Computer Systems Laboratory
Cornell University
CSL

*ASPLOS 2019*
**Shuang Chen**
**Christina Delimitrou**
**José F. Martínez**

# PARTIES: QoS-Aware Resource Partitioning for Multiple Interactive Service
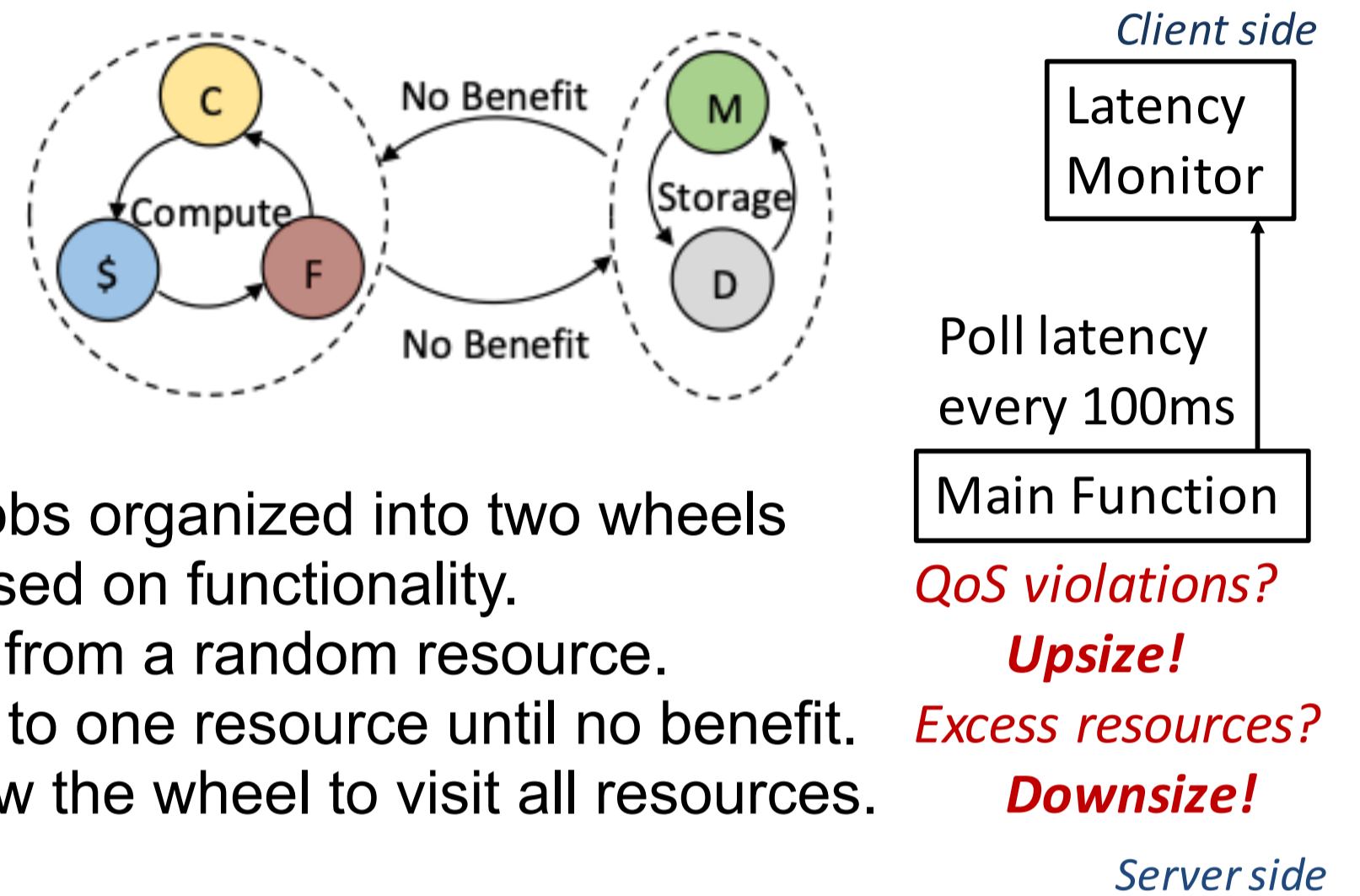
## Motivation

QoS violations ☹

Best-effort

Latency-critical

Hadoop

GraphLab

Private caches  Private caches  …  Private caches  Private caches

Last-level Cache

Google Maps

Bing

Google Translate

Performance unpredictability from Interference in shared resources leads to QoS violations for LC applications.

## Colocation of Multiple LC Applications

Monolith

Batch

hadoop

Bing

Latency-critical

Microservices

**1** LC + many BE

**many** LC + many BE

**Challenge**: all LC services have QoS targets, so none of them can be easily scarified for another.

## Isolation Mechanisms

PARTIES leverages all the existing software and hardware isolation
mechanisms to partition:
* Cores
    * Hyperthreads
    * Core counts
* Power budget
* Last-level cache capacity
    * LLC bandwidth
    * Memory bandwidth
* Memory capacity
* Disk bandwidth
* Network bandwidth

Private caches  Private caches  …  Private caches  Private caches

Last-level Cache

★ cgroup    ● ACPI frequency driver
▲ Intel CAT    ■ qdisc

Xapian

#Cores

Cache ways

>2.2GHz
2.2GHz
2.0GHz
1.8GHz
1.6GHz
1.4GHz
1.2GHz

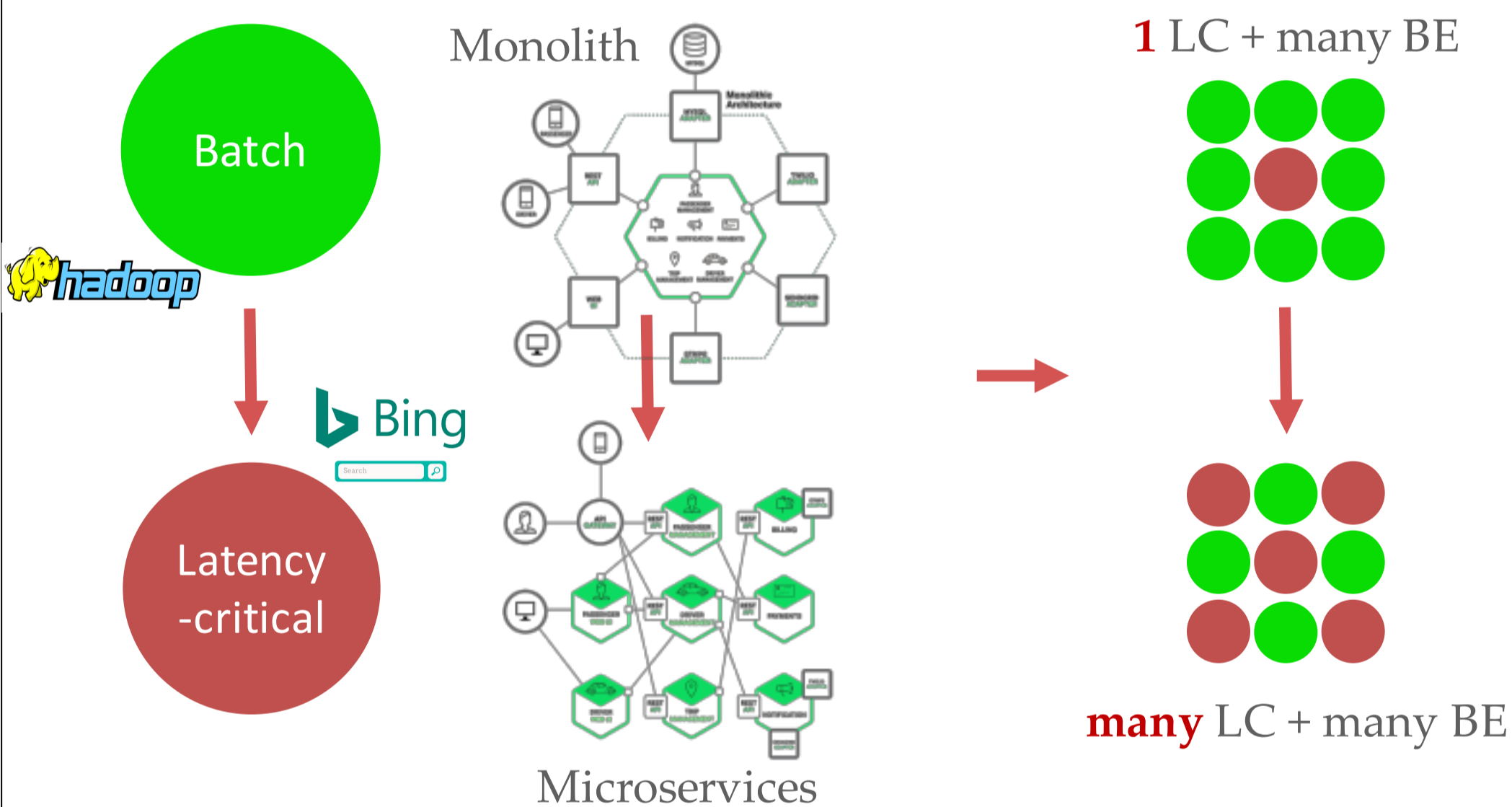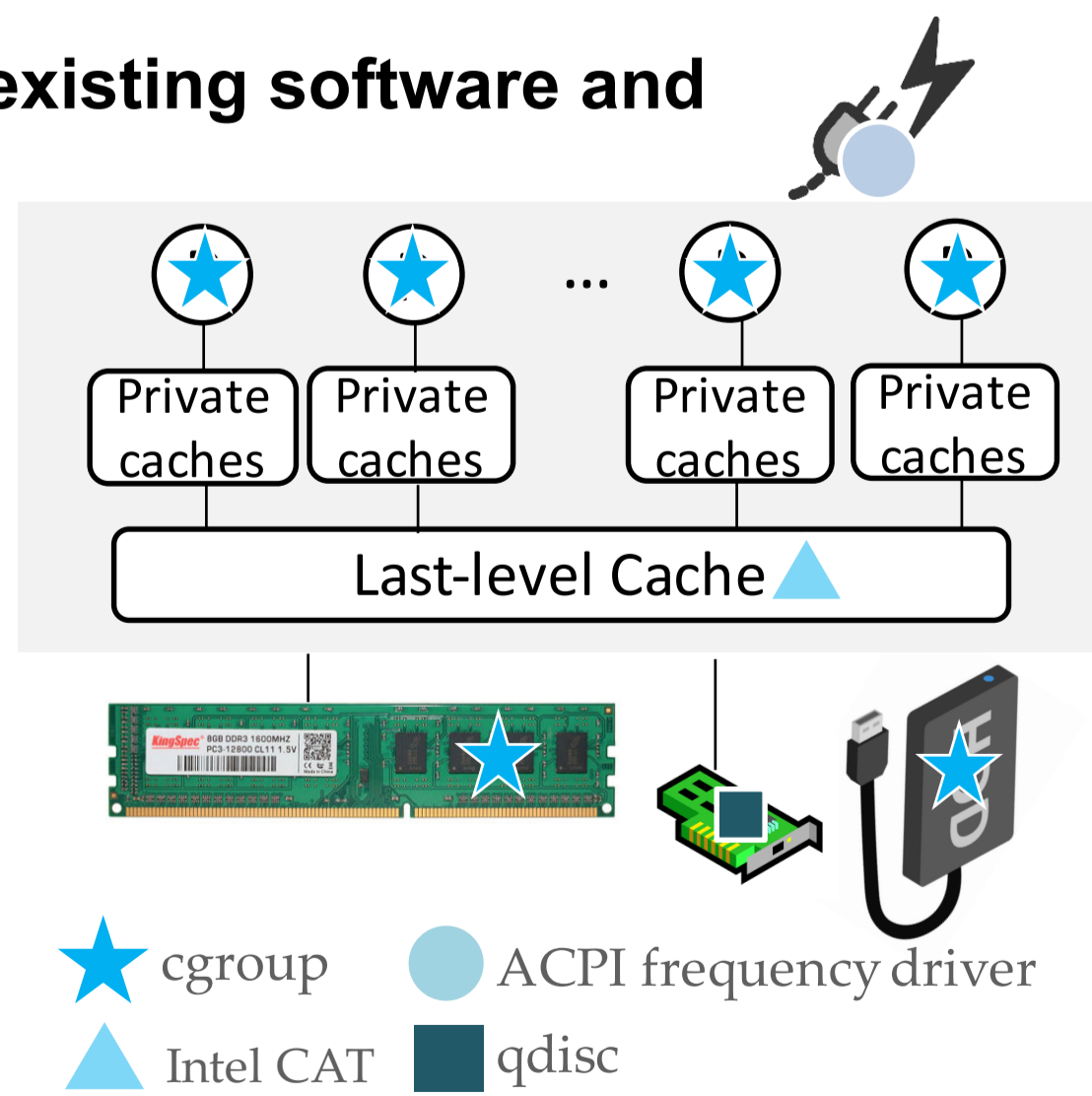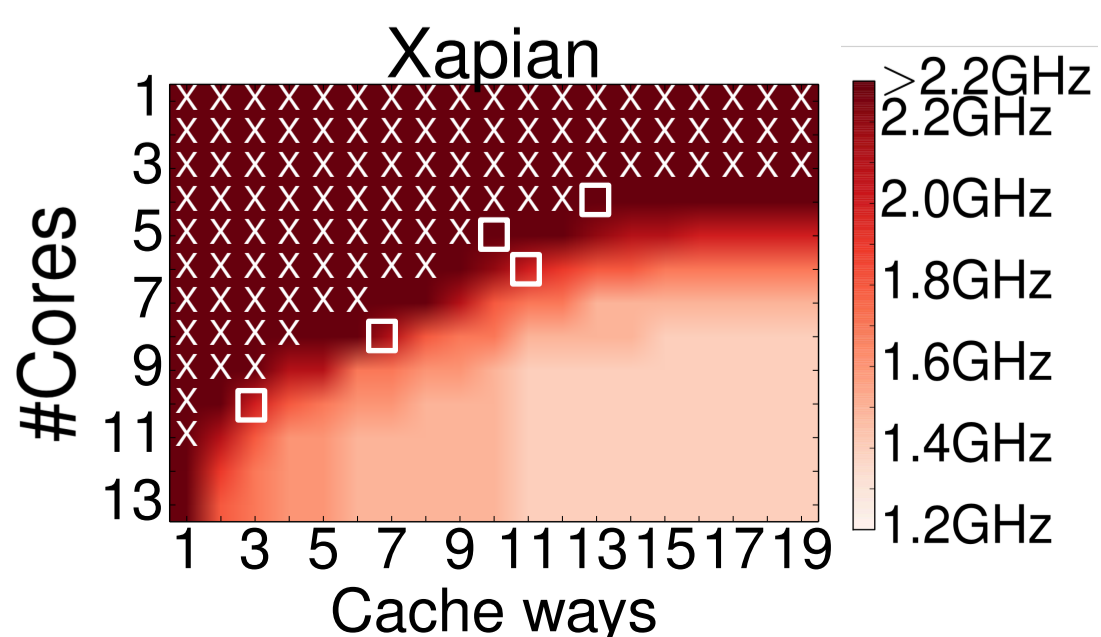**Resource fungibility**: resources can be traded with each other.

## PARTIES Design

**Design principles:**
1. All LC applications are equally important.
2. Allocation should be dynamic and fine-grained.
3. No a priori application knowledge or offline profiling is needed.
4. Recover quickly from incorrect decisions.
5. Migration is used as a last resort.

C
Compute
$ — F
No Benefit
M
Storage
D
No Benefit

1. 5 knobs organized into two wheels based on functionality.
2. Start from a random resource.
3. Stick to one resource until no benefit.
4. Follow the wheel to visit all resources.

*Client side*

Latency Monitor

Poll latency every 100ms

Main Function

*QoS violations?*
**Upsize!**
*Excess resources?*
**Downsize!**

*Server side*

## Evaluation

**Platform:** Intel E5-2699 v4
**Benchmarks:** Memcached; Xapian; NGINX; Moses; MongoDB; Sphinx

(a) Unmanaged    (b) Heracles    (c) PARTIES    (d) Oracle

Max Load of NGINX(%)
Max Load of Memcached(%)
Max Load of Xapian(%)

Heracles    PARTIES

EMU (%)
Number of Colocated Apps

Convergence Time (s)
#Colocated Apps

% of Max Load
moses
xapian
memcached
Time (s)

Norm. Latency w. Heracles
Time (s)

Norm. Latency w. PARTIES
Time (s)

#Cores
Time (s)