# Disaggregated Operating System
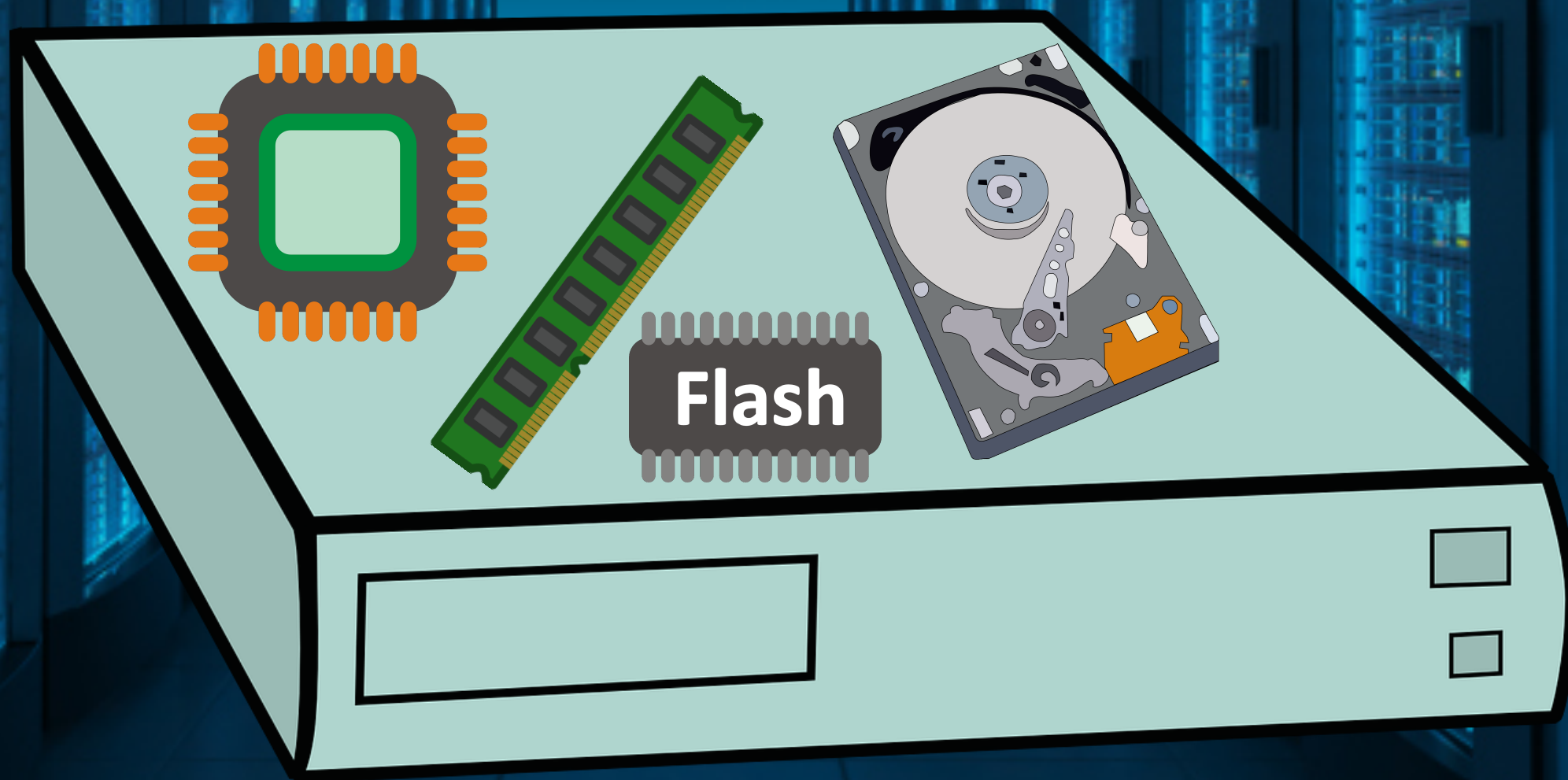
*Yiying Zhang*

*Yizhou Shan, Yilun Chen, Yutong Huang, Sumukh Hallymysore*
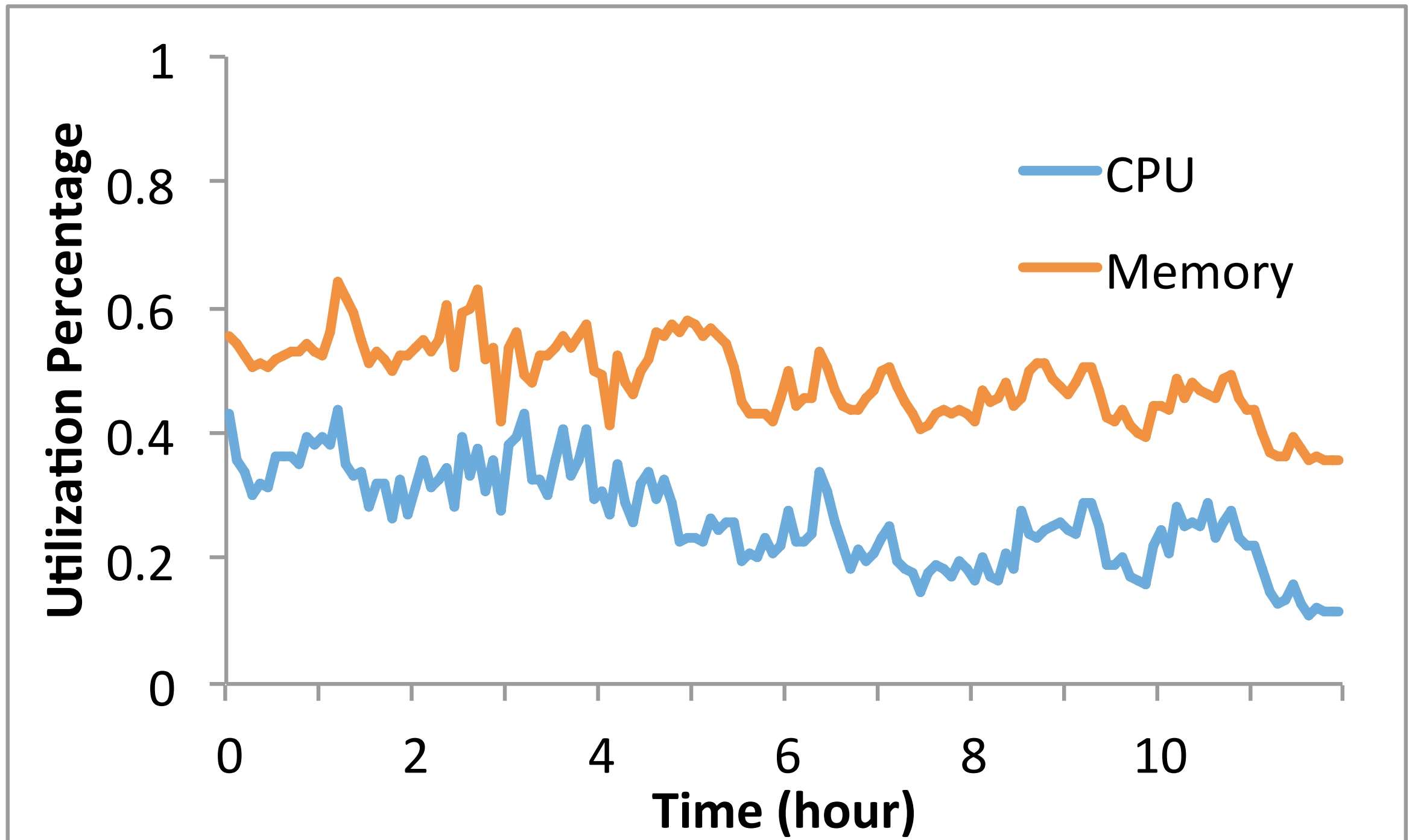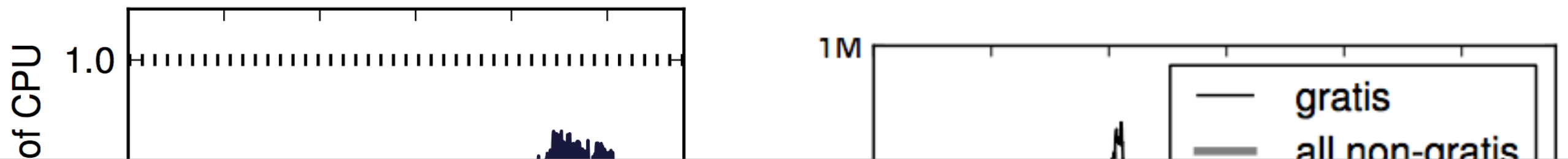
# Monolithic Server

- Resource utilization

- Failure

- Flexibility

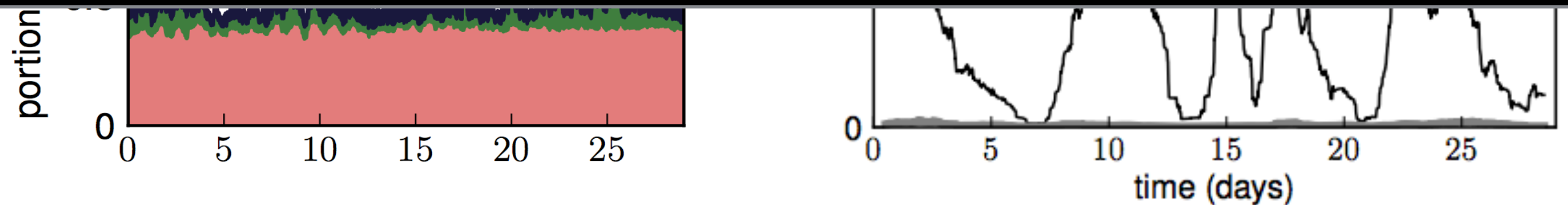# Alibaba Cluster Resource Utilization

# Google Cluster Resource Utilization



Resource can't be efficiently utilized

*REISS, C. etal. Heterogeneity and dynamicity of clouds at scale: Google trace analysis. SoCC'12*

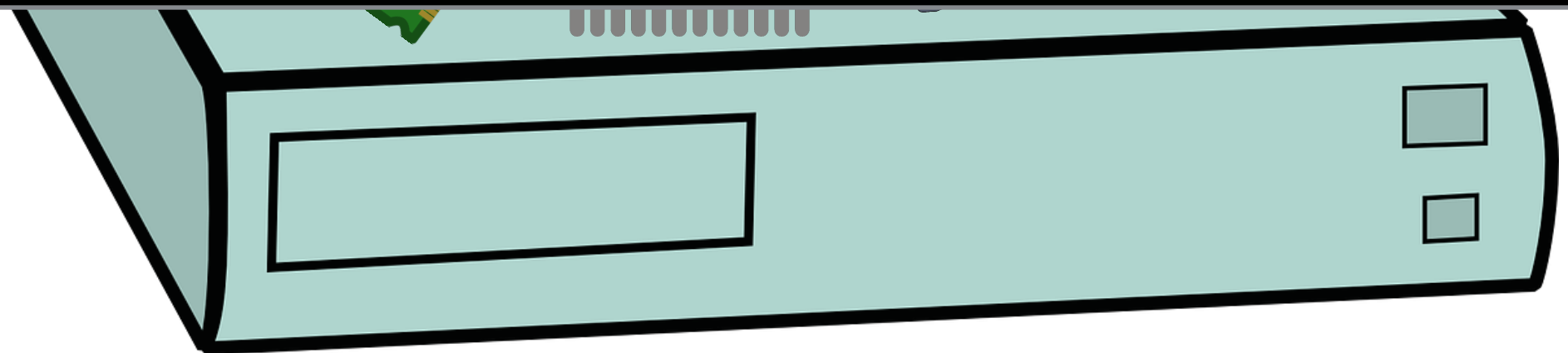No fine-grained failure handling

**App** **App** **App** **App**

**OS/Hypervisor**

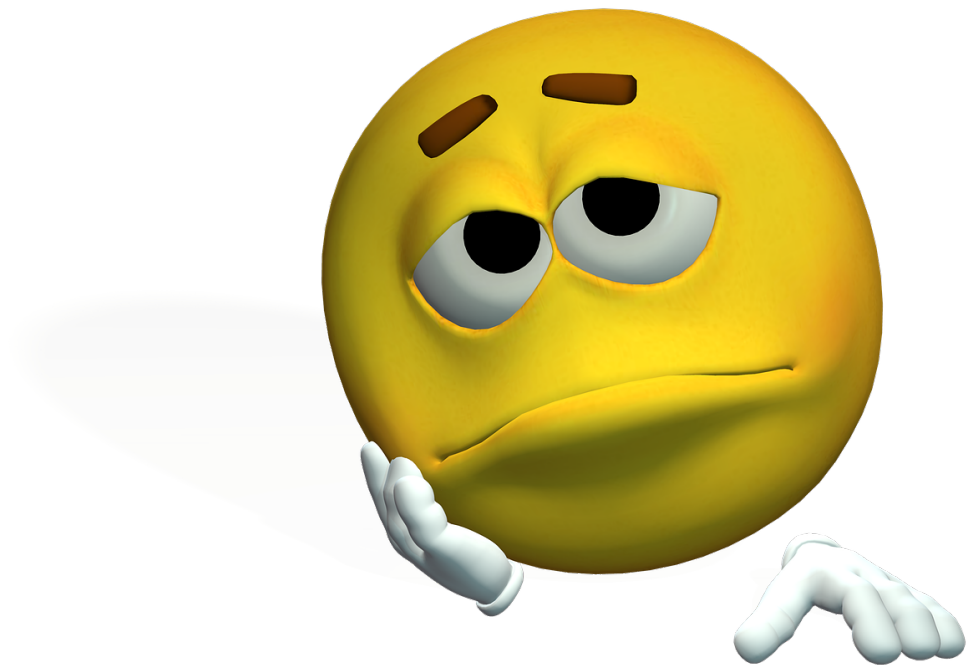Difficult to incorporate new hardware

# **Monolithic Server**

- Resource utilization

- Failure
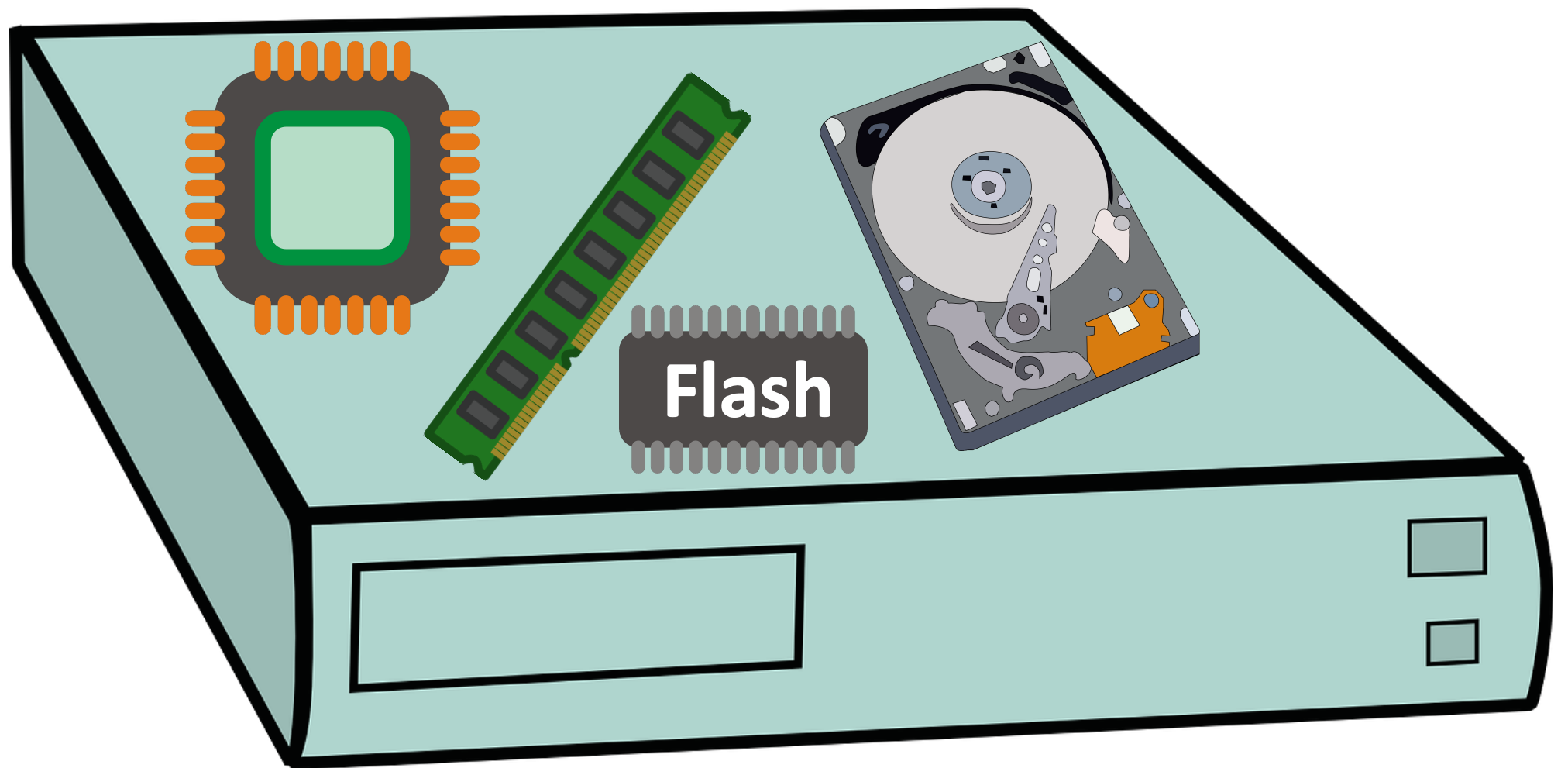
- Flexibility

- Memory capacity wall

When was the last time
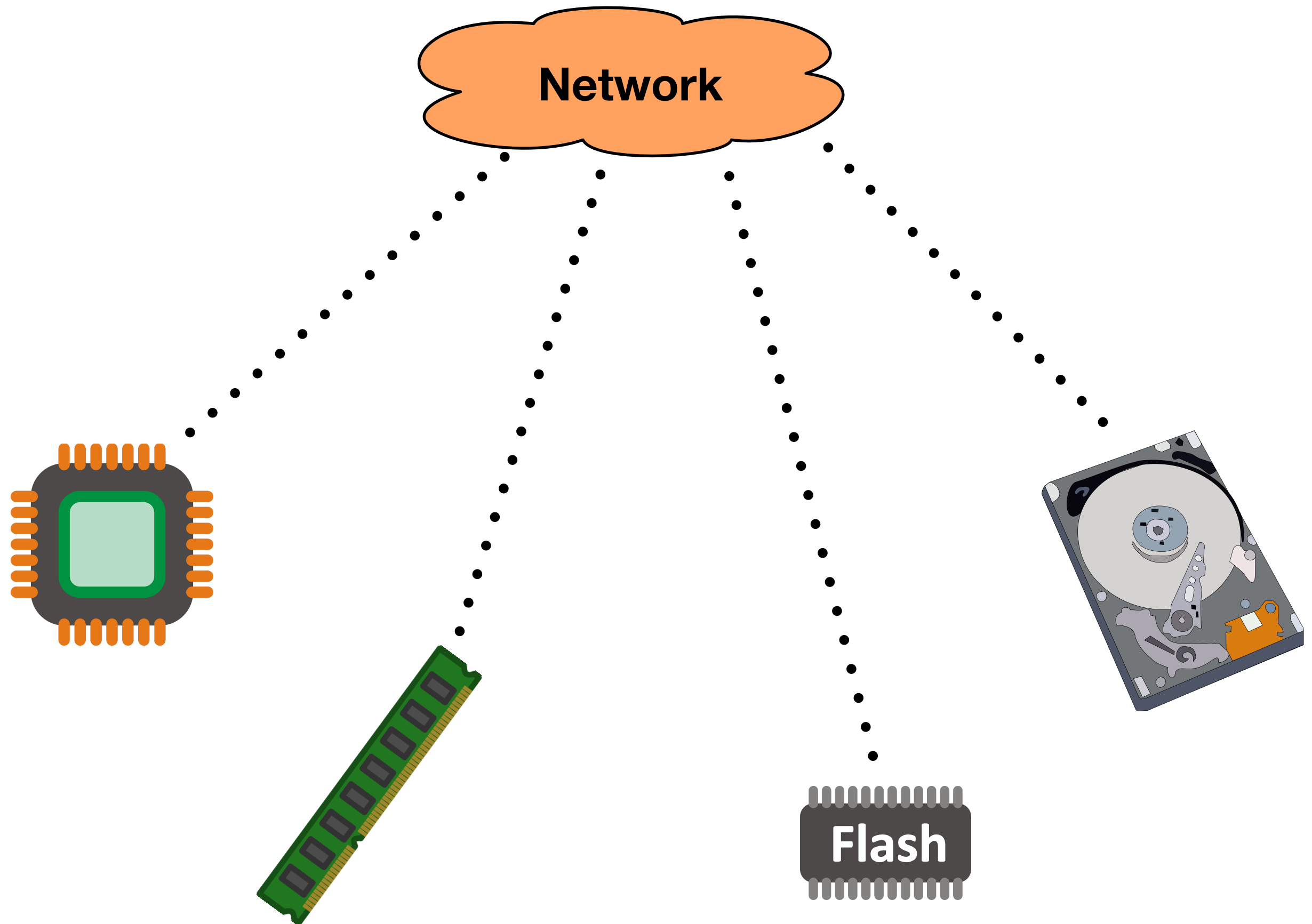
you did something for the first time?

# *Resource Disaggregation*:

**Breaking monolithic servers into <span style="color:red">network-attached</span>, <span style="color:red">independent</span> hardware components**

Flash

Network

Flash

# Gen-Z Consortium Formed: Developing a New Memory Interconnect

by Ian Cut

Posted in SoC



Gen-Z

High Bandw Low Lat

Advanc Worklo & Technol

Secu Compat Econom

10/11/2016

# HP Enterpris a single-me of addressin

DEAN TAKAHASHI     @DEANTAK

ME

HPE's vision for performan

Above: HPE's new Memory-Driven

Image Credit: HPE

dRedBox.eu demonstrates its progress in materializing its vision towards fully disaggregated datacenters and cloud.

Box
dRed

dReDBox
Disaggregated Memory
Prototype

dReDBox
Disaggregation OS

dReDBox
Server
Prototype

n materializing disaggregated datacenters

**Oct 2nd 2017**

# **Why Now?**

- Faster network

10,000

Gb/s

# Mellanox: We're gonna make InfiniBand great again – 200Gbps great

So great, offload as much as possible from CPUs,

## ConnectX®-6 Single/Dual-Port Adapter Supporting 200Gb/s Ethernet
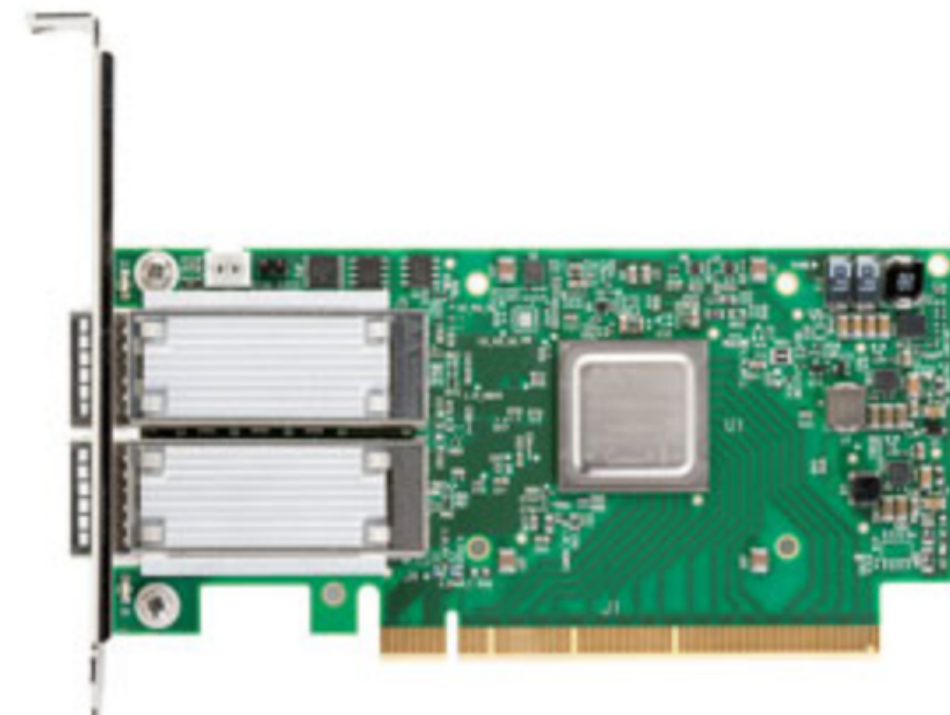
✉ Contact Sales for Availability

Intelligent ConnectX-6 adapter cards, the newest additions to the Mellanox Smart Interconnect suite and supporting Co-Design and In-Network Compute, introduce new acceleration engines for maximizing Cloud, Web 2.0, Big Data, Storage and Machine Learning applications.

ConnectX-6 EN supports two ports of 200Gb/s Ethernet connectivity, sub-600 nanosecond latency, and 200 million messages per second, providing the highest performance and most flexible solution for the most demanding applications and markets.

ConnectX-6 offers Mellanox Accelerated Switching And Packet Processing (ASAP2) Direct technology to offload the vSwitch/vRouter by handling the data plane in the NIC hardware while maintaining the central plane unmodified. As a result, significantly higher

90ns and aggregate capacity is 16Tbps.

ConnectX·6

# Why Now?

- Faster network

- More powerful hardware controller


- Dynamic application resource requirement

- Quickly changing, heterogeneous hardware

# Resource Disaggregation

- Better resource utilization

- Fine-grained failure

- Heterogeneity

- Embracing hardware innovations

# Using Existing Kernels?

- Monolithic/micro kernel: built for single monolithic server

- Multikernel: (vertically) replicated kernel across cores

- Distributed OS [*Sprite*, *V*, *MOSIX*, *Charlotte*]:

  manages distributed monolithic servers

  *Amoeba*: manages resource pool, but not in modern days

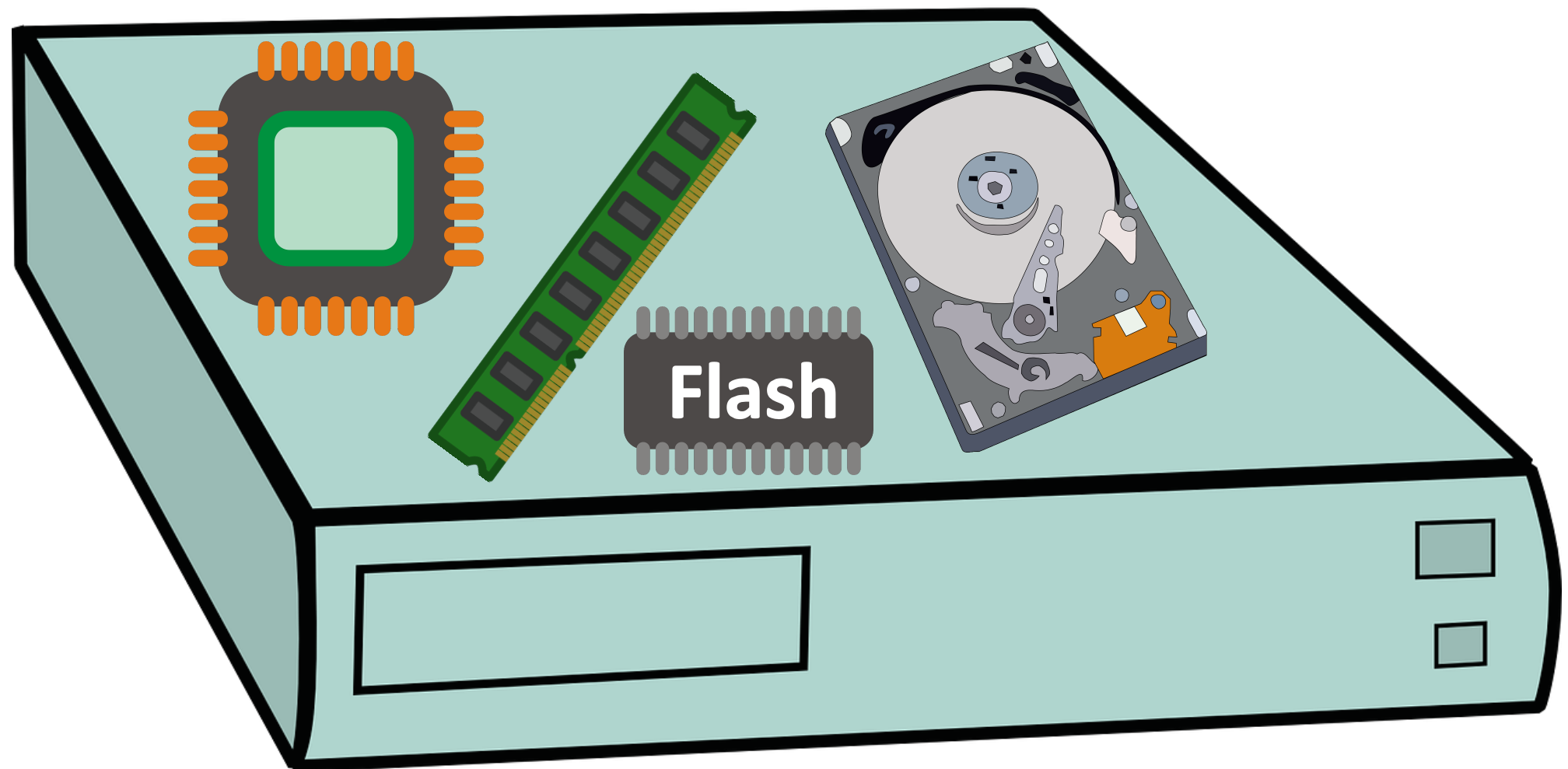# When hardware is disaggregated, the OS should be also!
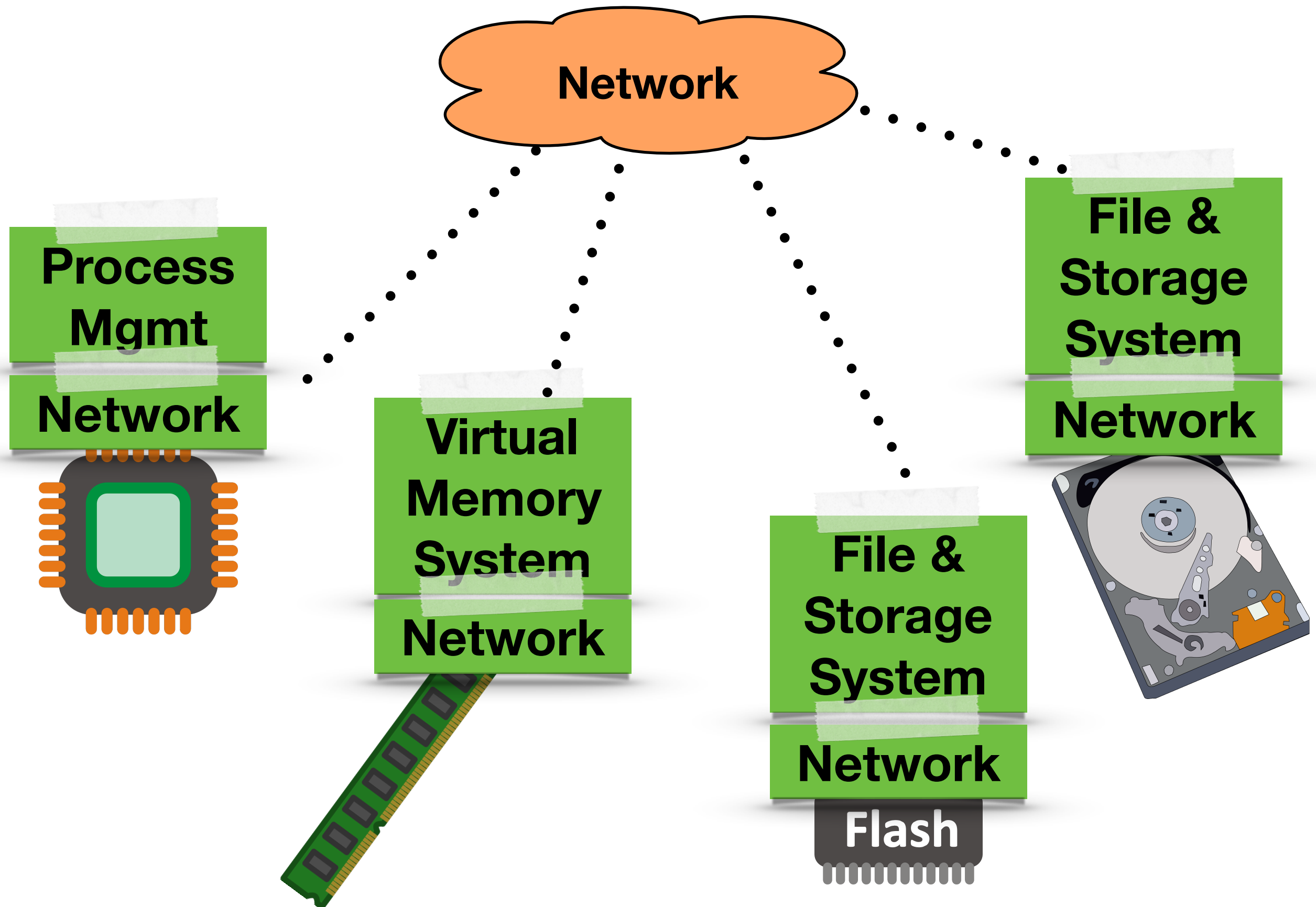
# Lego

Processor

Memory

Storage

NVM

# **Challenges**

- Cleanly separate OS services

  - *Stateless, minimal dependencies*

- Fit hardware constraints

  - Processor: no or limited local DRAM

  - Memory: limited processing power

User Space

**Application processes**

**Linux interface, state session**

Kernel Space

**Process/thread scheduling**

**L4 cache management**

**Process checkpointing**

**Process Mgmt**

**File & Storage System**

**Network**

**Network**

**Virtual Memory System**

*Virtual Addr*

*Virtual Addr*

*Virtual Addr*

*Virtual Addr*

| Core | Core | Core |
|------|------|------|
| L1 $ | L1 $ | L1 $ |
| L2 $ | L2 $ | L2 $ |

**L3 $**

**Storage**

**L4 DRAM Cache**

• Coarse $ line

• High associativity

• Software managed

L1 to L4 all
*Virtual Cache:*
virtually indexd
virtually tagged

- No globally shared memory

- No coherence traffic across network

- RDMA-based RPC

**Network**

**Process Mgmt**

Kernel Space
(running in mem controller)

**Network**

**Virtual Memory System**

**Network**

Virtual address (de)allocation

Physical address (de)allocation

Memory-mapped file mgmt

Memory replication

**System**

**Network**

**Memory Controller**

**Cached Index (e.g., TLB)**

**DRAM**

**Virtual to Physical Mapping**

# Challenges

- Cleanly separate OS services

- Fit hardware constraints

- Handle failures

- Global resource management
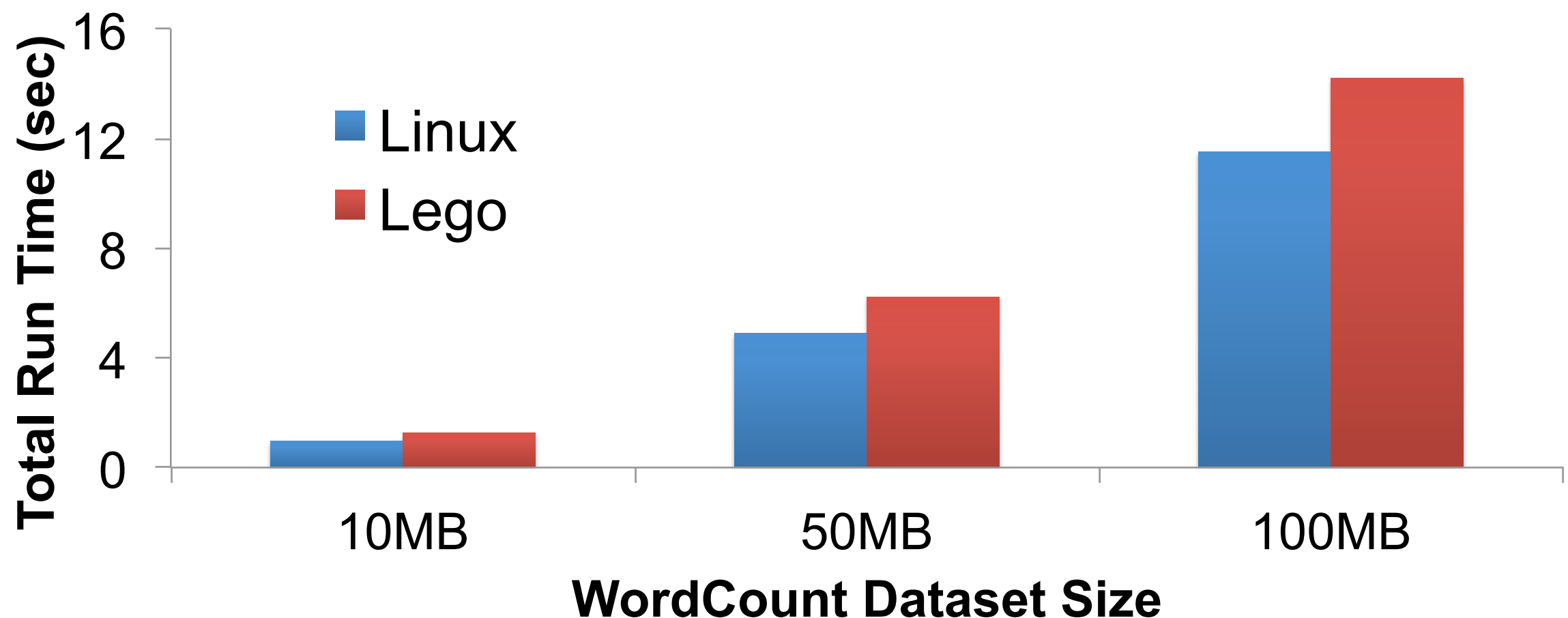
# Status Report

- 170K LOC so far

- Simple processor, memory, storage managers

- Support X86-64

- Backward compatible with common Linux interface

- Run unmodified datacenter application binaries

- Emulate hardware devices using commodity servers

*We will open source!*

# Initial Results are Encouraging

Phoenix (single-node MapReduce), unmodified statically-linked binary
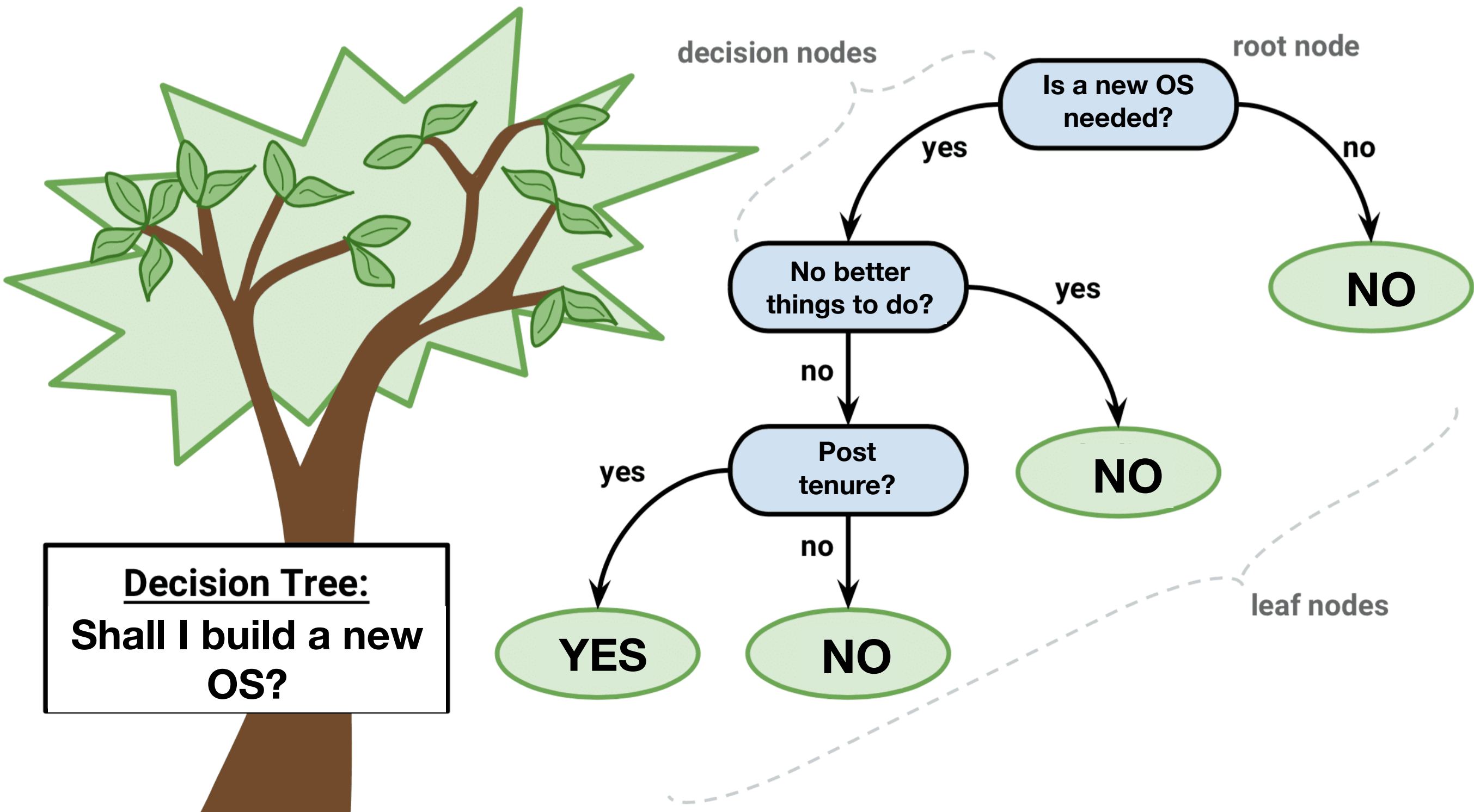
Compare one commodity server running Linux with Lego running on one proc, one mem, one storage, emulated using three servers

# Conclusion - A Bunch of Questions

- Time to change datacenters?

- Do you believe in resource disaggregation?

- New OS for new hardware?

- Are we reinventing the wheel?

- Killer applications?

# Conclusion [hidden version]



decision nodes

root node

**Is a new OS needed?**

yes

no

**No better things to do?**

yes

no

**NO**

**Post tenure?**

yes

no

**NO**

yes

no

**YES**

**NO**

leaf nodes

**Decision Tree:**
**Shall I build a new OS?**

# Thank You
# Questions?



wuklab.io