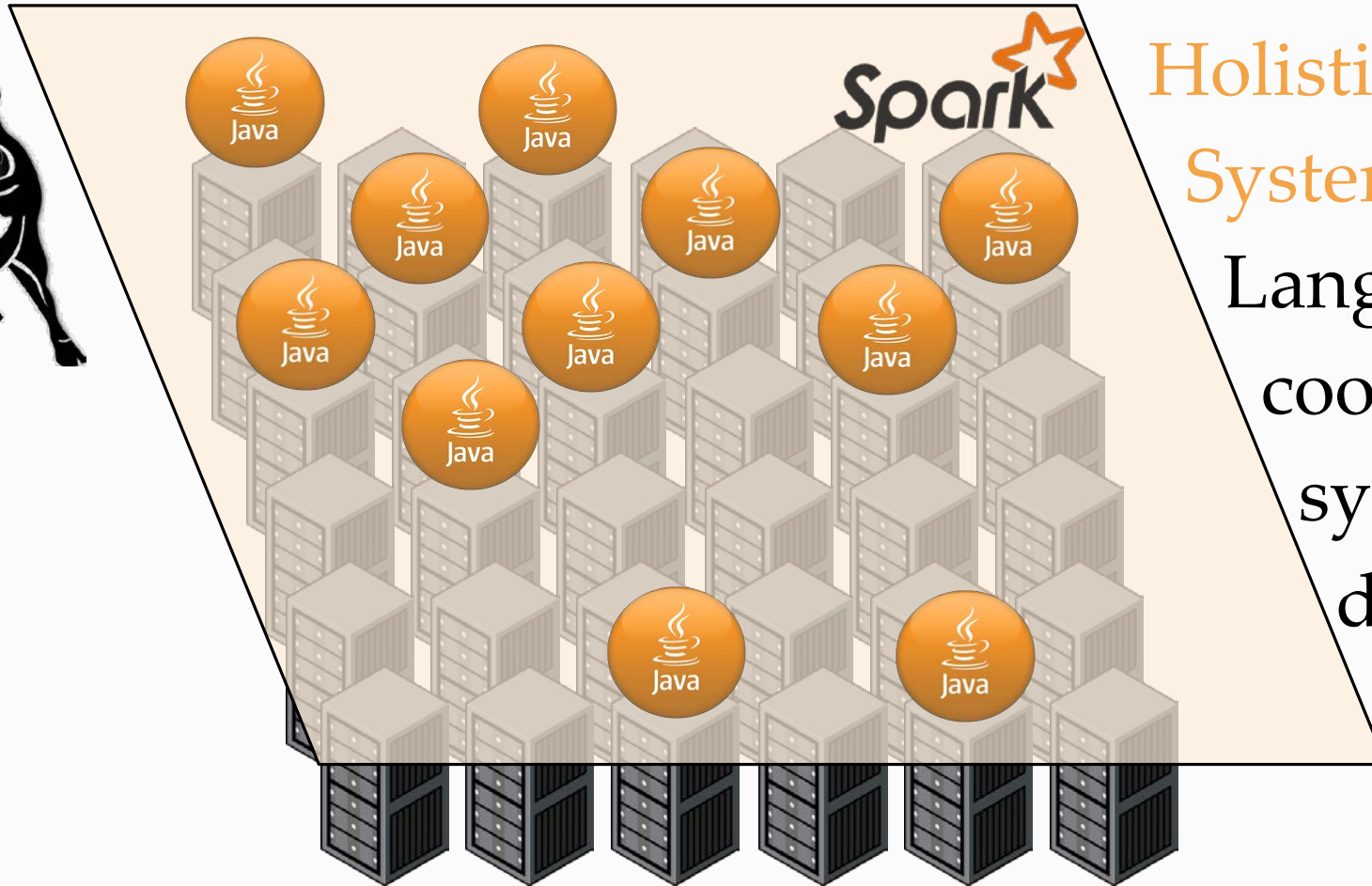→ Explored OS noise (and other overheads)

**Conclusions**

→ "Embrace the Noise"

- Better to be tail-tolerant than perfectly predictable

→ Path to performance without new OS

- New kernel-bypass techniques

- Self-virtualizing hardware (SR-IOV)



Traditional VM

User App

Exec Framework

Unikernel

Wabash Policy

Cluster Manager

User App

Linux

LibOS

Wabash (modified Xen)

Cell Scheduler

Resource Monitor

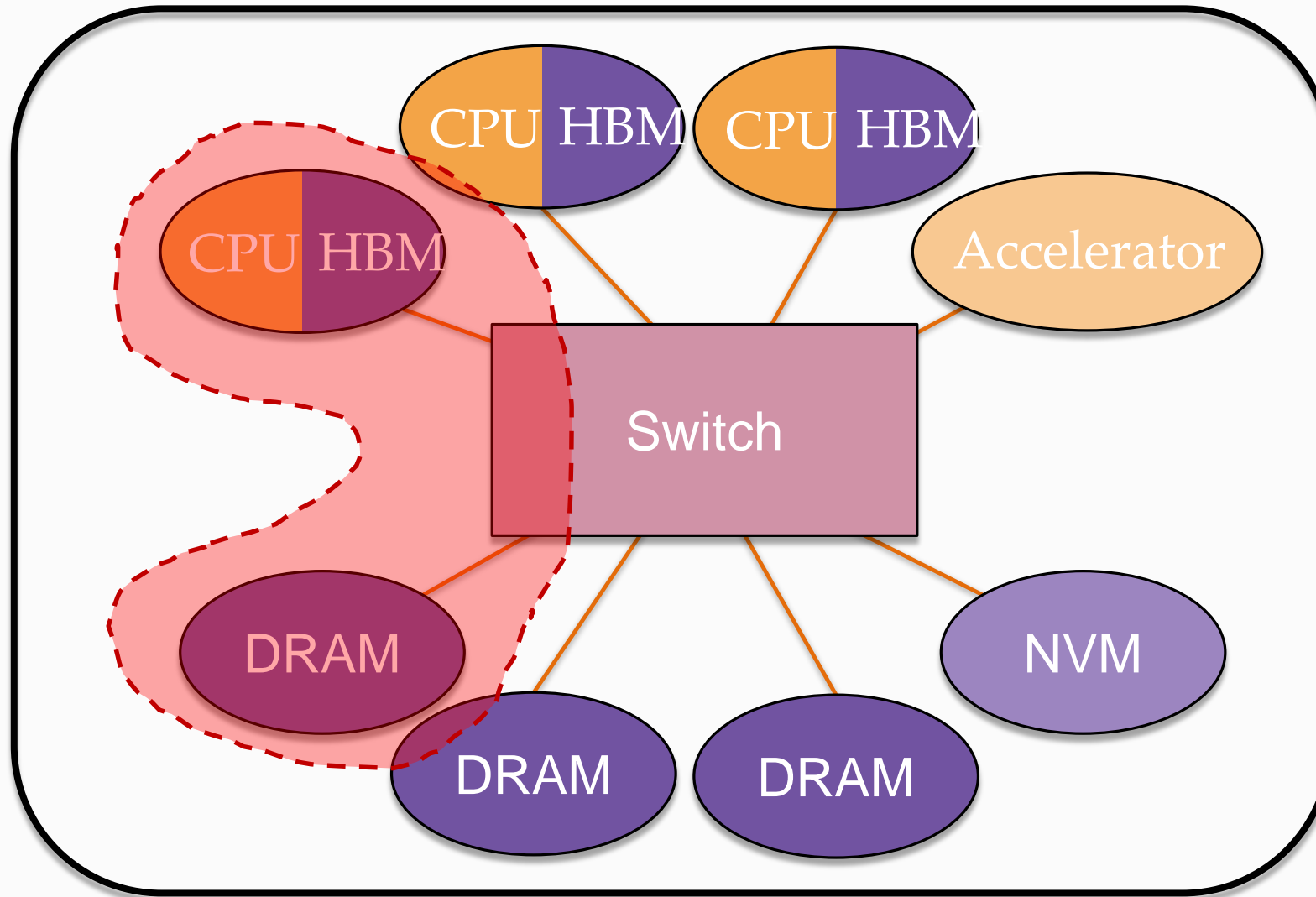Cell Communication Channels

Physical Node

## Taurus and the Holistic Runtime

**Holistic Language Runtime System**: A Distributed Language Runtime to coordinate the runtime systems underlying a distributed application
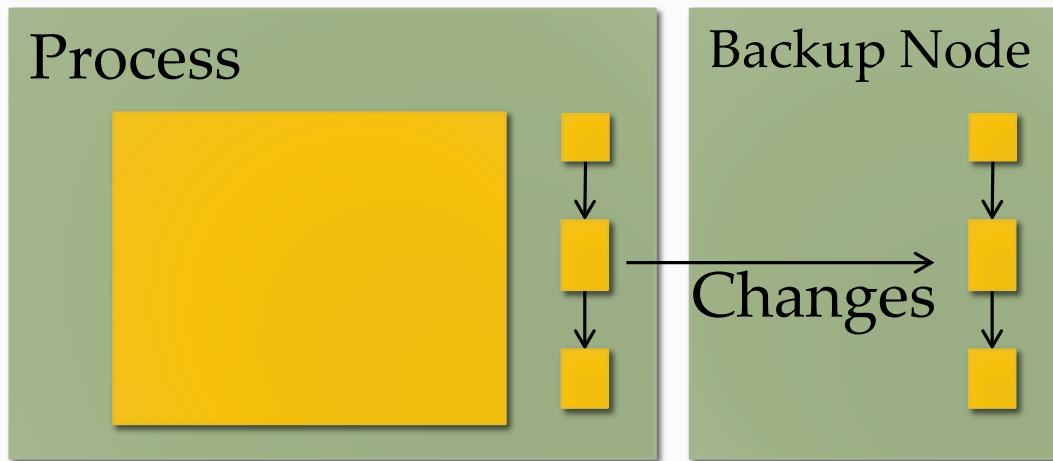
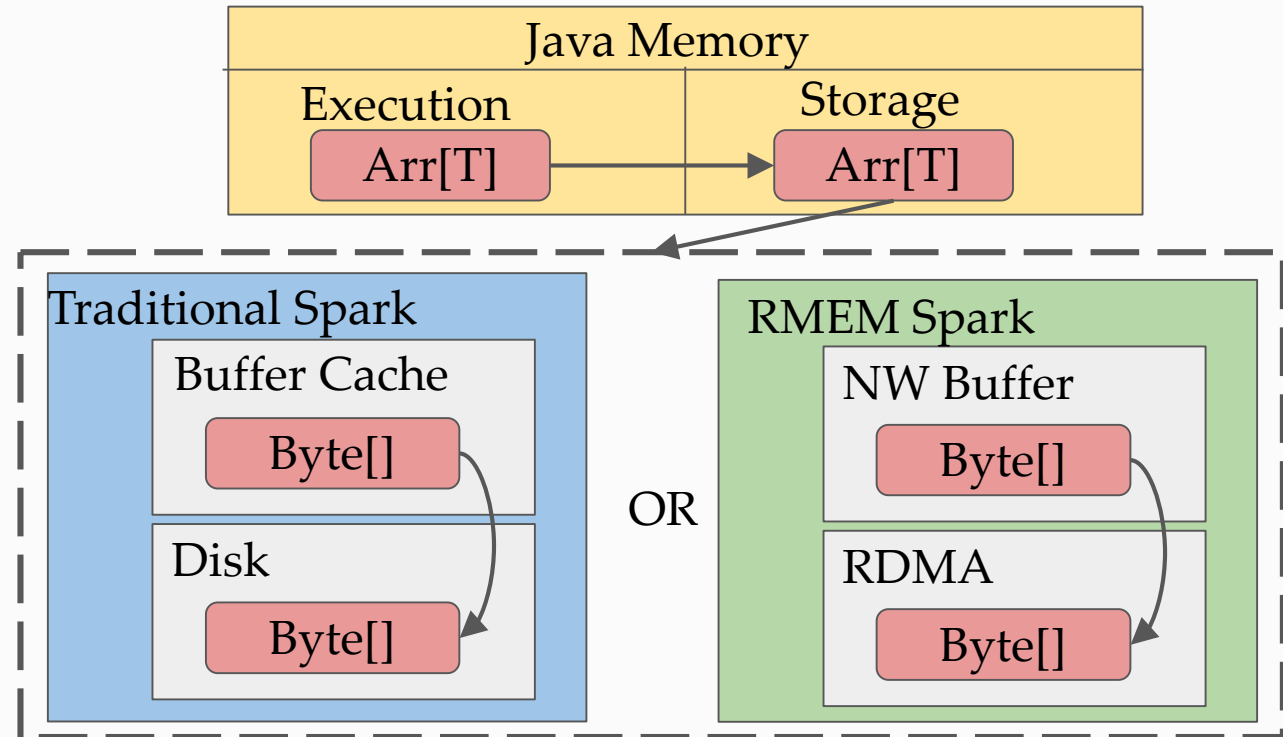## What is the right interface?

# Remote Memory
## Explicit Interfaces



**Nephele**
→ Transactional memory backup to NVM

**Spark over RDMA**
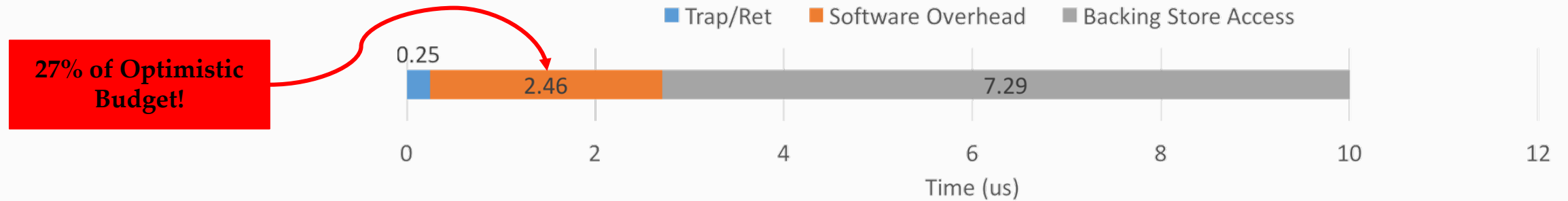→ Use remote memory to cache Spark objects

# Remote Memory
## Cache-Like Interfaces

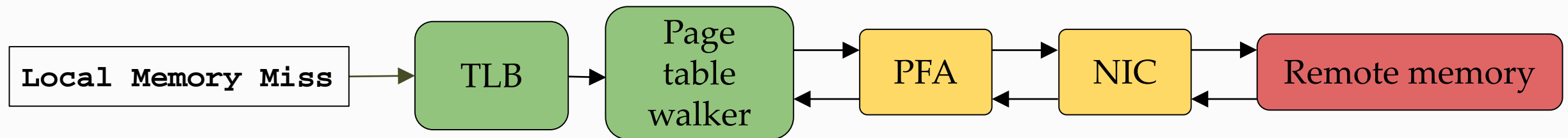## Network Requirements for Disaggregation
→ Found 100Gbps@3us may be sufficient for memory disaggregation
→ But! SW overheads are significant

**27% of Optimistic Budget!**

Legend: ■ Trap/Ret ■ Software Overhead ■ Backing Store Access

0.25 | 2.46 | 7.29

Time (us)

## Page-Fault Accelerator
→ Handle page-faults in HW
→ Manage evictions asynchronously in OS

```
Local Memory Miss → TLB → Page table walker ⇄ PFA ⇄ NIC ⇄ Remote memory
```

# The Future of Firebox

**Biggest Challenge so far: Lack of scalable and flexible experimentation platform**
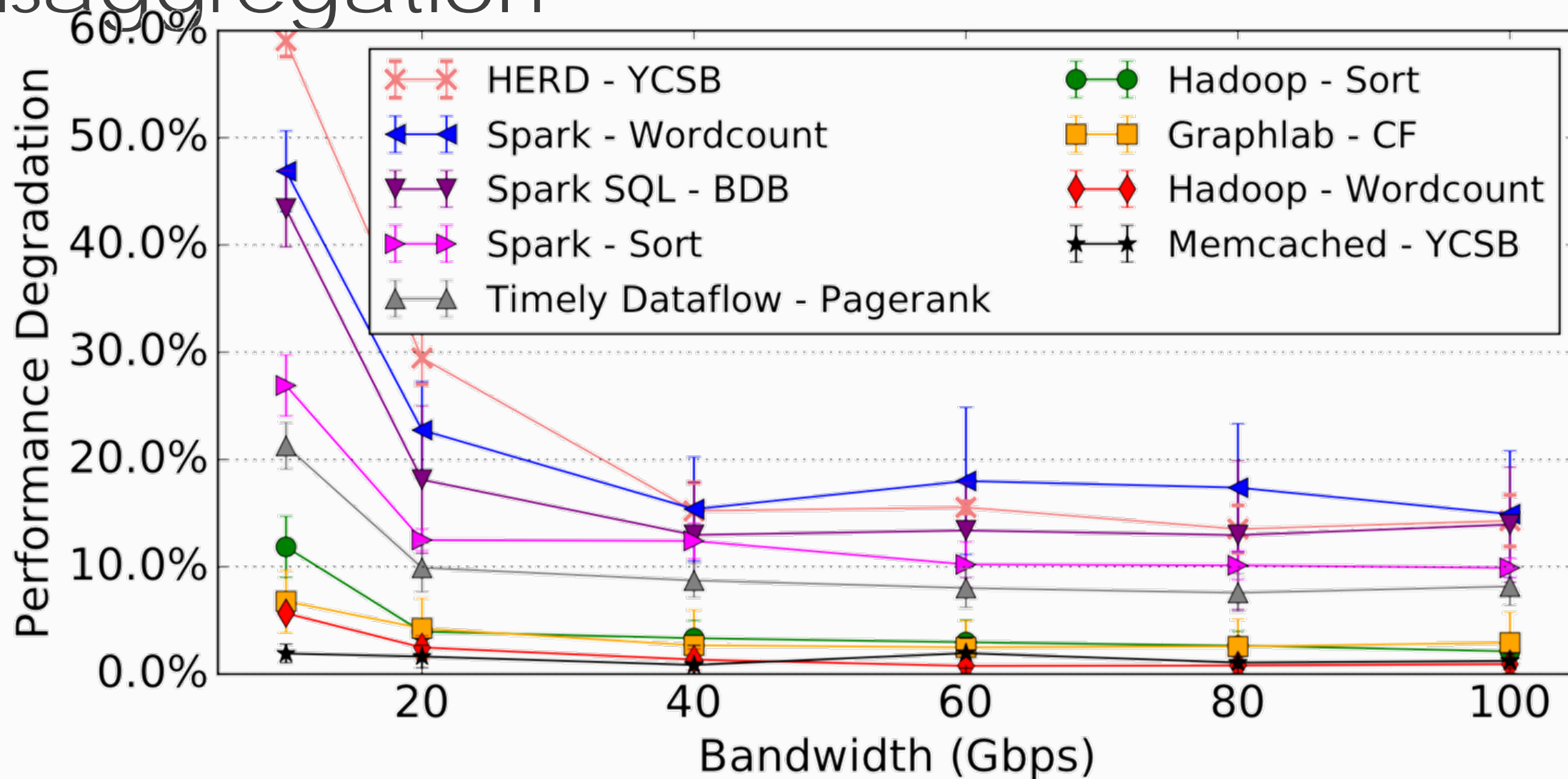
**What if we could simulate thousands of nodes with arbitrary RTL?**

→ Evaluate impact of swapping on NW

→ Place accelerators at network endpoints ("bump in the wire")

→ Experiment with memory blade design (compute-in-memory?)

→ Extend hardware enclaves across NW and heterogeneous devices

→ Gather detailed traces without impacting performance

→ …

# BACKUP AND REFERENCE SLIDES

Berkeley
UNIVERSITY OF CALIFORNIA

# Network Requirements for Disaggregation

# Page-Fault Accelerator