# *Autothrottle:*

## A Practical Bi-Level Approach to Resource Management for SLO-Targeted Microservices

**Zibo Wang**[1,2], Pinghe Li[3], Chieh-Jan Mike Liang[1], Feng Wu[2], Francis Y. Yan[1]

[1] *Microsoft Research*

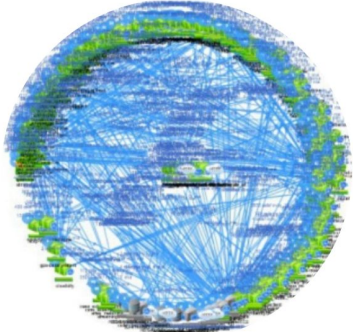[2] *University of Science and Technology of China*

[3] *ETH Zurich*

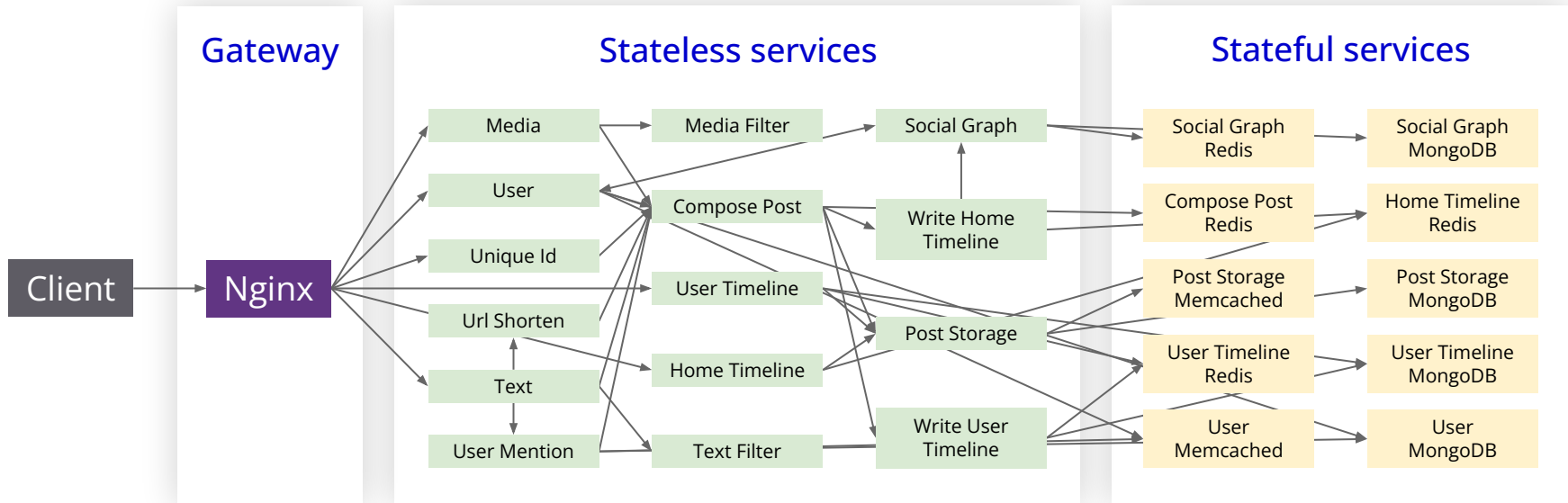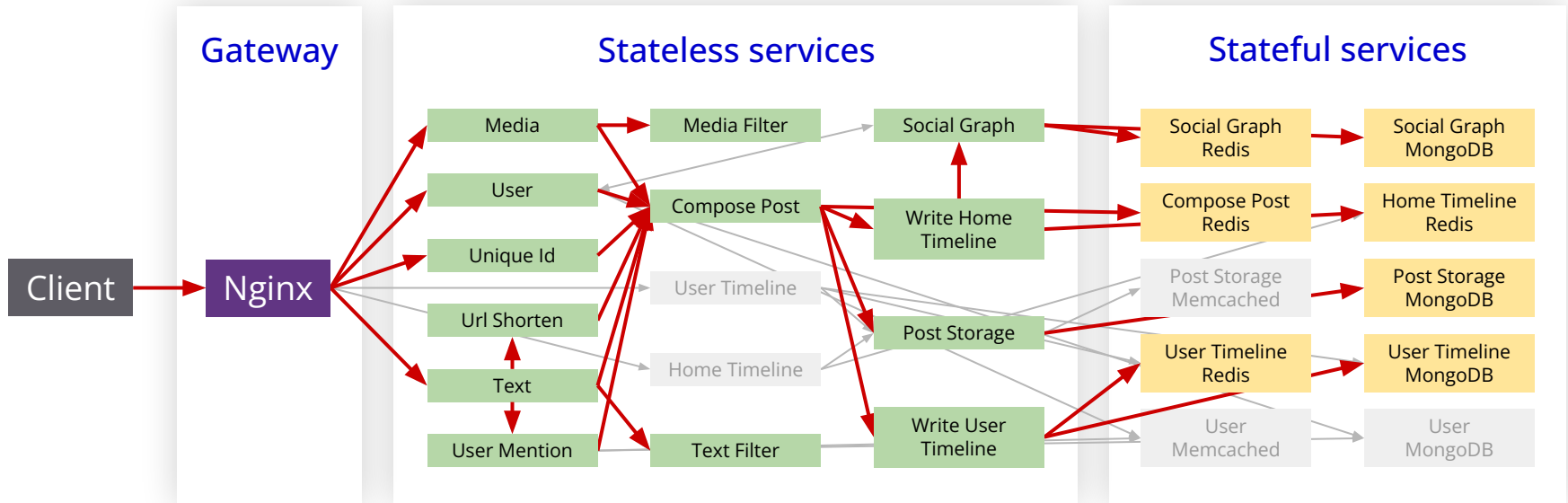# Cloud applications are shifting toward microservices
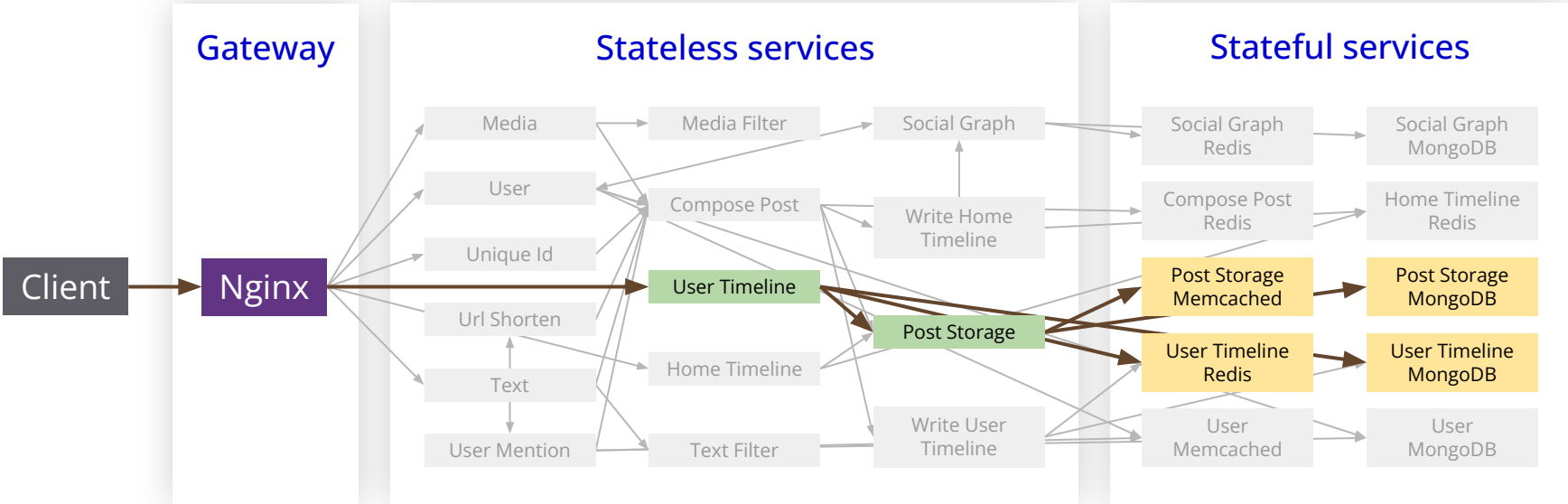


Netflix

Twitter

# What microservice applications look like
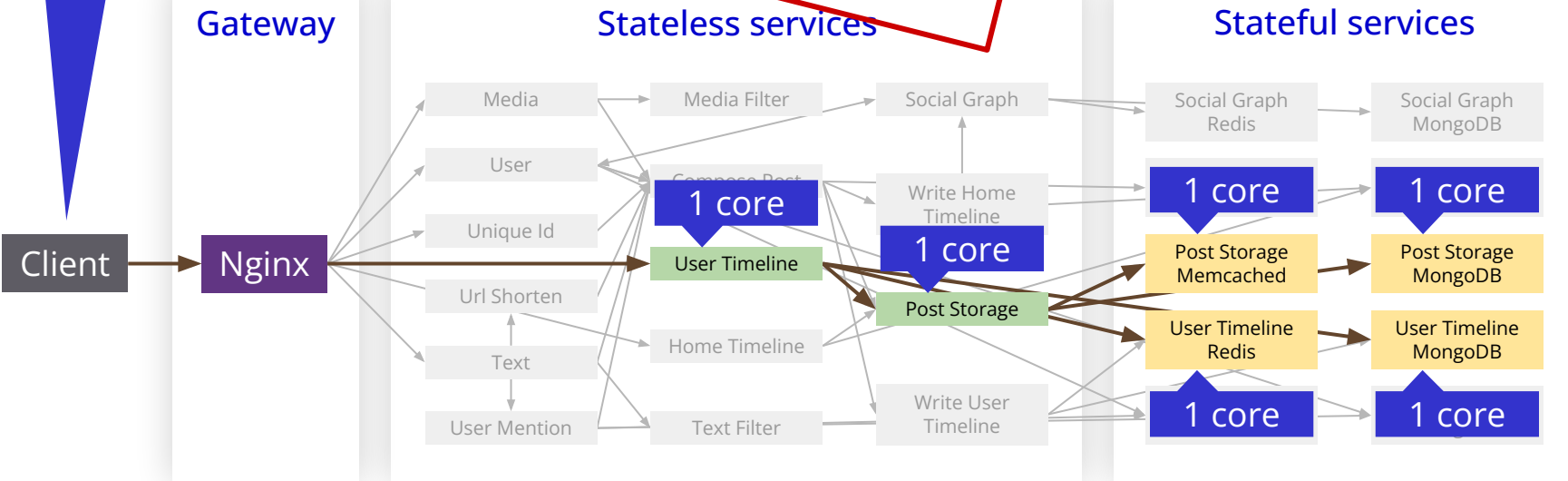
# A client request traverses many services
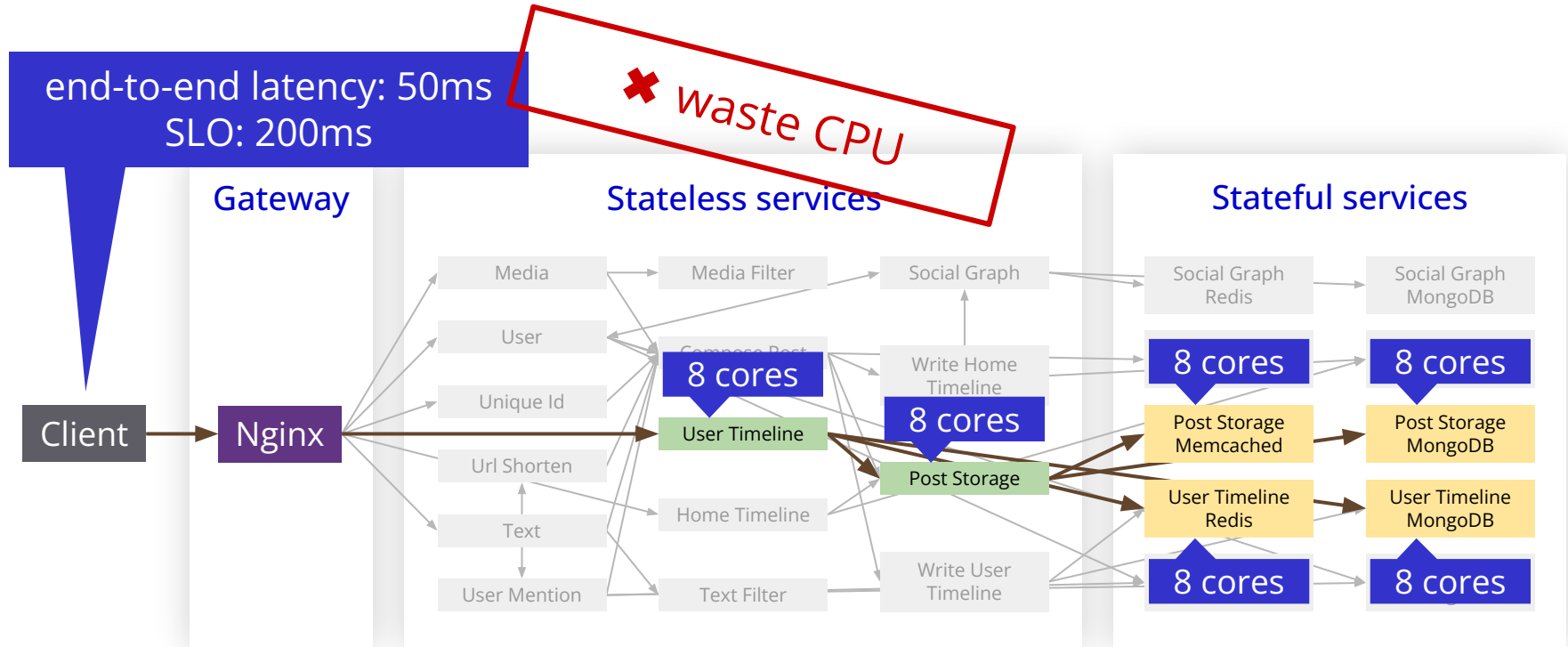
# Different requests have different trajectories

# Inadequate CPU allocations => high application latency
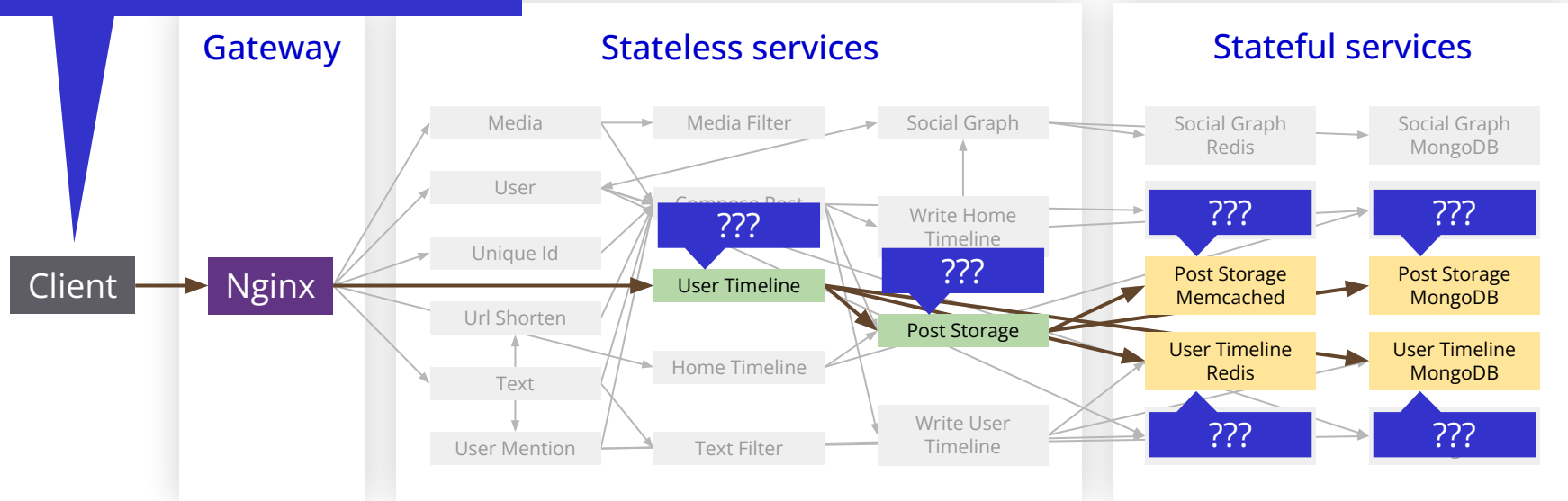
end-to-end latency: 800ms
SLO: 200ms

✖ SLO violation

**Gateway**

**Stateless services**

**Stateful services**

Media

Media Filter

Social Graph

Social Graph Redis

Social Graph MongoDB

User

Compose Post

Write Home Timeline

1 core

1 core

Unique Id

1 core

Post Storage Memcached

Post Storage MongoDB

Client → Nginx →

User Timeline

Url Shorten

1 core

Post Storage

User Timeline Redis

User Timeline MongoDB

Text

Home Timeline

User Mention

Text Filter

Write User Timeline

1 core

1 core

# Excessive CPU allocations => waste of resources



end-to-end latency: 50ms
SLO: 200ms

❌ waste CPU

Gateway

Stateless services

Stateful services

Client

Nginx

Media

User

Unique Id

Url Shorten

Text

User Mention

Media Filter

Compose Post

8 cores
User Timeline

Home Timeline

Text Filter

Social Graph

Write Home Timeline

8 cores
Post Storage

Write User Timeline

Social Graph Redis

Social Graph MongoDB

8 cores

8 cores

Post Storage Memcached

Post Storage MongoDB

User Timeline Redis

User Timeline MongoDB

8 cores

8 cores

# Minimizing CPU allocation while meeting SLO

- Search space grows exponentially with number of services
- Mapping from CPU allocations to latency is unclear

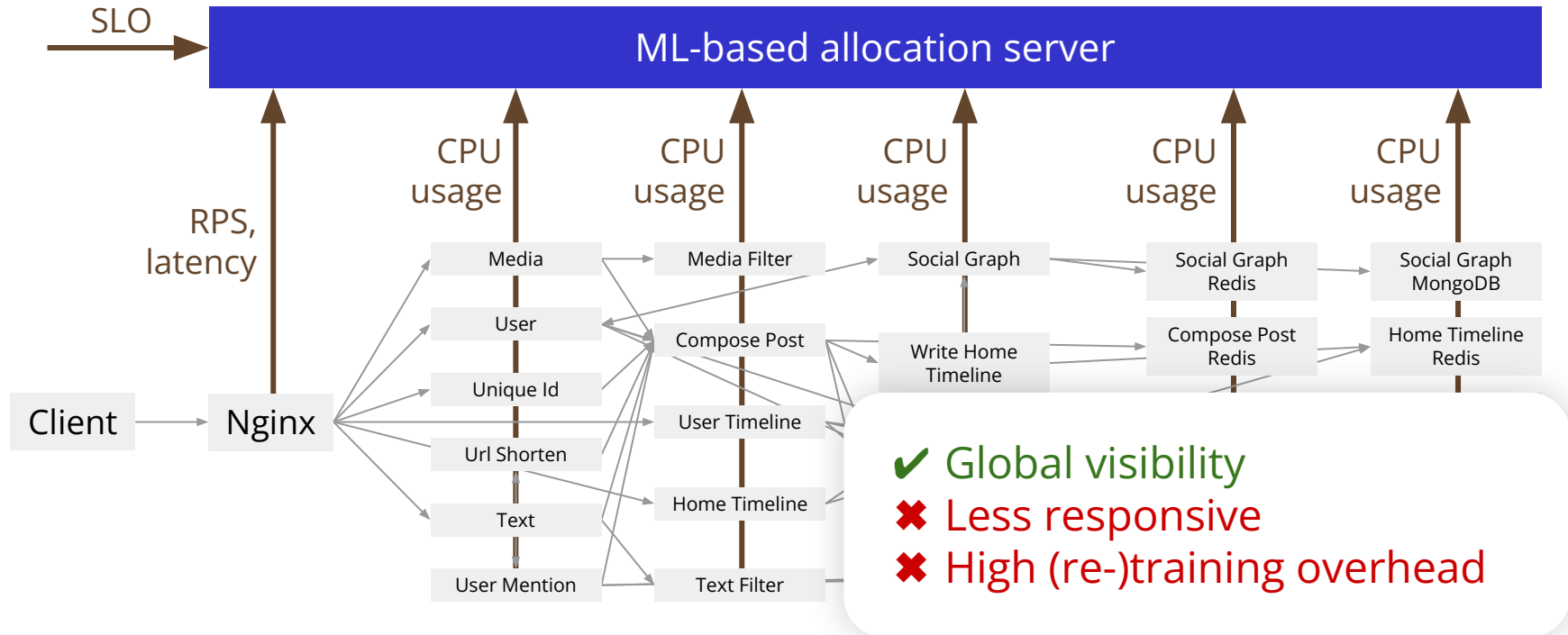# Existing approach: service-level allocation

- Example: Kubernetes' default heuristics

**Per-service allocation algorithm**

Parameter: target CPU utilization = 50%
Input: CPU usage = 2 cores
Output: CPU allocation = 2 / 50% = 4 cores

Client → Nginx

User
Unique Id
Url Shorten
Text
User Mention

Comp...
User Timeline
Home Timeline
Text Filter

Write Home Timeline

Social Graph Redis → Social Graph MongoDB
Compose Post Redis → Home Timeline Redis

✔ Low overhead
✔ Fast reaction
✖ No global visibility

# Existing approach: application-level allocation

- Example: Sinan (ASPLOS '21)
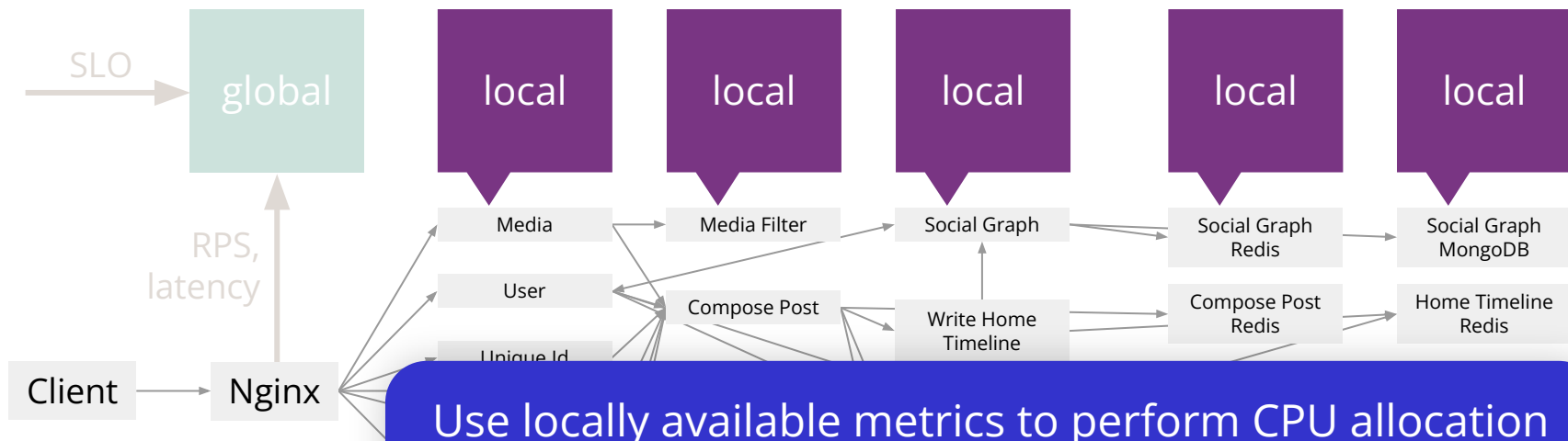


SLO

ML-based allocation server

CPU usage

RPS, latency

Client → Nginx

Media
User
Unique Id
Url Shorten
Text
User Mention

Media Filter
Compose Post
User Timeline
Home Timeline
Text Filter

Social Graph
Write Home Timeline

Social Graph Redis
Compose Post Redis

Social Graph MongoDB
Home Timeline Redis

✔ Global visibility
✖ Less responsive
✖ High (re-)training overhead

# How to obtain the best of both worlds?

|  | Service-level | Application-level | ??? |
|---|:---:|:---:|:---:|
| Low overhead Fast reaction | ✅ | ❌ | ✅ |
| Global visibility | ❌ | ✅ | ✅ |

# Our bi-level approach to resource management

# Our bi-level approach to resource management



SLO

global

local | local | local | local | local

RPS, latency

| Media | Media Filter | Social Graph | Social Graph Redis | Social Graph MongoDB |
| User | Compose Post | Write Home Timeline | Compose Post Redis | Home Timeline Redis |
| Unique Id | | | | |

Client → Nginx

Use locally available metrics to perform CPU allocation

✔ Low overhead
✔ Fast reaction

# Our bi-level approach to resource management



SLO

global

local      local      local      local      local

RPS, latency

Media    Media Filter    Social Graph    Social Graph Redis    Social Graph MongoDB

User    Compose Post    Write Home Timeline    Compose Post Redis    Home Timeline Redis

Unique Id

Client → Nginx

**Monitor RPS, end-to-end latencies, and SLO violations**

**✔ Global visibility**

# Our bi-level approach to resource management

# Implementing bi-level approach with **Autothrottle**

# Interface: throttle ratio

# Interface: throttle ratio

- Example: Linux CFS (Completely Fair Scheduler)



$$\text{throttle ratio} = \frac{1 \text{ period}}{3 \text{ periods}}$$

# Throttle ratio has a higher correlation with latency
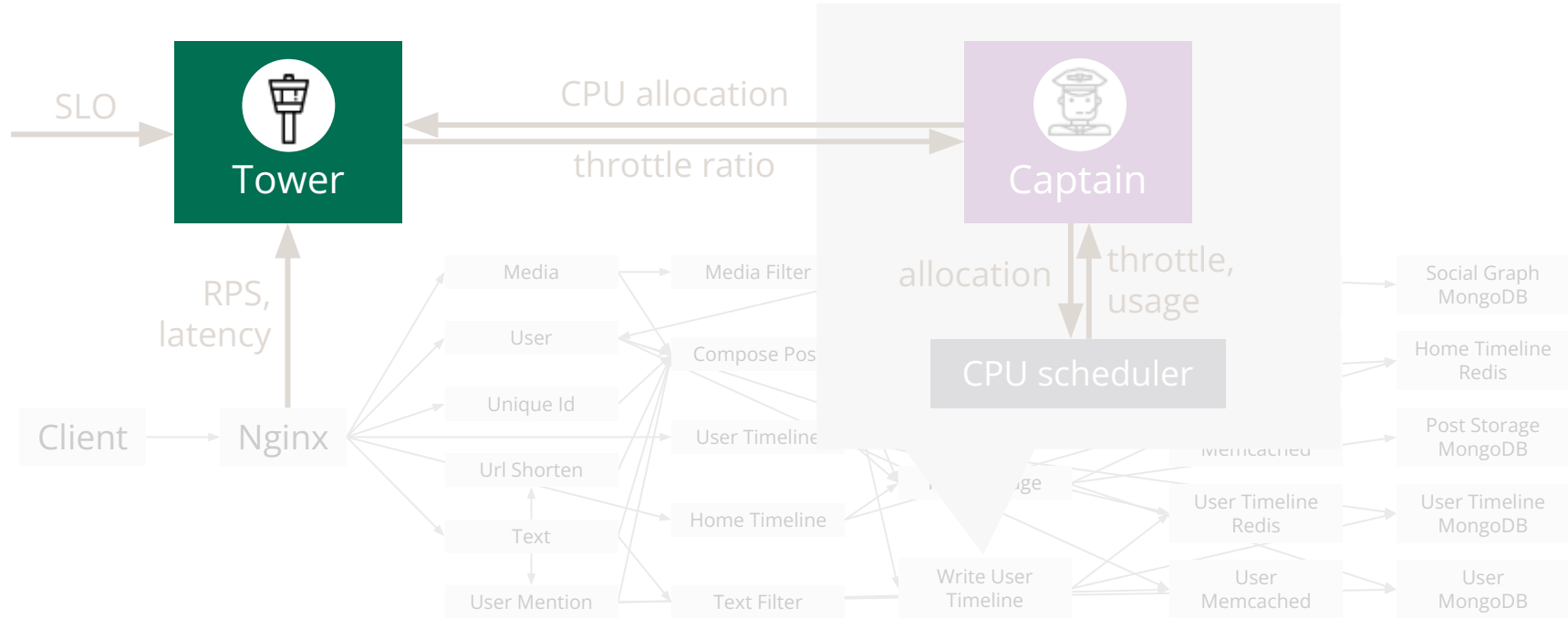
# Service-level: fast and lightweight Captains

# Service-level: fast and lightweight Captains

- Closed-loop control based on throttle ratio target
- Collect data every 100ms, adjust allocation every 1s

# Application-level: online learning Tower

# Application-level: online learning Tower

- Determine the best throttle targets for Captains to achieve

- Lightweight online learning: contextual bandit algorithm
  - One step per minute, each step runs in ~100ms

context: RPS

Tower

action #1

action #2

action #3

......

action: all performance targets for Captains

cost: computed with CPU allocation, end-to-end latency, and SLO

# Evaluation methodology

- **Testbed:** 5 Azure VMs, 160 CPU cores in total

- **4 workload traces**
  - with patterns commonly observed in production environments
  - e.g. Puffer's streaming requests, Google's cluster usage, and Twitter tweets

- **3 benchmark applications**
  - Train-Ticket
  - Hotel-Reservation from DeathStarBench
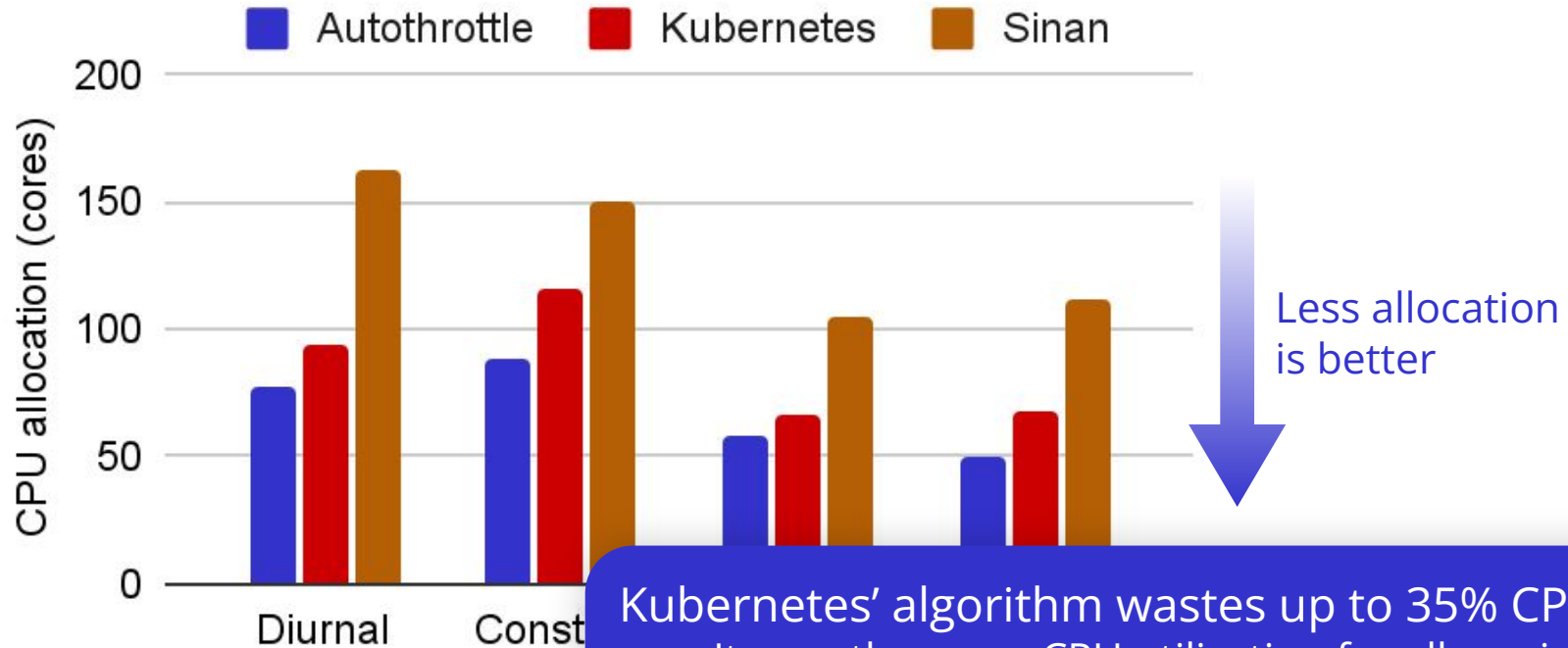  - **Social-Network** used in Sinan



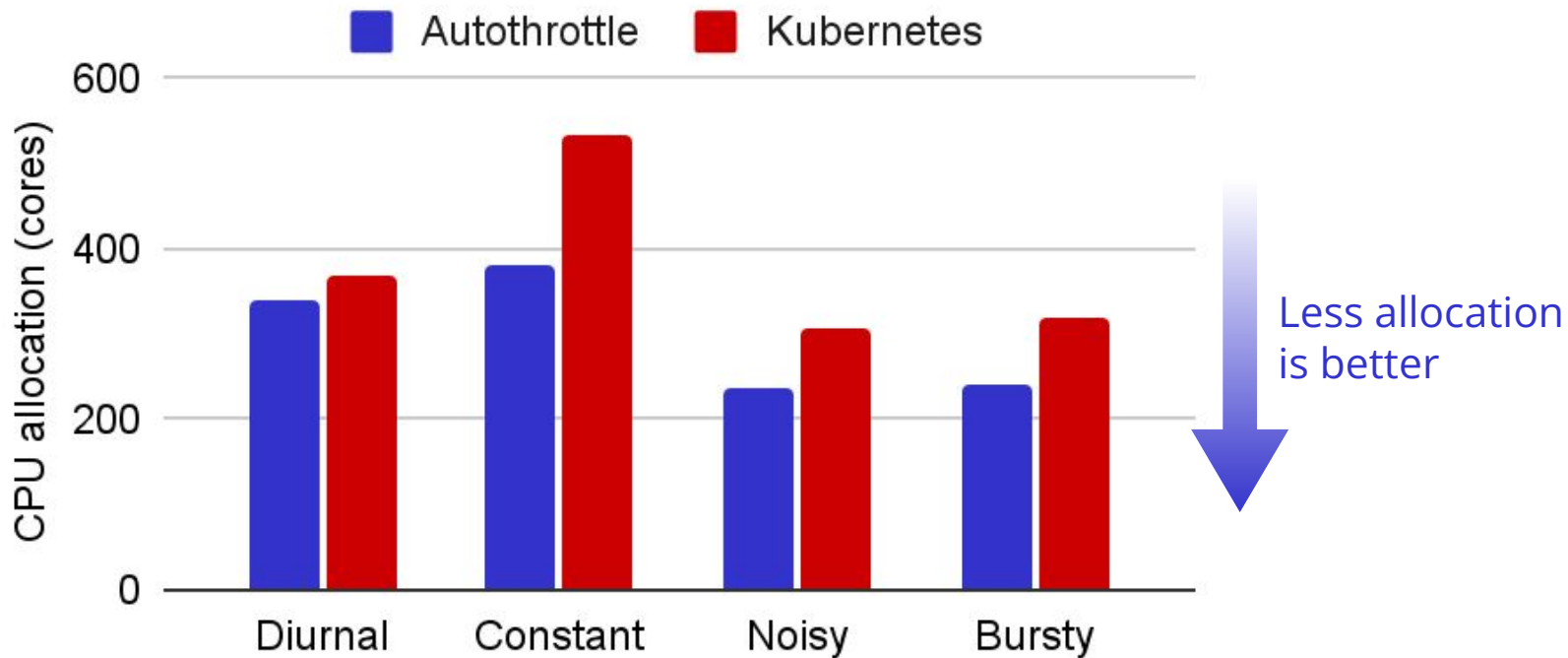(a) Diurnal  (b) Constant  (c) Noisy  (d) Bursty

# Evaluation results



Less allocation is better

Kubernetes' algorithm wastes up to 35% CPU
- It uses the same CPU utilization for all services

Sinan wastes even more CPU
- Its search space is too large to explore

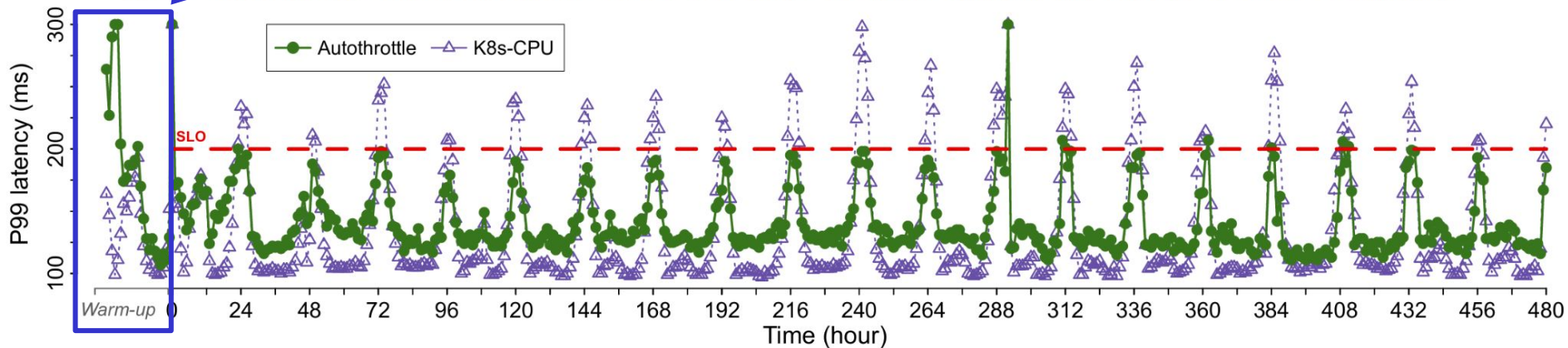# Large-scale evaluation on a 512-core cluster



Less allocation is better

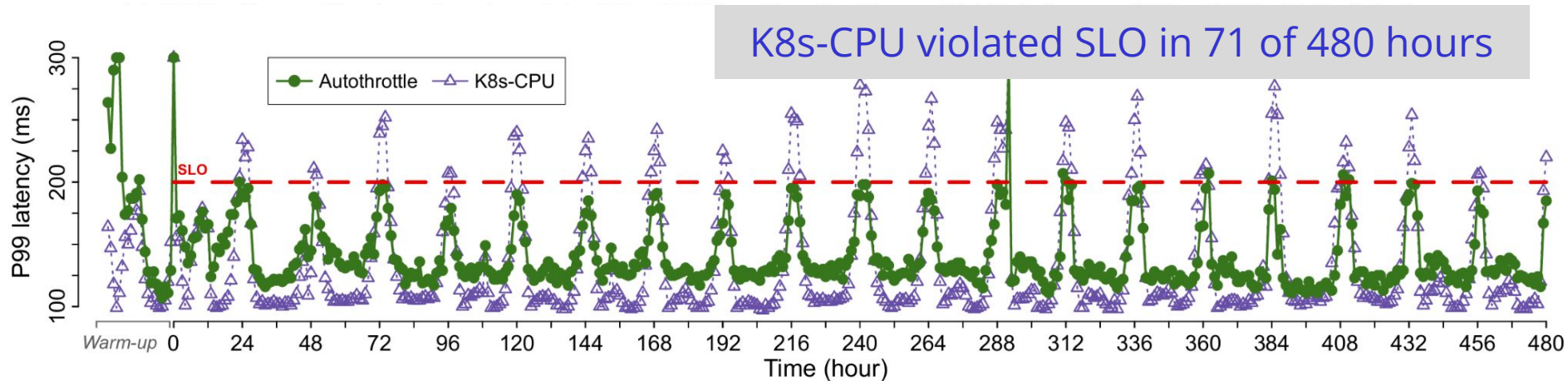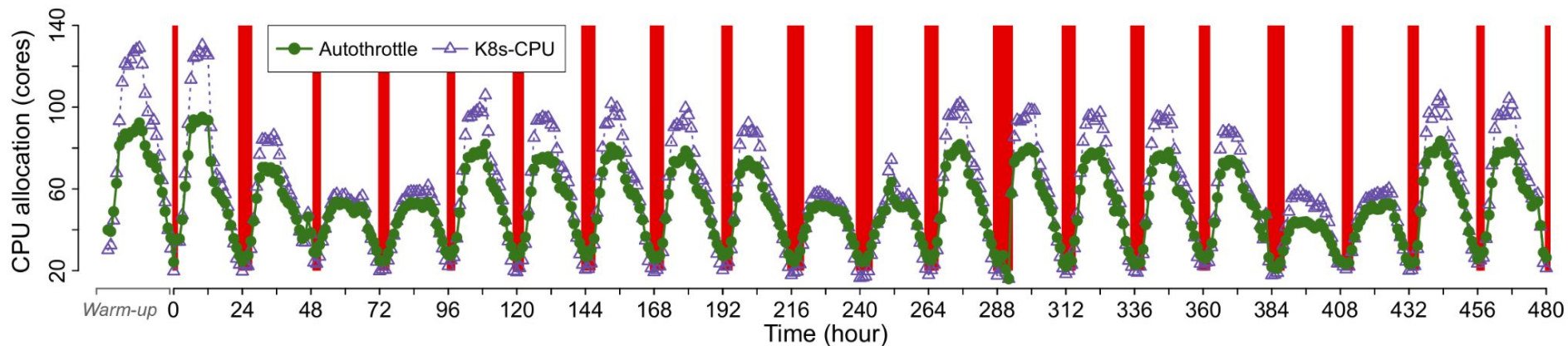Kubernetes' algorithm wastes up to 39% CPU

# A 21-day comparison

# A 21-day comparison

# A 21-day comparison

# A 21-day comparison



K8s-CPU violated SLO in 71 of 480 hours

# A 21-day comparison



K8s-CPU violated SLO in 71 of 480 hours
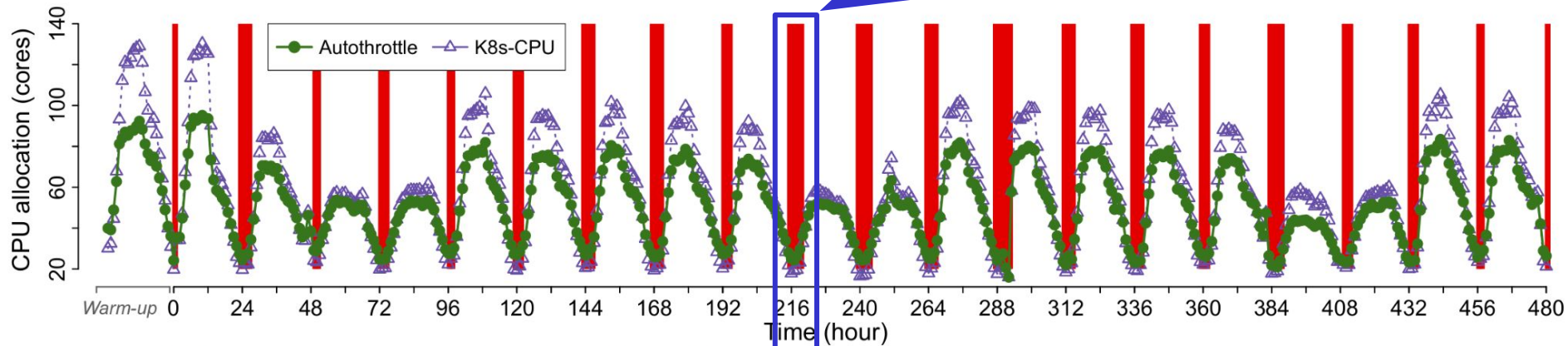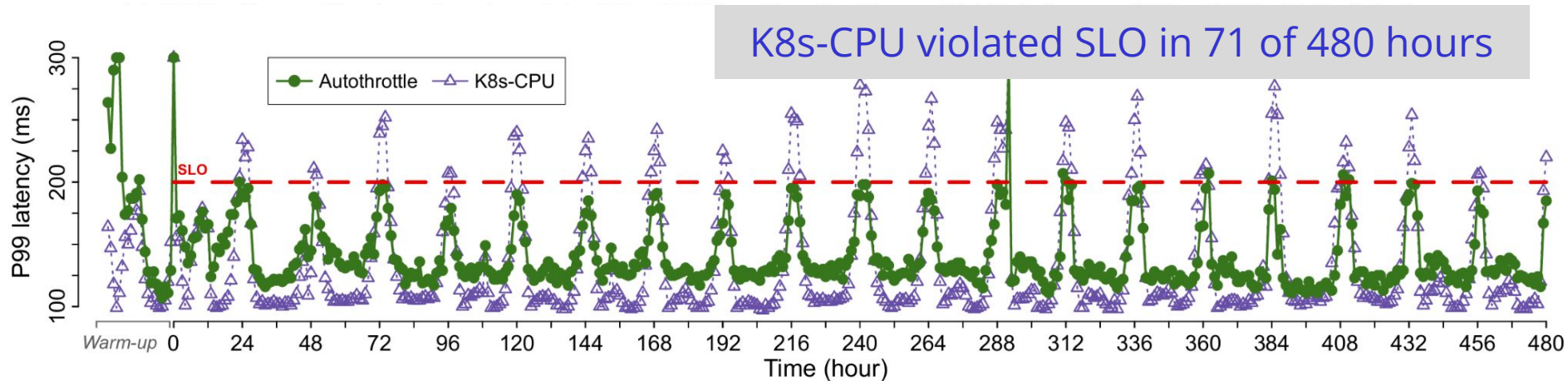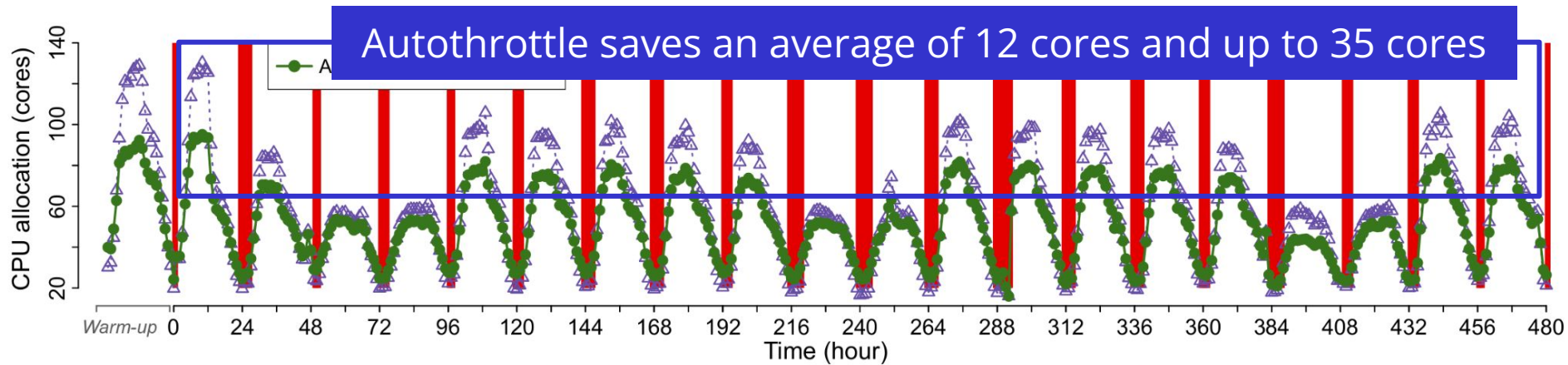
# A 21-day comparison



Red boxes mark K8s-CPU's SLO violations

K8s-CPU violated SLO in 71 of 480 hours

# A 21-day comparison
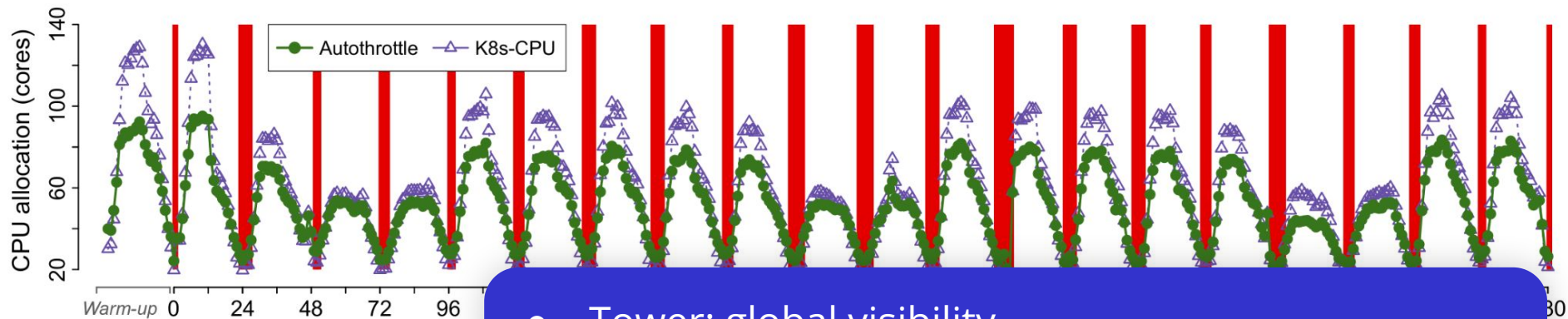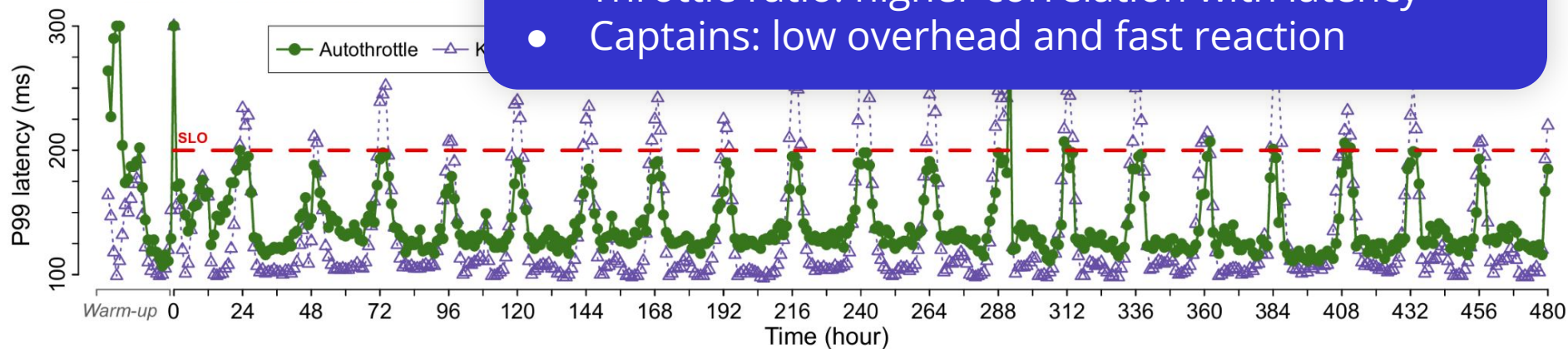


Autothrottle saves an average of 12 cores and up to 35 cores

K8s-CPU violated SLO in 71 of 480 hours

# A 21-day comparison



- Tower: global visibility
- Throttle ratio: higher correlation with latency
- Captains: low overhead and fast reaction

# Conclusion

- Autothrottle: a bi-level learning-assisted resource management framework for SLO-targeted microservices.

- Results show a CPU saving up to 26% while satisfying SLO

- Open-sourced at https://github.com/microsoft/autothrottle