

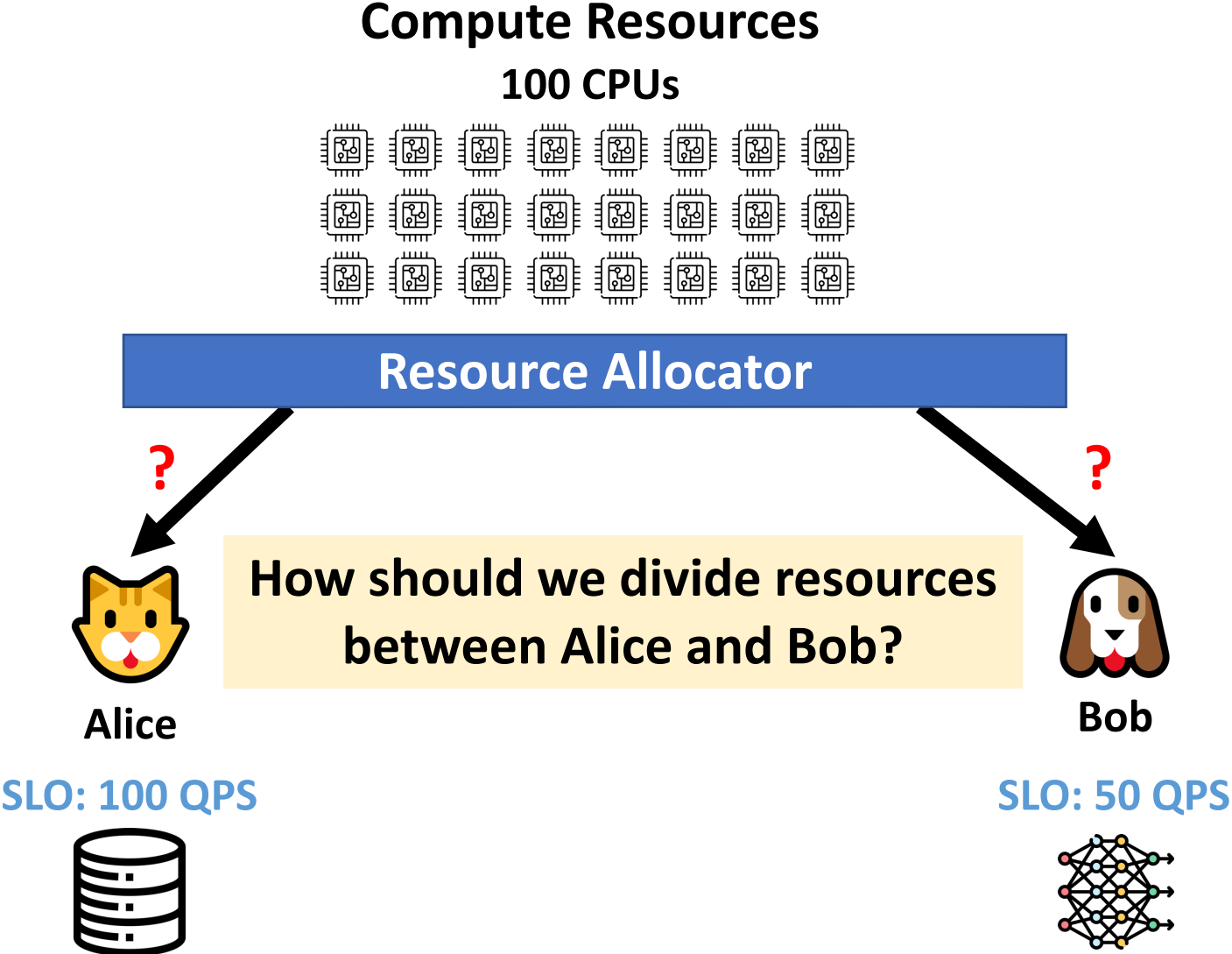
Cilantro

Performance-Aware Resource Allocation for General Objectives via Online Feedback

Romil Bhardwaj, Kirthevasan Kandasamy, Asim Biswal, Wenshuo Guo,
Benjamin Hindman, Joseph Gonzalez, Michael Jordan, Ion Stoica



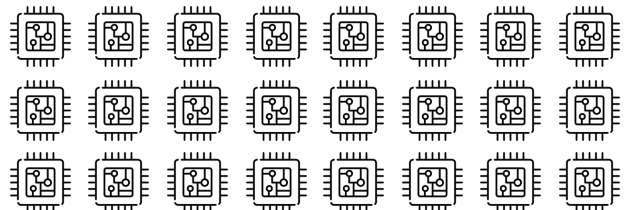
Resource allocation in multi-tenant clusters



Equal Allocation

Compute Resources

100 CPUs



Resource Allocator

50 CPUs



Alice

SLO: 100 QPS

I needed only 20 CPUs for 100 QPS

50 CPUs



Bob

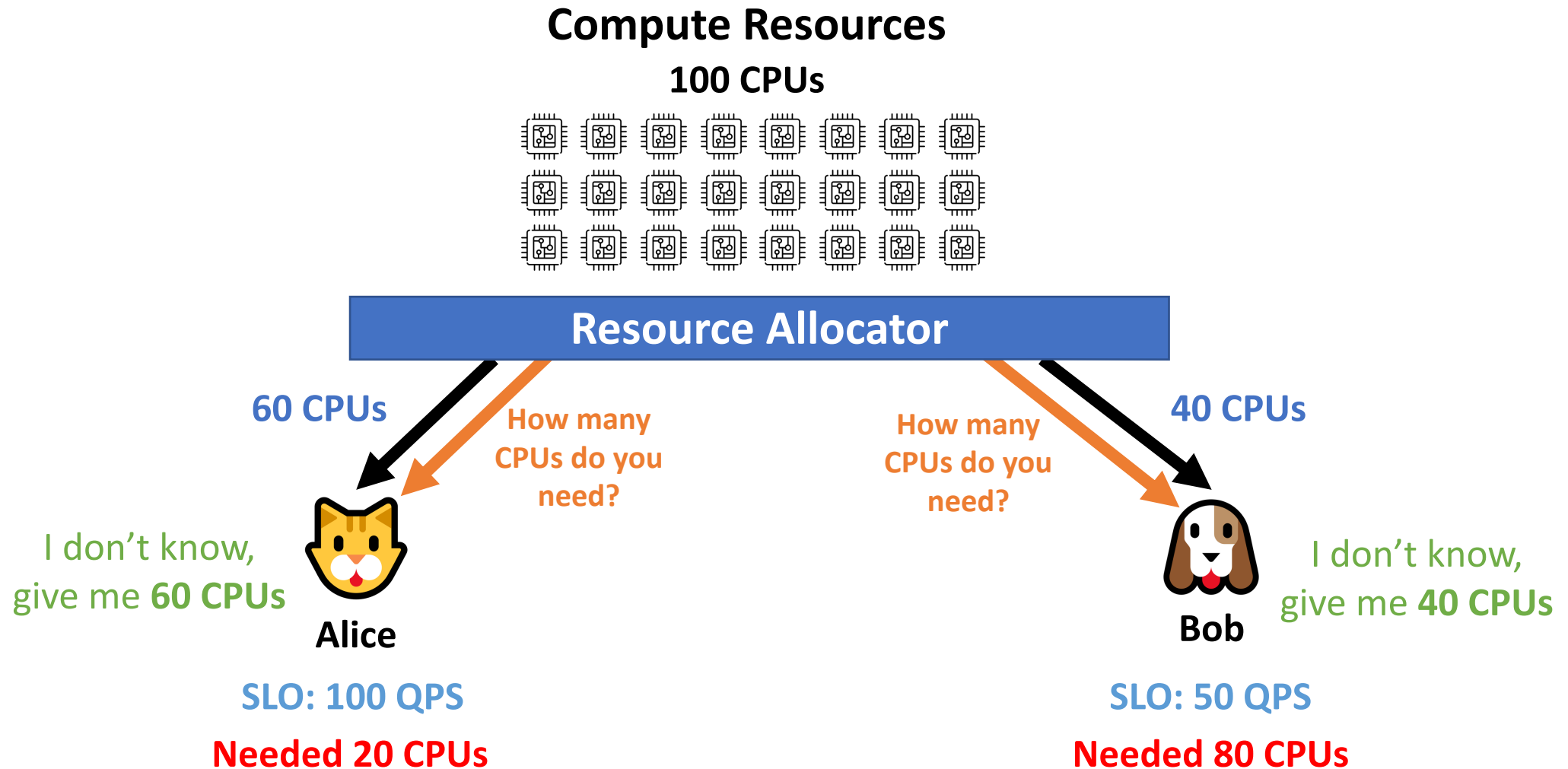
SLO: 50 QPS

I needed 80 CPUs for 50 QPS

Equal shares are fair*, but inefficient

³
*fair = proportional shares

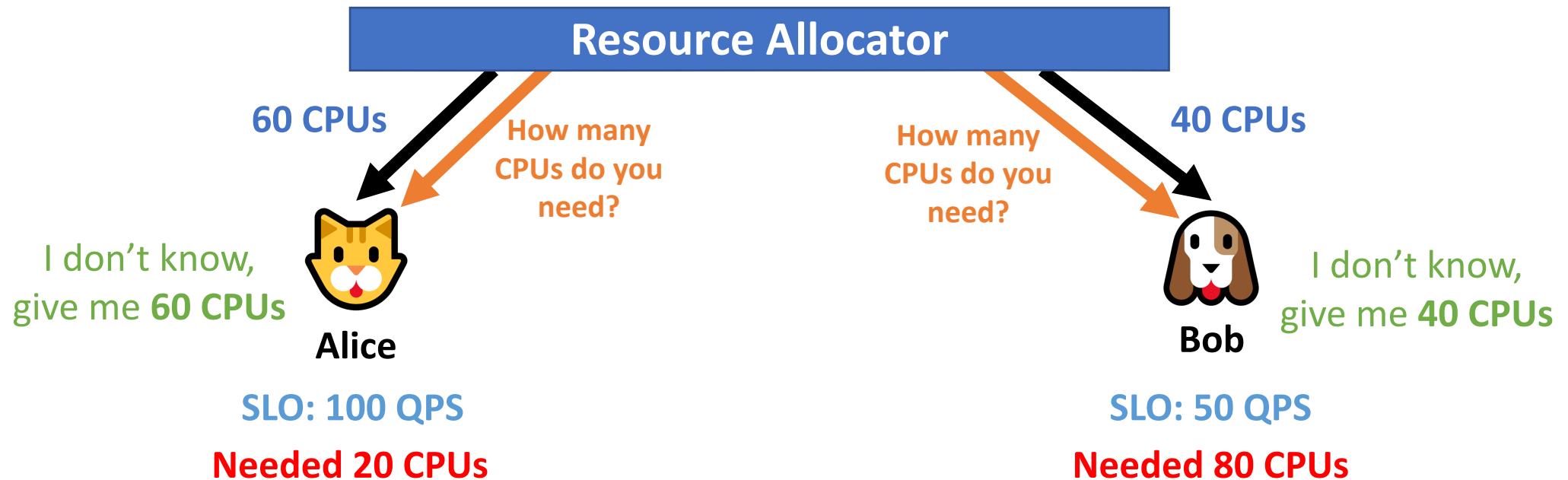
Asking for resource demands



Asking for resource demands

Compute Resources
100 CPUs

Inaccurate resource demands result in inefficient allocations

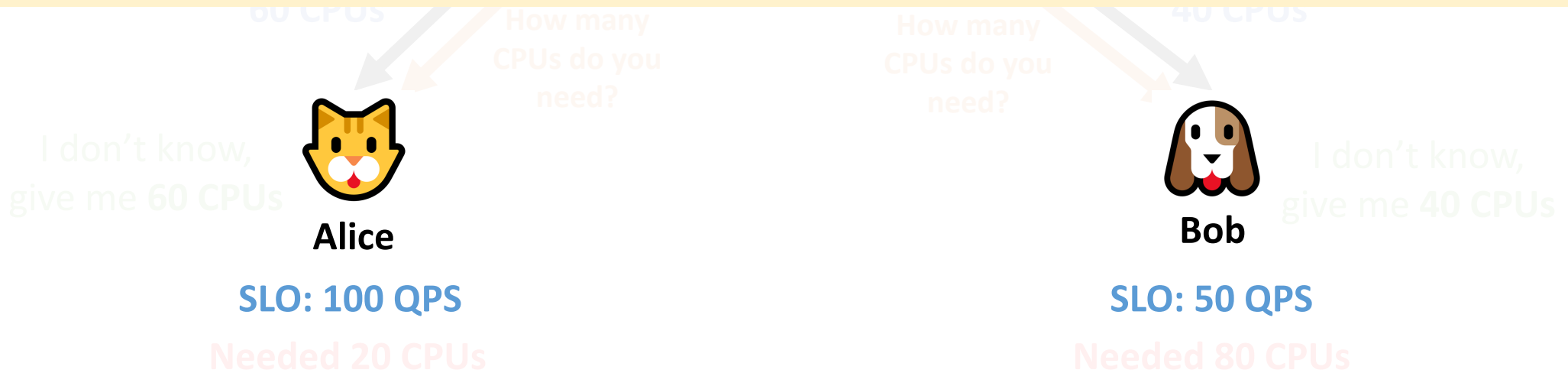


Asking for resource demands

Compute Resources

Users have performance objectives, not resource demands

How can we do performance-aware scheduling?

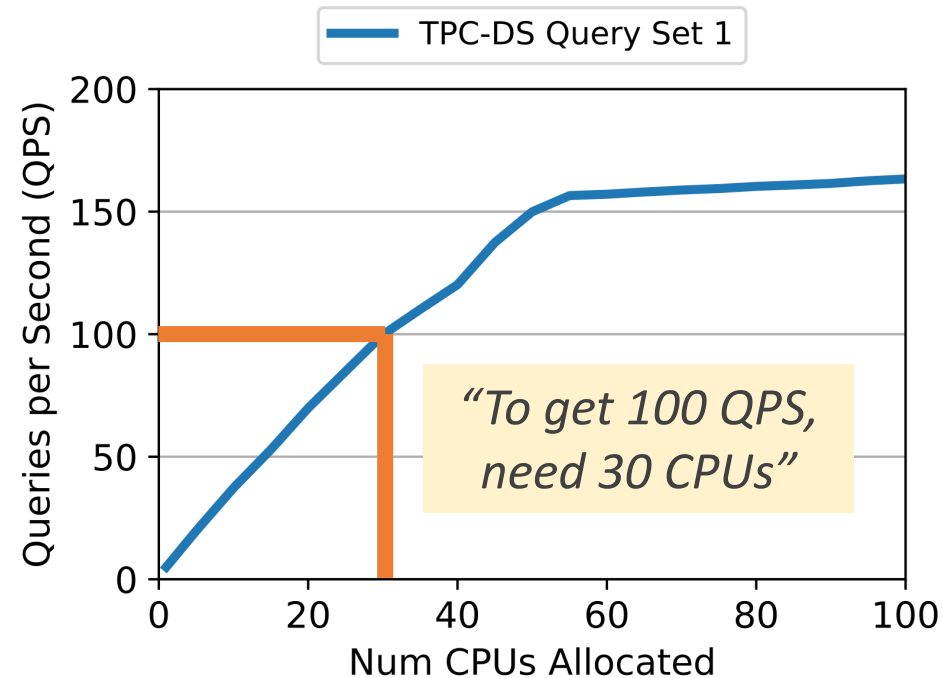


Towards performance-aware scheduling

- Users know their **performance objectives (SLOs)**, but not the resources required to achieve SLOs
- A **resource-performance mapping** can map SLOs to resource demands



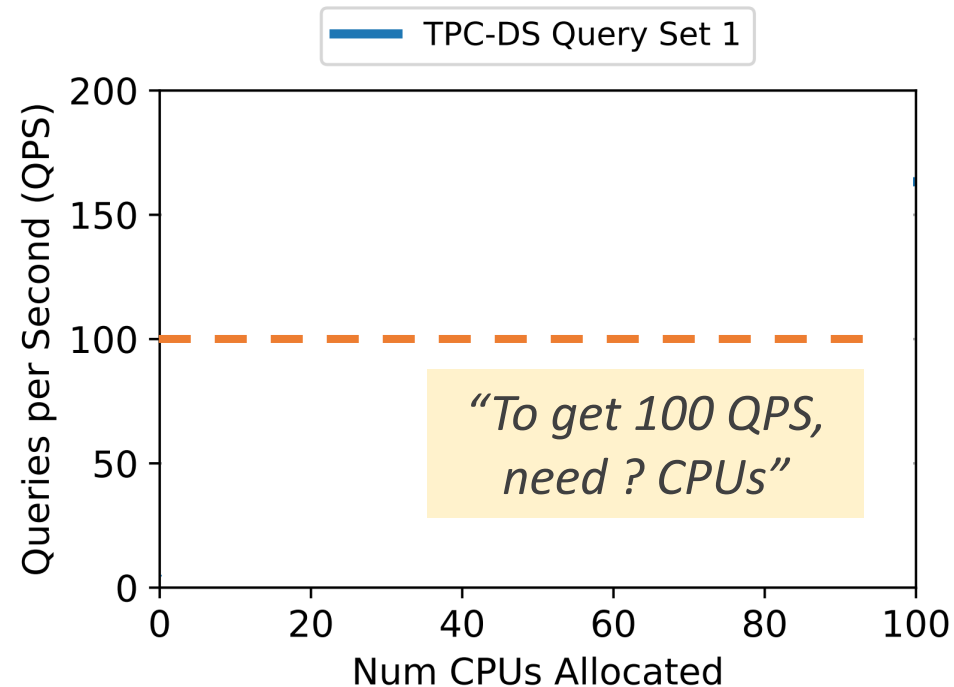
Alice
SLO: 100 QPS



Towards performance-aware scheduling

Challenge 1 - Resource-performance curves are not known

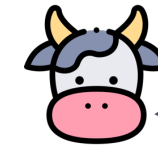
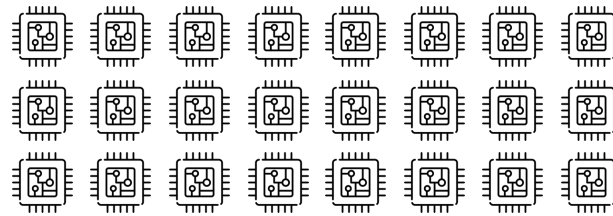
- A resource-performance mapping can map SLOs to resource demands



Resource allocation objectives are diverse

Compute Resources

120 CPUs



Admin

Objective

Maximize the average QPS over all users

Challenge 2 – Scheduling objectives are diverse and require a general resource allocator to satisfy objectives

I need **20 CPUs**
for 100 QPS



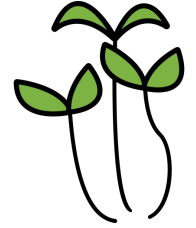
Alice



Bob

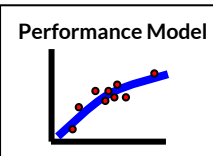
I need **80 CPUs**
for 50 QPS

Cilantro



- Cilantro is a **general framework** to enable **performance-aware scheduling** in clusters

Challenge 1 – Resource-performance curves are unknown



Use **online learning** to generate **progressively accurate** resource-performance models

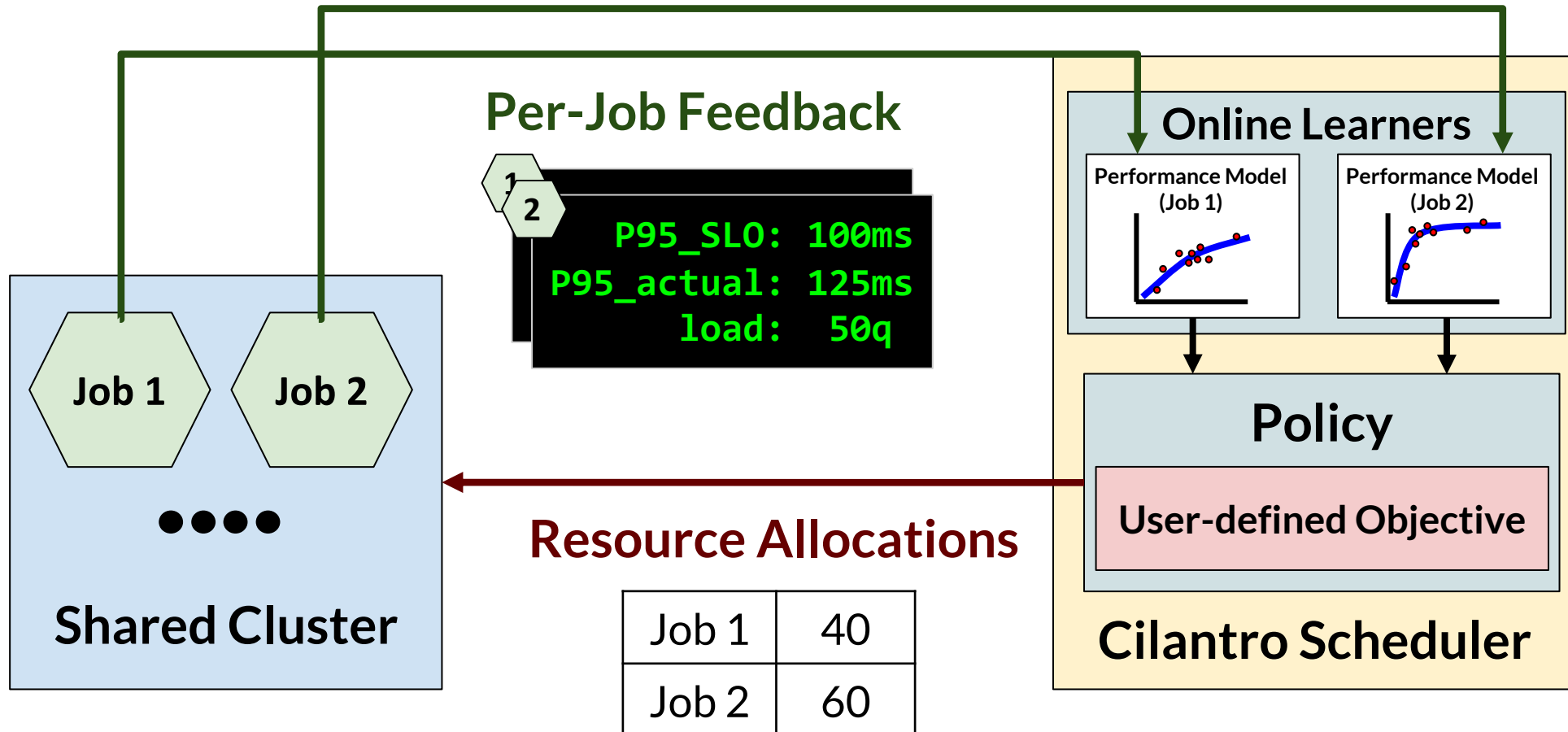
Challenge 2 - Support a diverse set of objectives



Decouple allocation policy from learning mechanisms by using resource-performance models

Cilantro workflow

Creating a closed scheduling loop of resource allocations and job feedback



Challenges in Online Learning



Large and expensive search space



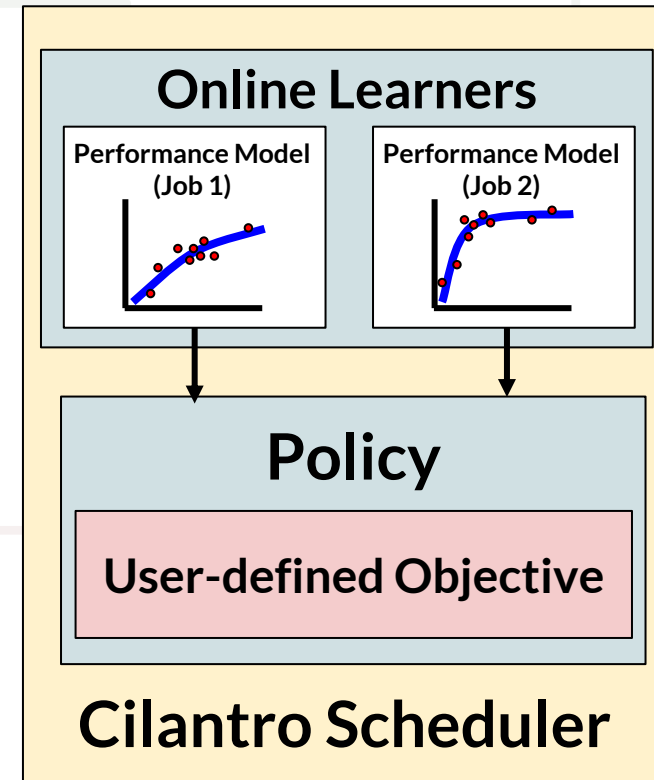
Reduce search space by defining a constrained utility model



Model estimates have uncertainty



Make policies uncertainty-aware



Cilantro utility model

Each job has an unknown resource-performance curve

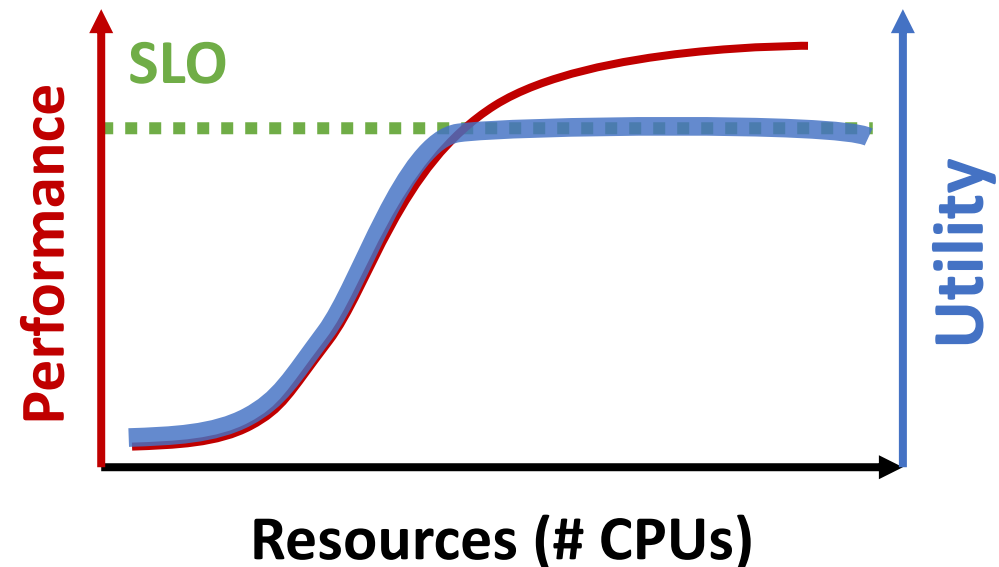
SLO (Service level objective) is the performance goal specified by the user

Utility is a function of performance. It must be:

1. **Monotonically non-decreasing**
2. **Clipped at SLO**

These properties are desirable because

1. **Makes exploration tractable**
2. **Reflective of real-world “utility”**



Common Scheduling Objectives

General Performance objectives

Social Welfare

"Maximize mean utility across all jobs"

$$a = \operatorname{argmax}_{a_j} \frac{1}{n} \sum_{j=1}^n u_j(a_j, l_j)$$

Egalitarian Welfare

"Maximize minimum utility across all jobs"

$$a = \operatorname{argmax}_{a_j} (\min(u_j(a_j, l_j)))$$

Application Performance objectives

Minimize end-to-end latency
of a collection of microservices

$$a^r = \operatorname{argmax}_{a \in A^r} (u(a, l))$$

Demand-based objectives

No Justified Complaints (NJC) Fairness
"Guarantee at least fair-share"

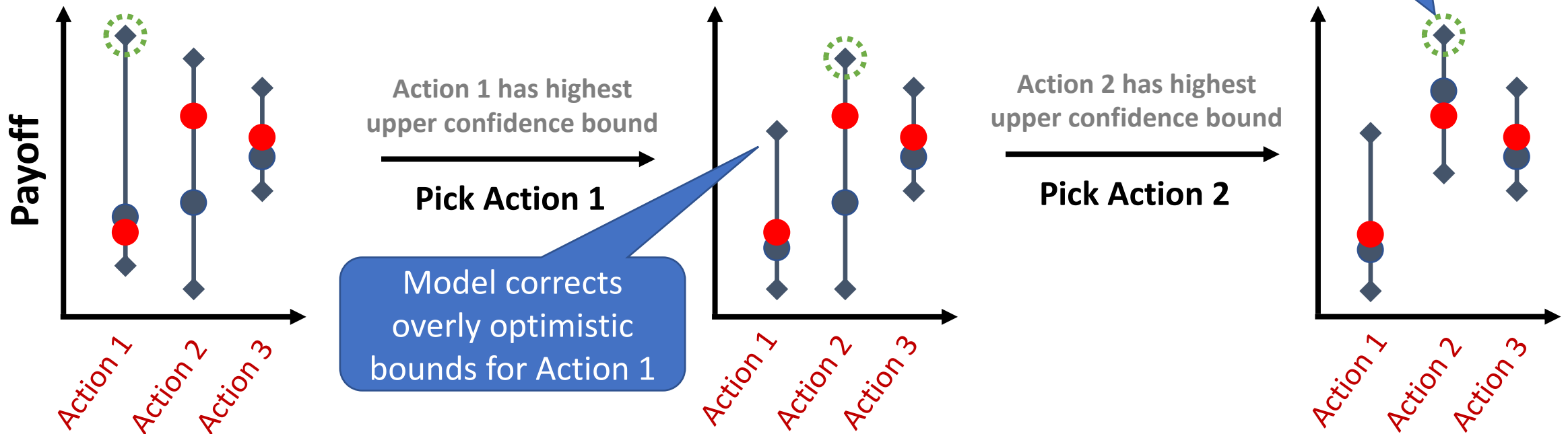
$$a = \operatorname{argmax}_{a_j} \left(\min_{j \in \{1, \dots, n\}} \left(\frac{u_j(a_j, l_j)}{u_j(R/n, l_j)} \right) \right)$$

Making policies uncertainty-aware

“Optimism in the Face of Uncertainty”^[1]
(OFU) principle

Prioritize actions with the highest potential utility

- Real payoff
- Expected payoff
- ◆ Confidence bounds



Making policies uncertainty-aware

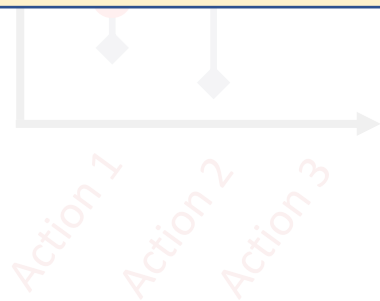
“Optimism in the Face of Uncertainty”^[1]
(**OFU**) principle

Prioritize actions with the highest potential utility

- Real payoff
- Expected payoff
- ◆ Confidence bounds

Model updates
bounds for Action 2

OFU allows policies to use **confidence bounds** to
balance explore-exploit tradeoff



OFU-Adapted Objectives

Performance-based objectives

Social Welfare

“Maximize mean utility across all jobs”

$$a = \operatorname{argmax}_{a_j} \frac{1}{n} \sum_{j=1}^n \hat{u}_j(a_j, \hat{l}_j)$$

Egalitarian Welfare

“Maximize minimum utility across all jobs”

$$a = \operatorname{argmax}_{a_j} (\min(\hat{u}_j(a_j, \hat{l}_j)))$$

Application-specific objectives

Minimize end-to-end latency
of a collection of microservices

$$a^r = \operatorname{argmax}_{a \in A^r} (\hat{u}(a, \hat{l}))$$

Demand-based objectives

NJC Fairness

“Guarantee at least fair-share”

Details in the paper!

Applications of Cilantro

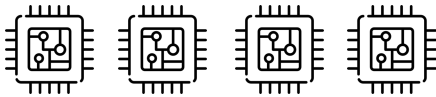
Multi-tenant Cluster Sharing

Allocating resources to independent jobs to maximize stated objective

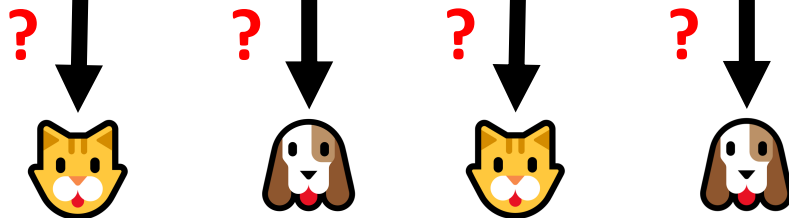
Objective



Resources



Resource Allocator



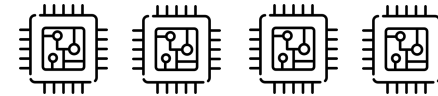
In the paper

This Talk

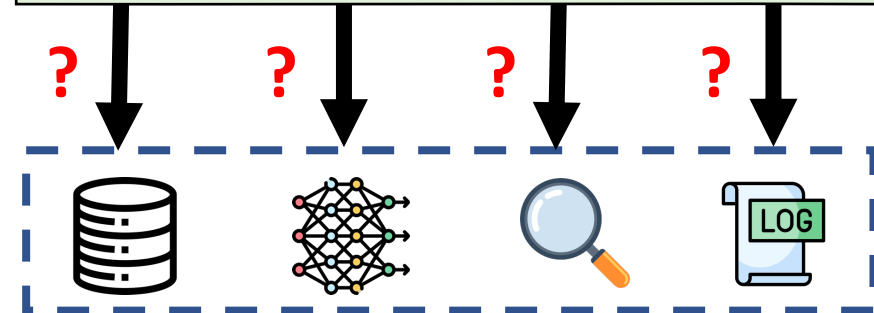
Microservices Resource Allocation

Allocating resources to microservices to minimize end-to-end application latency

Resources



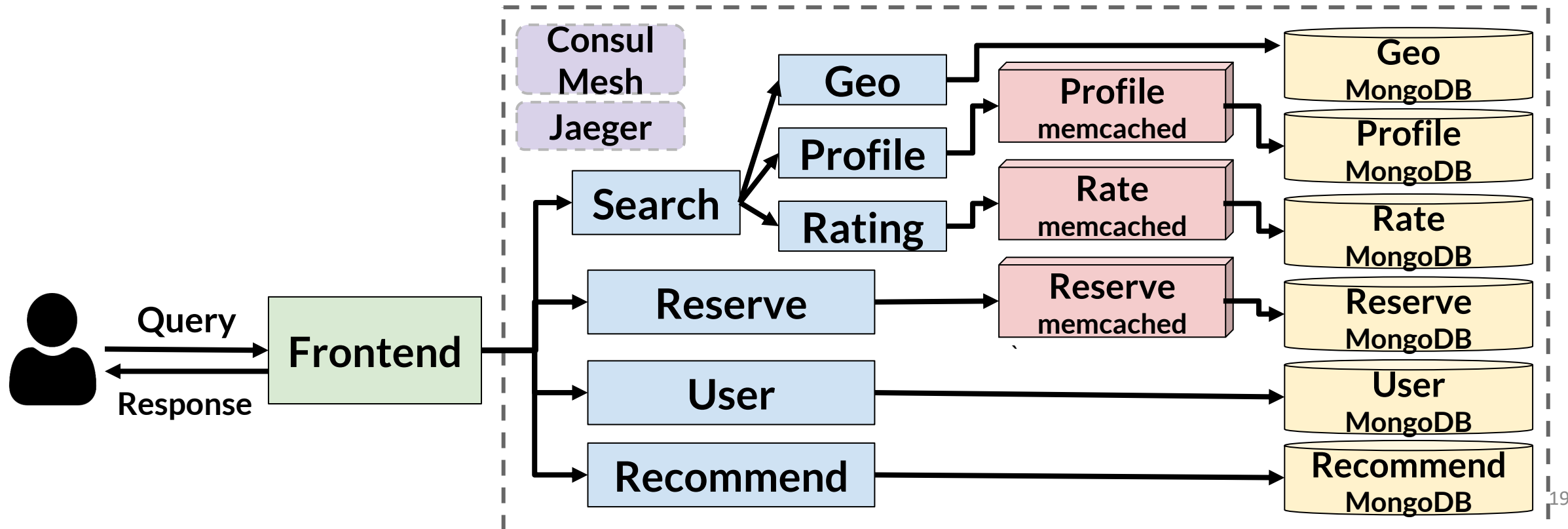
Resource Allocator



Application

Evaluation – Microservices

- Hotel Reservation benchmark from DeathStarBench
- Contains 19 microservices serving one endpoint

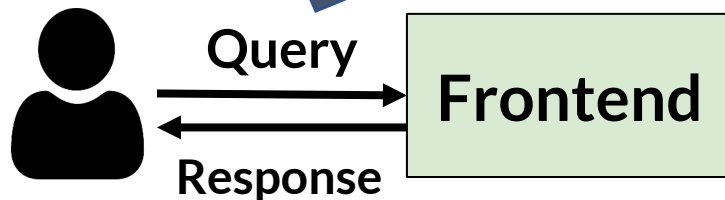


Evaluation – Microservices

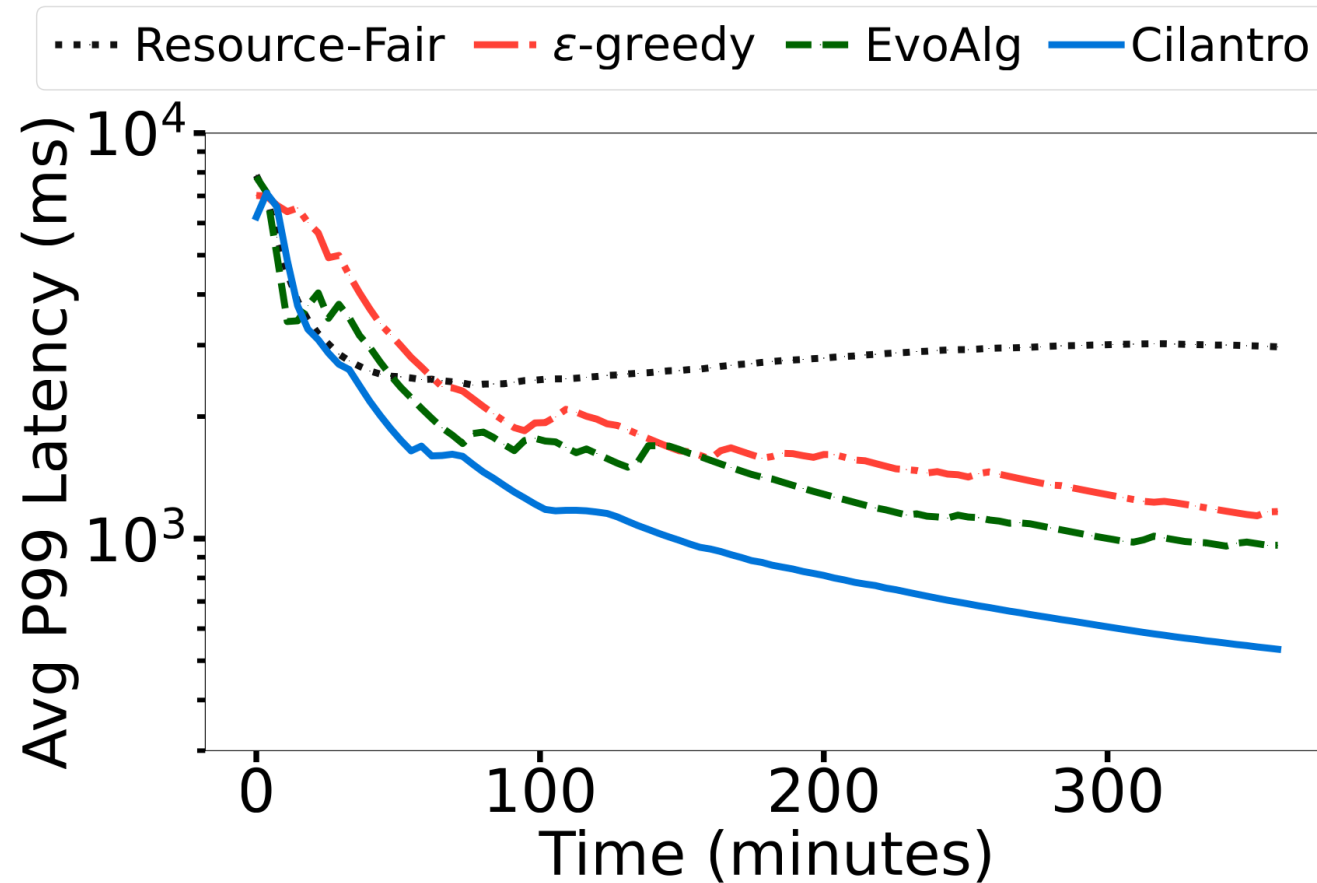
- Hotel Reservation benchmark from DeathStarBench
- Contains 19 microservices serving one endpoint

... to distribute 160 CPUs
between 19 microservices...

Cilantro observes end-
to-end query latency...



Evaluation - Microservices



Cilantro reduces P99 latency to 0.57x of the best baseline
without any visibility into the performance or structure of component microservices

Cilantro Summary



github.com/romilbhardwaj/cilantro

romilb@eecs.berkeley.edu

- Cilantro is a **general, performance-aware** cluster scheduling framework
- Users do not state their resource demands - they provide **SLOs and feedback**.
- Online learning estimates **resource-performance curves** and **policies adapt to uncertainty in estimates**
- Improves utilities up to 1.2 – 3.7× in cluster sharing and reduces latencies by 43% for microservices

