# WORKLOAD CHARACTERIZATION OF INTERACTIVE CLOUD SERVICES ON BIG AND SMALL SERVER PLATFORMS

**Shuang Chen***, Shay Galon**, Christina Delimitrou*,
Srilatha Manne**, and José Martínez*

*Cornell University
**Cavium Inc.

- **How to achieve low tail latency for interactive cloud services?**
  - Tail latency more important and challenging
  - The entire stack from SW to HW is involved

- **Understand how tail latency reacts to application and system changes**
  - Quantify how current designs work
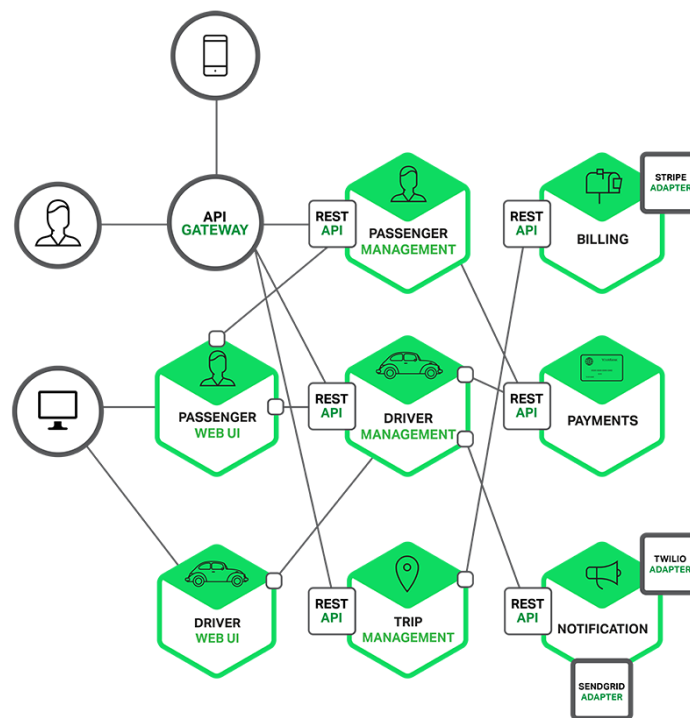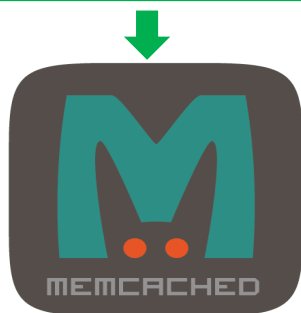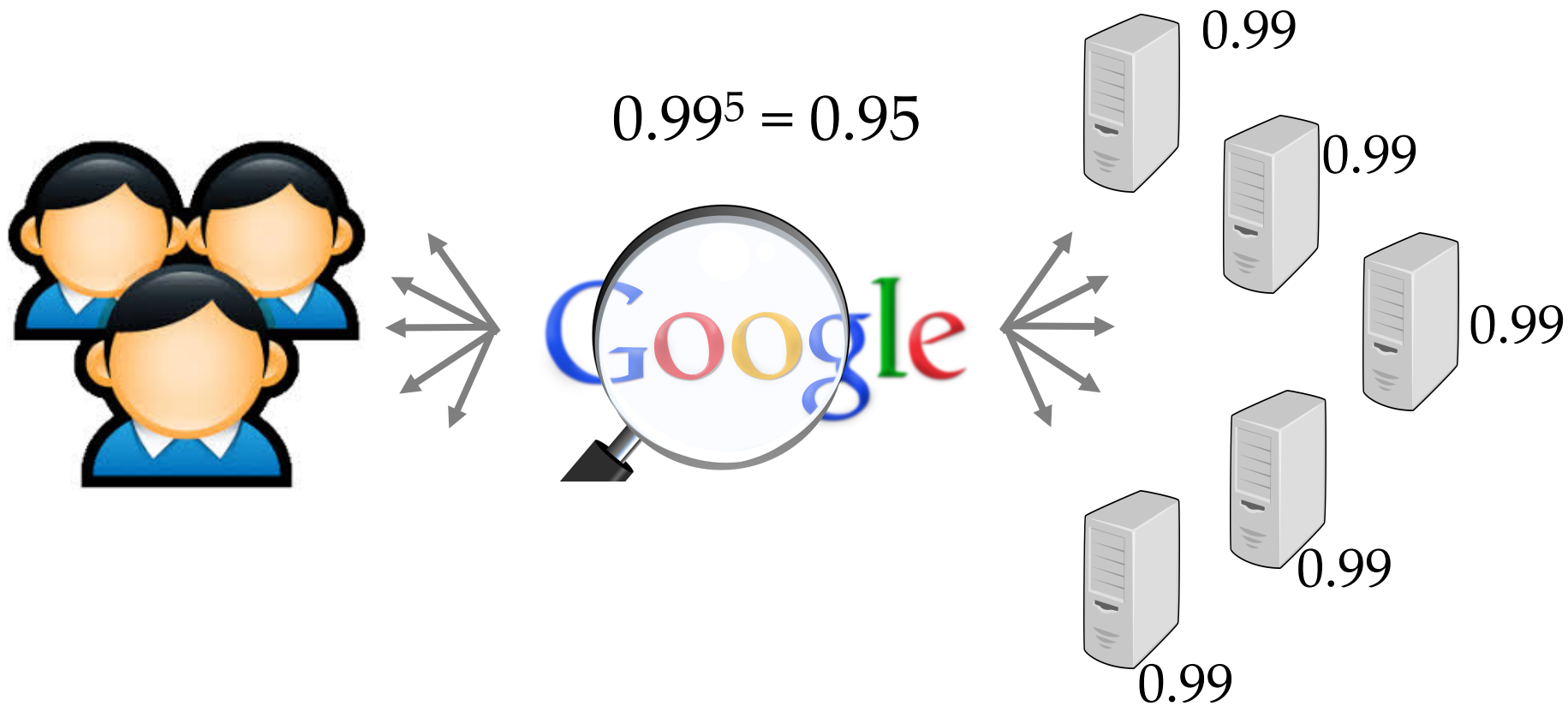  - Get insights on future designs

## ▪ Tail latency

- e.g., QoS defined as $99^{th}$ %ile in 1ms



$$0.99^5 = 0.95$$

0.99
0.99
0.99
0.99
0.99
0.99

▪ **The entire stack from SW to HW is involved**

Application

Resource Manager

Virtualization

OS

Hardware

- Application bottleneck
- Different user cases
- Scalability

- Overhead of virtualization
- SW isolation mechanisms
- Overhead of context switching

- HW isolation mechanisms
- Hyperthreading

## By QoS Strictness
- us:  memcached
- ms: web server, in-memory database
-  s:    persistent database

## By Statefulness
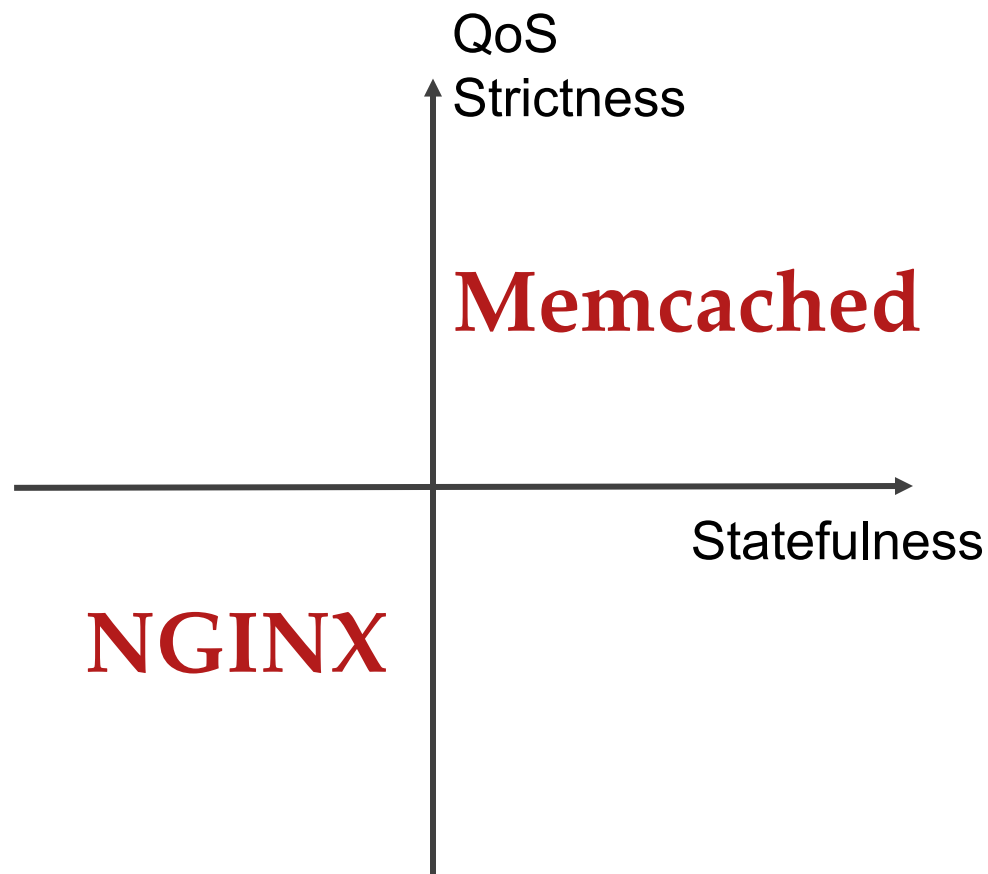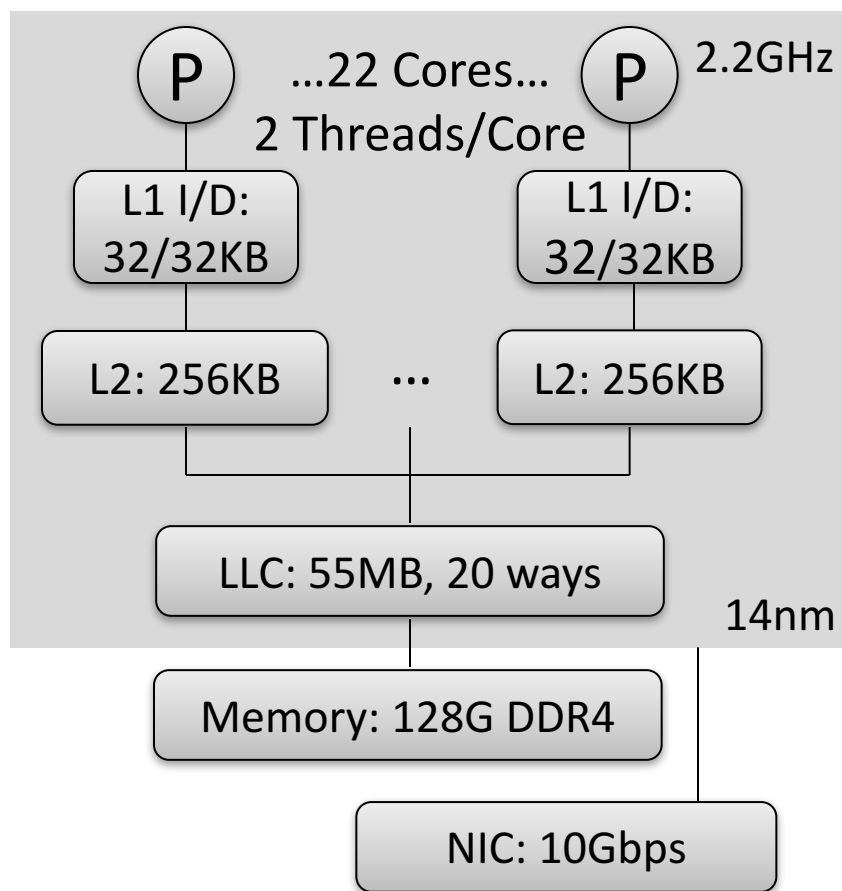- Stateful:  memcached
- Stateless: web server

## NGINX

- Web server
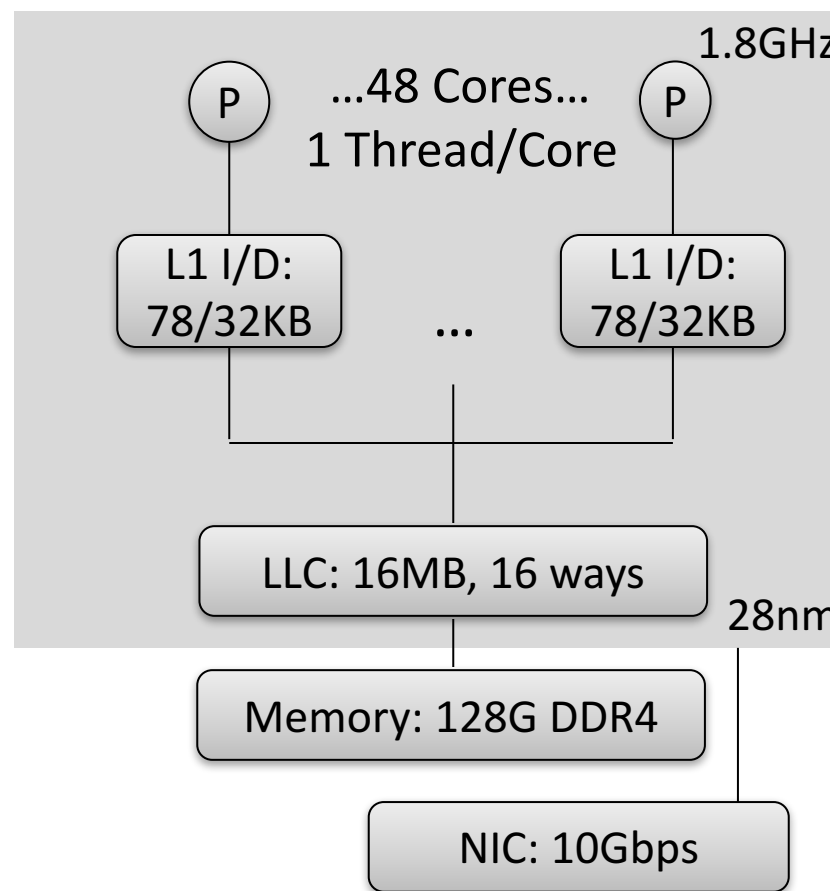- Stateless
- 99th% in tens of ms

## Memcached

- Key-value store
- Stateful
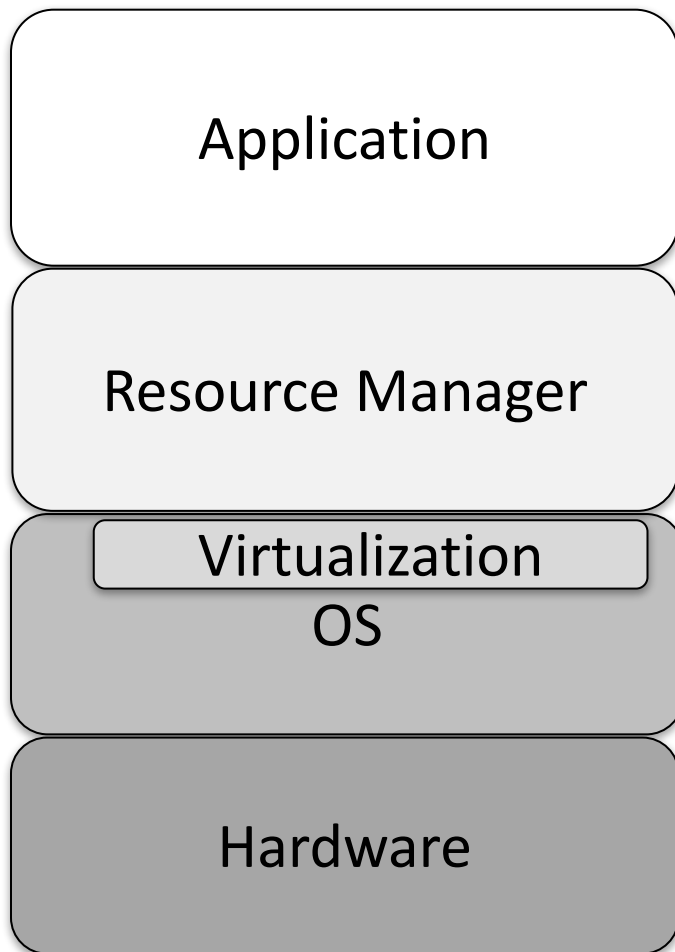- 99th% in hundreds of us

QoS Strictness

**Memcached**

Statefulness

**NGINX**
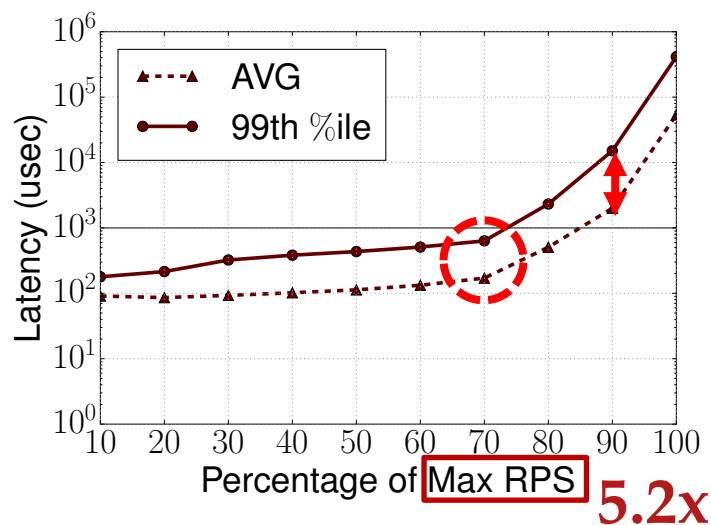
Intel Xeon E5-2699 v4
$4,115

Cavium ThunderX
$785

# STUDIED PARAMETERS



- **Application bottleneck**
- Different user cases
- Scalability

- Overhead of virtualization
- SW isolation mechanisms
- Overhead of context switching
- HW isolation mechanisms
- Hyperthreading

Cornell University
Computer Systems Laboratory

Xeon

ThunderX

Memcached

NGINX

5.2x

5x

# MEMCACHED LATENCY DECOMPOSITION



**Little user-space processing**

At **10%** of max throughput

Xeon    | 6 | 31 | 11 | 1 | 5 |

ThunderX | 14 | 4 | 5 | 9 | 7 | 24 |

**Network delay**

**2x slower than Xeon**

At **90%** of max throughput

Xeon    | 6 | 782 | 1,009 | 3 | 1 | 5 |

ThunderX | | 1,290 | 1,650 | | 7 | |

14                                                              20  24

**Queuing delay**

| Application |
|---|

| Resource Manager |
|---|

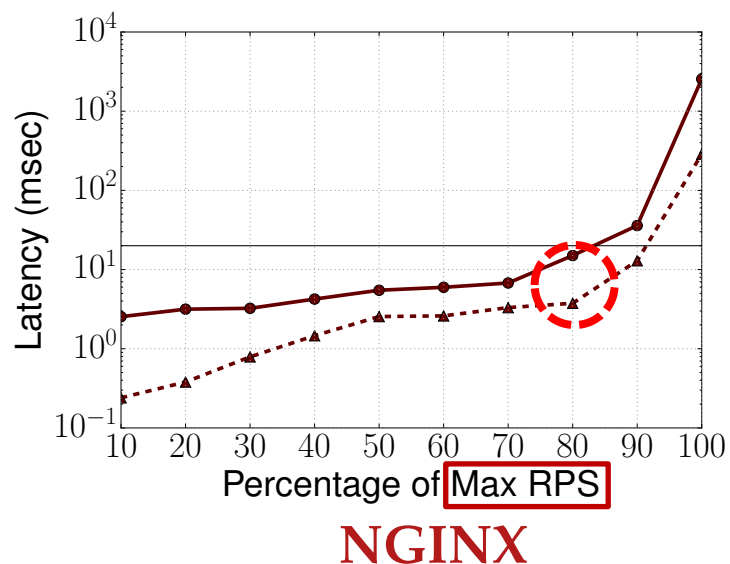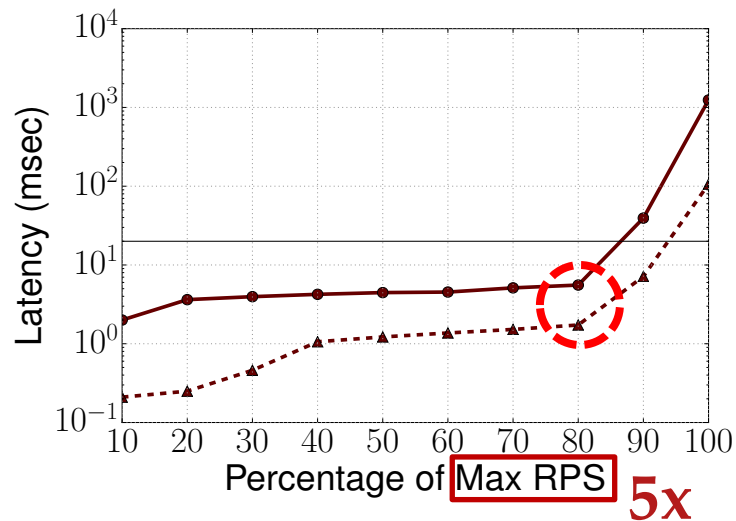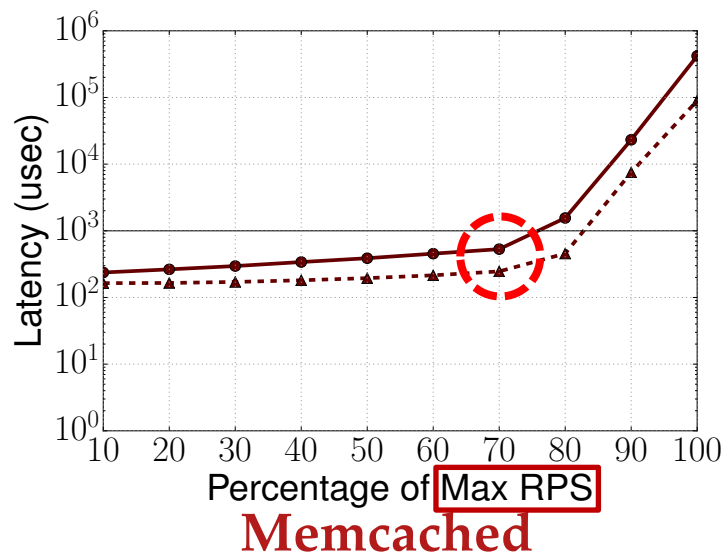| Virtualization |
|---|
| OS |

| Hardware |
|---|

- Application bottleneck
- **Different user cases**
- Scalability

- Overhead of virtualization
- SW isolation mechanisms
- Overhead of context switching
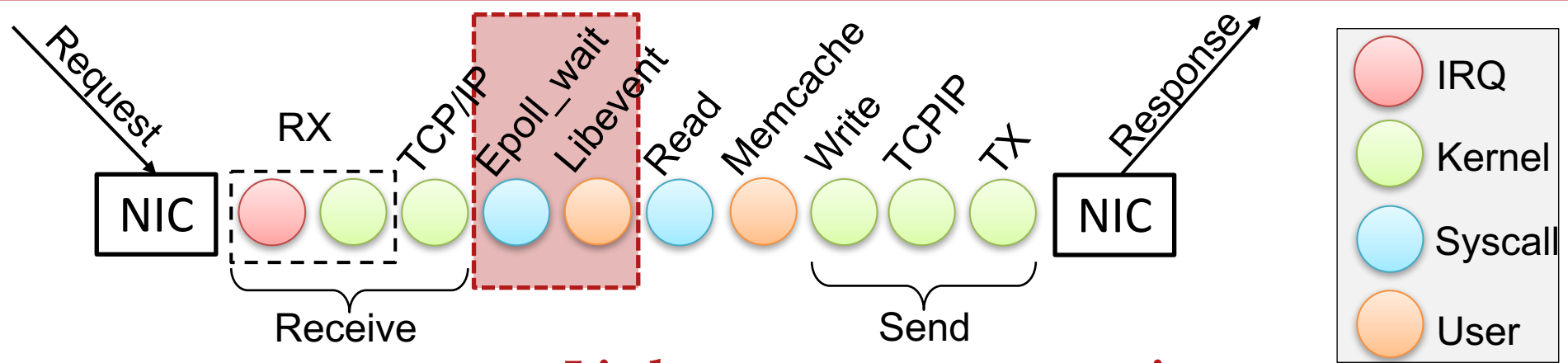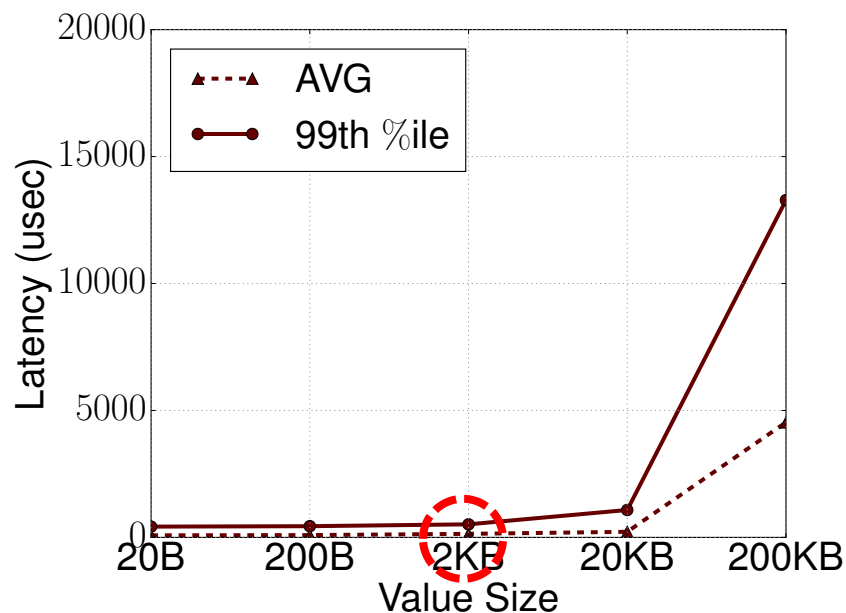- HW isolation mechanisms
- Hyperthreading

Xeon

ThunderX

- Memory copy
- Network processing and transmission
- ThunderX is more sensitive

CSL



**Xeon**

**ThunderX**

- Cache capacity
- ThunderX is more sensitive

| | |
|---|---|
| **Application** | • Application bottleneck<br>• Different user cases<br>• **Scalability** |
| **Resource Manager** | |
| **Virtualization**<br>**OS** | • Overhead of virtualization<br>• SW isolation mechanisms<br>• Overhead of context switching |
| **Hardware** | • HW isolation mechanisms<br>• Hyperthreading |

## Memcached

## NGINX



- Interrupt handling
- Load imbalance
- Lock contention

# STUDIED PARAMETERS

**CSL**

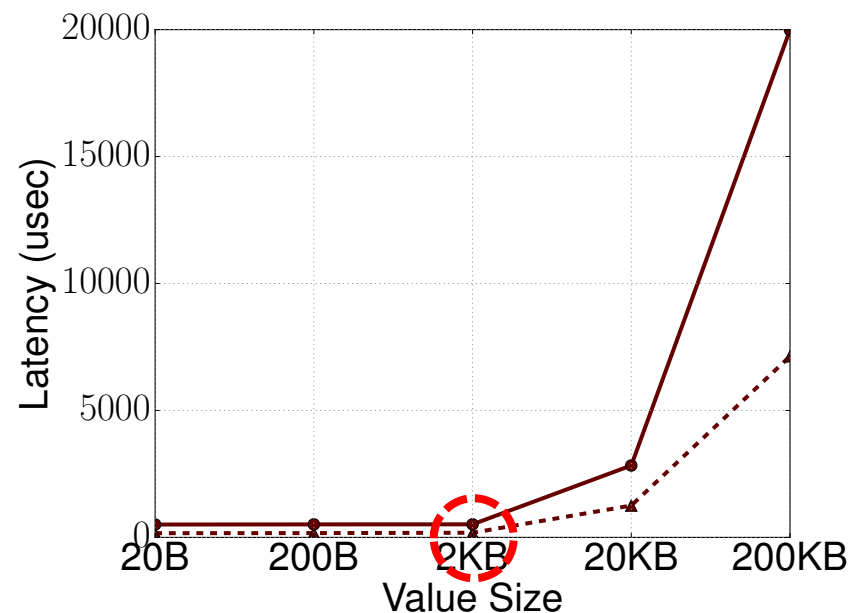Application

Resource Manager

Virtualization

OS

Hardware

- Application bottleneck
- Different user cases
- Scalability


- Overhead of virtualization
- SW isolation mechanisms
- **Overhead of context switching**
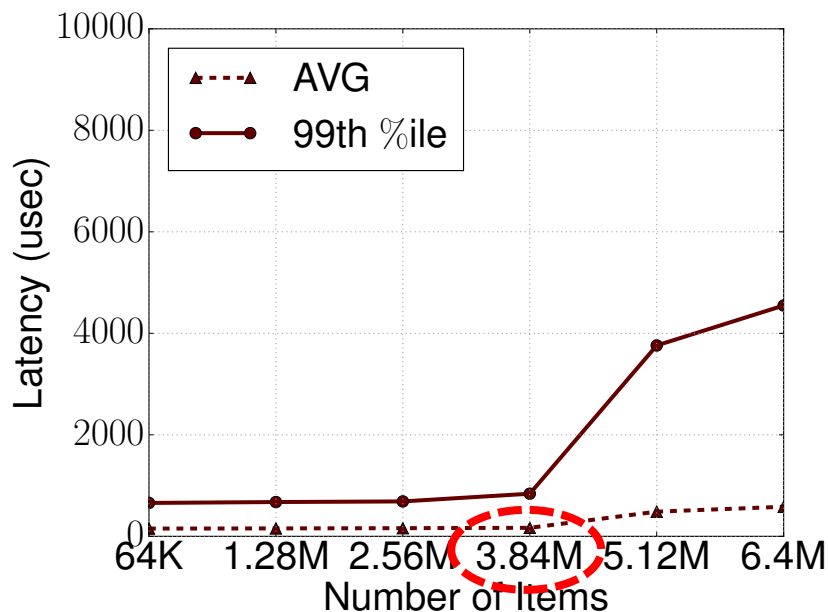- HW isolation mechanisms
- Hyperthreading

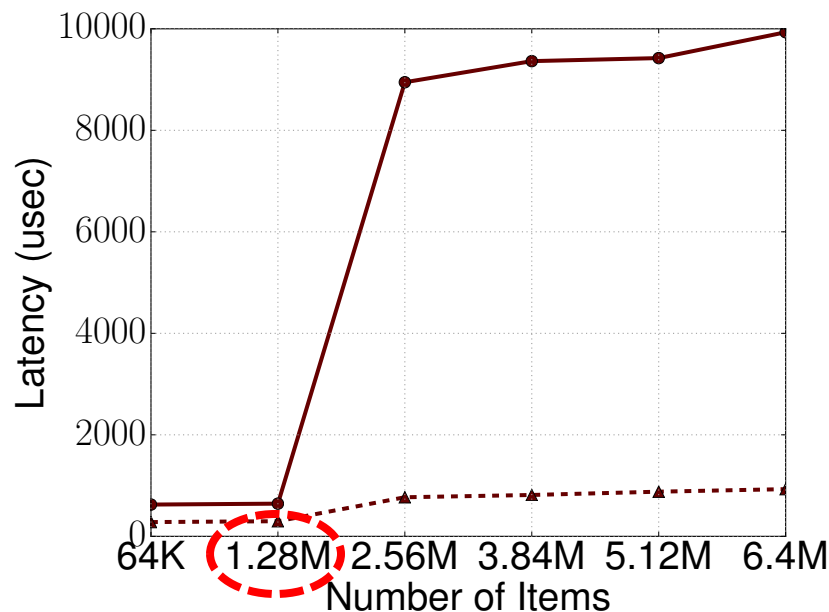**Memcached on Xeon**     **Memcached on ThunderX**

- Statically spawned threads VS dynamically allocated cores
- ThunderX is more sensitive
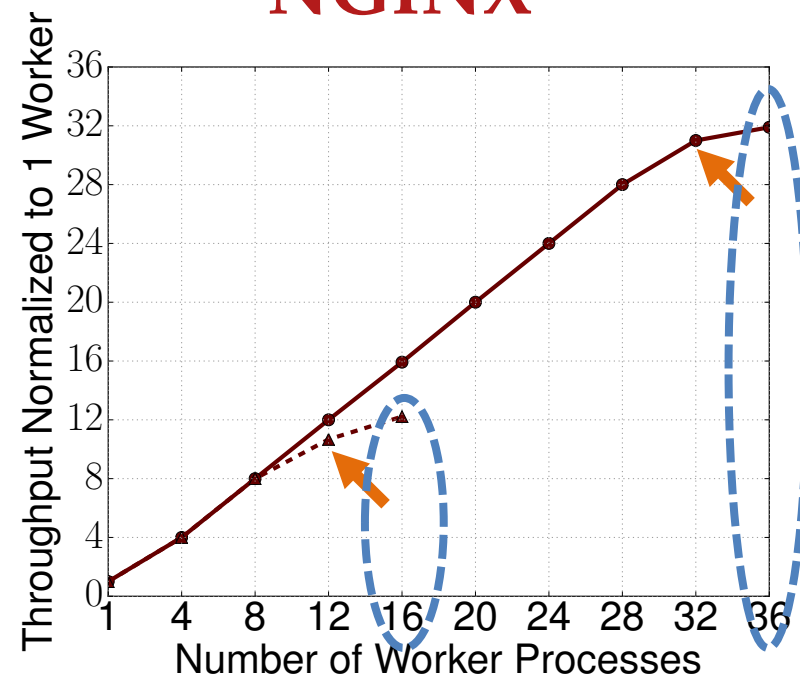
Application

Resource Manager

Virtualization
OS

Hardware

- Application bottleneck
- Different user cases
- Scalability

- Overhead of virtualization
- SW isolation mechanisms
- Overhead of context switching
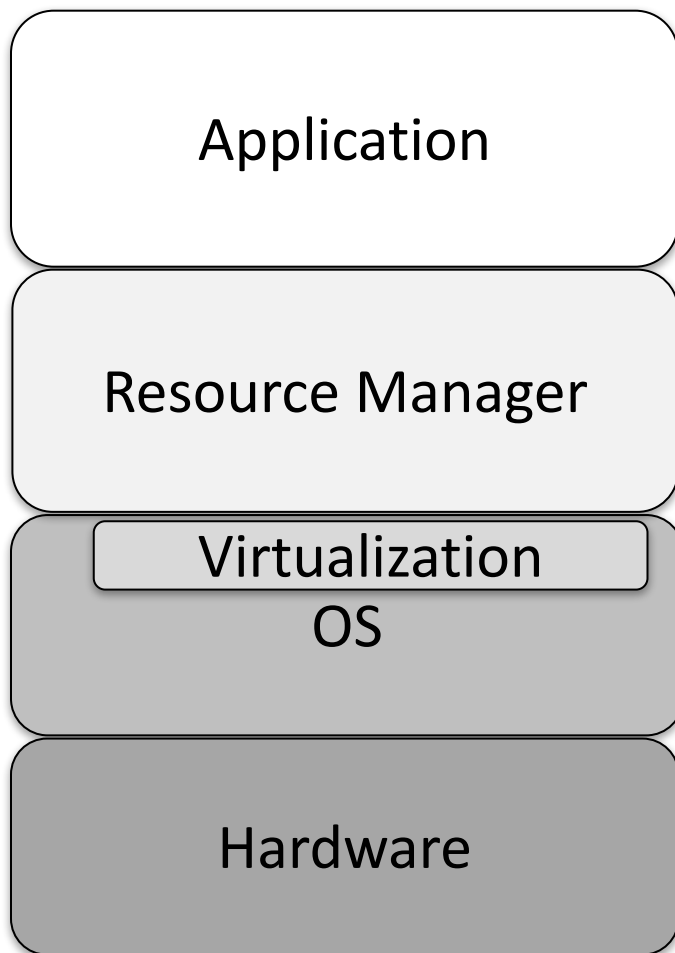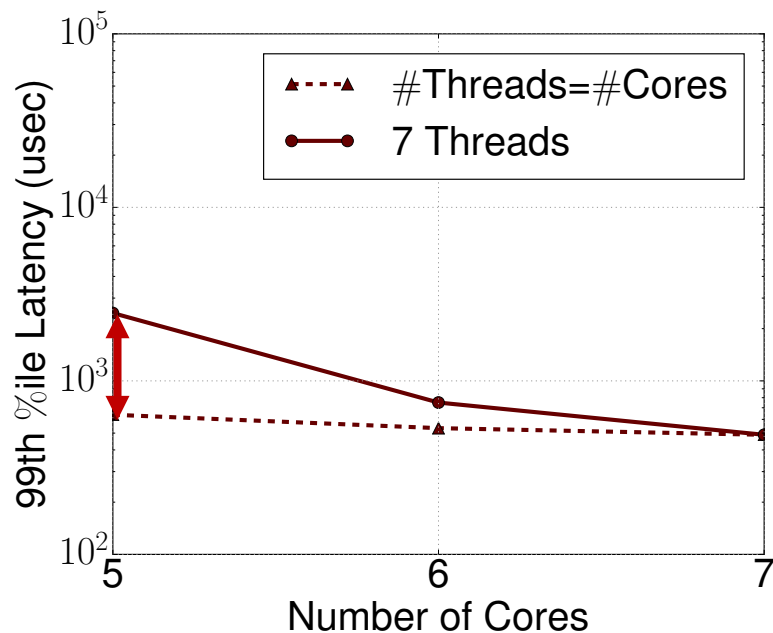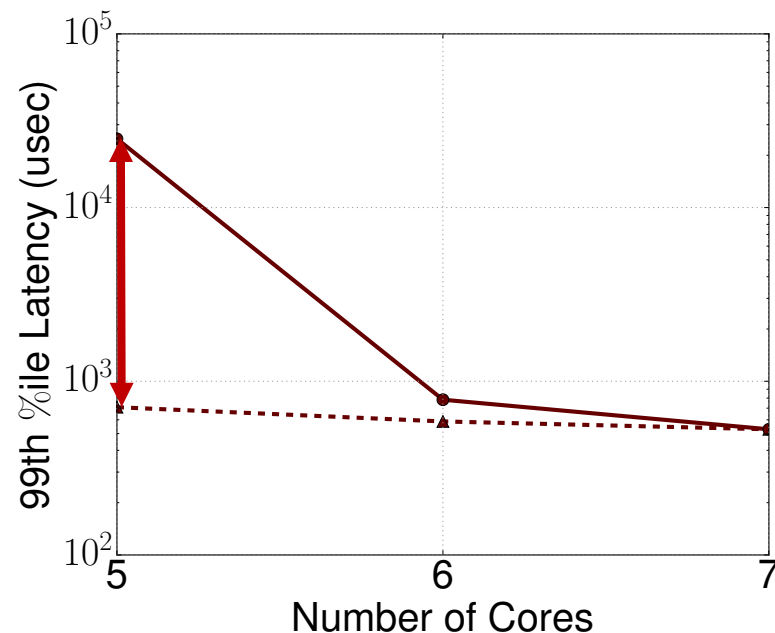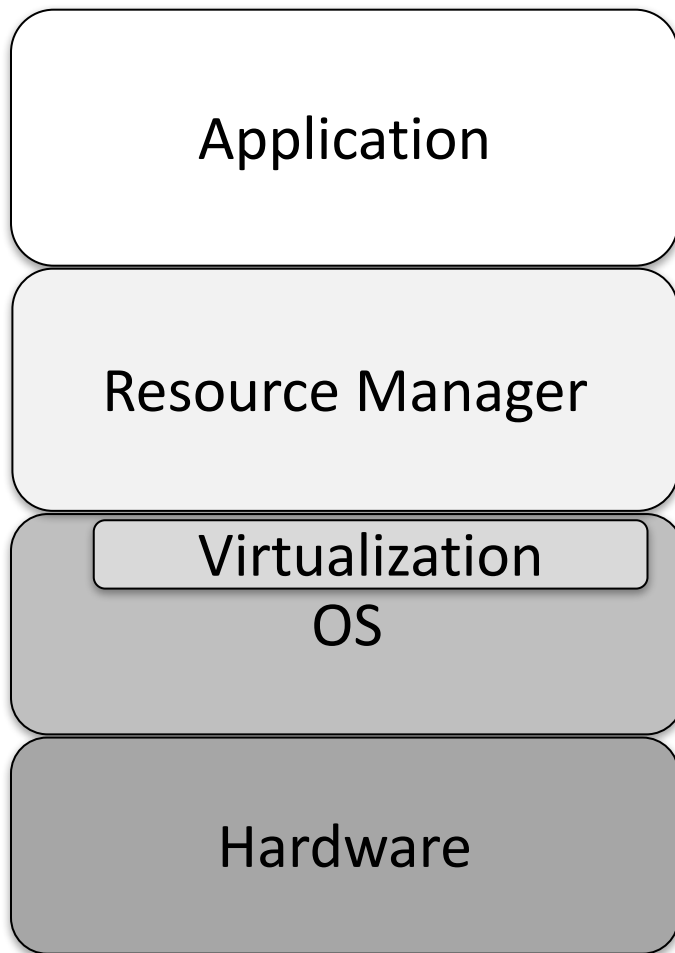- HW isolation mechanisms
- **Hyperthreading**

- **Reduce ... nd of context switching**
  - Allo ... s on two hyperthreads
- **Ma ... execution units**
  - Co-locate different applications

|      | 10% | 20% | 30% | 40% | 50% | 60% | 70% |
|------|-----|-----|-----|-----|-----|-----|-----|
| 10%  | MN  | MN  | MN  | MN  |     |     |     |
| 20%  | MN  | MN  | MN  |     |     |     |     |
| 30%  | MN  | MN  |     |     |     |     |     |
| 40%  | MN  | N   |     |     |     |     |     |
| 50%  | N   | N   |     |     |     |     |     |
| 60%  | N   |     |     |     |     |     |     |
| 70%  |     |     |     |     |     |     |     |

Memcached & Nginx on
the same hyperthreads

|      | 10% | 20% | 30% | 40% | 50% | 60% | 70% |
|------|-----|-----|-----|-----|-----|-----|-----|
| 10%  | MN  | MN  | MN  | MN  | MN  | MN  | M   |
| 20%  | MN  | MN  | MN  | MN  | MN  | M   | M   |
| 30%  | MN  | MN  | MN  | MN  | M   | M   | M   |
| 40%  | MN  | MN  | MN  | MN  | M   | M   | M   |
| 50%  | MN  | MN  | N   | N   |     |     |     |
| 60%  | N   | N   | N   |     |     |     |     |
| 70%  | N   | N   |     |     |     |     |     |

Memcached & Nginx on
different hyperthreads

# QUESTIONS?

**Application**

**Resource Manager**

**Virtualization**
**OS**

**Hardware**

- Reduce network/queuing delays
- Optimize common user cases
- Improve elasticity
  - Lock alternatives
  - Load balance

- Reduce the overhead of virtualization
- Reduce context switching
- Make best use of SW isolation mechanisms

- Big VS Small Cores
- Make best use of HW features