

Mu: An Efficient, Fair and Responsive Serverless

Framework for Resource-Constrained Edge Clouds

Viyom Mittal*, **Shixiong Qi***, Ratnadeep Bhattacharya⁺, Xiaosu Lyu⁺, Junfeng Li[§], Sameer G Kulkarni[†], Dan Li[§], Jinho Hwang^{*}, K. K. Ramakrishnan*, Timothy Wood⁺

*University of California, Riverside, ⁺George Washington University, [§]Tsinghua University,

[†]Indian Institute of Technology, Gandhinagar, ^{*}Facebook Inc.

November 1st, 2021



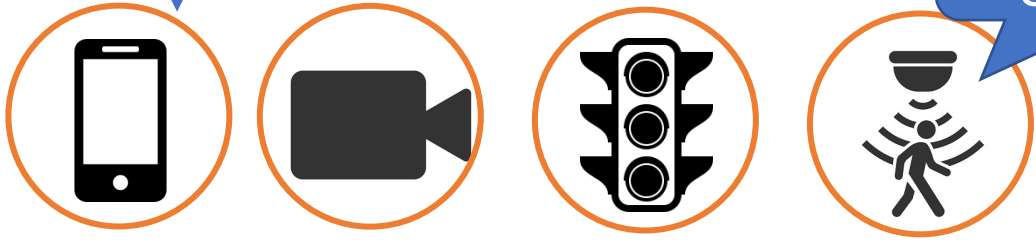
Serverless Cloud



*Current designs of serverless p...
yet a viable option for Edge*

🙄 resource constrained

Closer to the
end users



Challenges to Using Existing Approaches in Edge Clouds

Challenge 1: Imprecise Resource Provision

Existing autoscaling design depends on user inputs

- 🙄 Users are unaware of runtime features of functions
- 🙄 Inappropriate parameter choices

Single metric-based autoscaling (e.g., CPU utilization)

- 🙄 Insufficient resources to achieve optimal autoscaling

Slow resource provision in case of traffic bursts

- 🙄 Long response time, SLO violations

Missed SLO

Challenges to Using Existing Approaches in Edge Clouds

Challenge 2: Unfair Function Placement

Existing placement algorithms

- ☹️ With no constraints
- ☹️ Unfair resource allocation

Function-1: More resource provisioned and Better SLO

🤔 If two functions that have the same SLO are going to be placed...

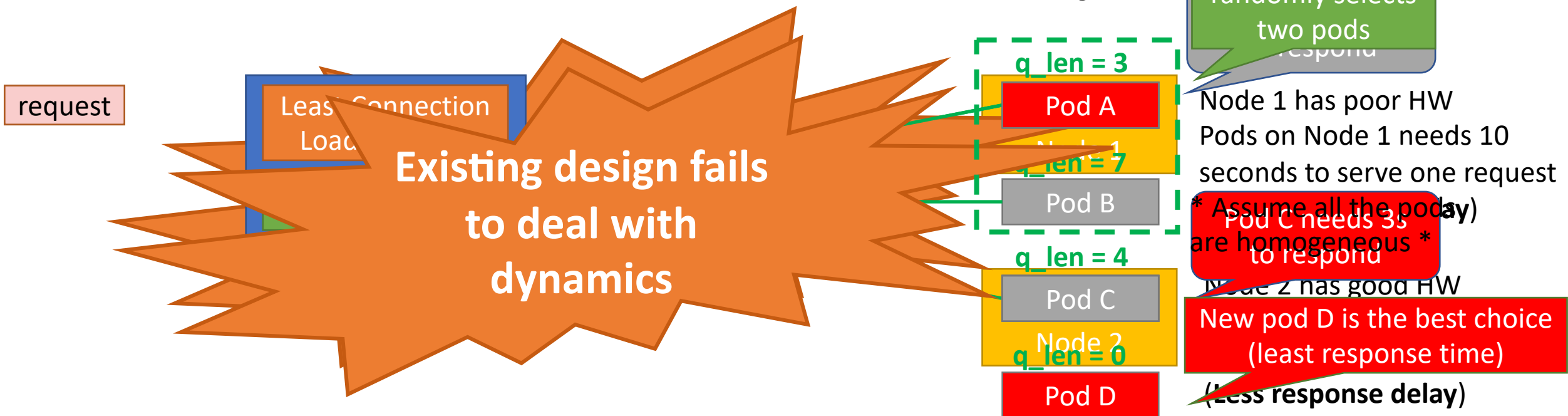
Function-2: Less resource provisioned and Worse SLO

Challenges to Using Existing Approaches in Edge Clouds

Challenge 3: Unawareness of Resource Heterogeneity and System Dynamics

☹️ **Resource Heterogeneity and System Dynamics can lead to poor load balancing decision**

Least connection LB: Track the queue length at backend pods and distribute the request to the pod with minimum queue length



Challenges to Using Existing Approaches in Edge Clouds

Challenge 4: Approximate metrics collection

Serverless platform relies on pod metrics to guide resource management

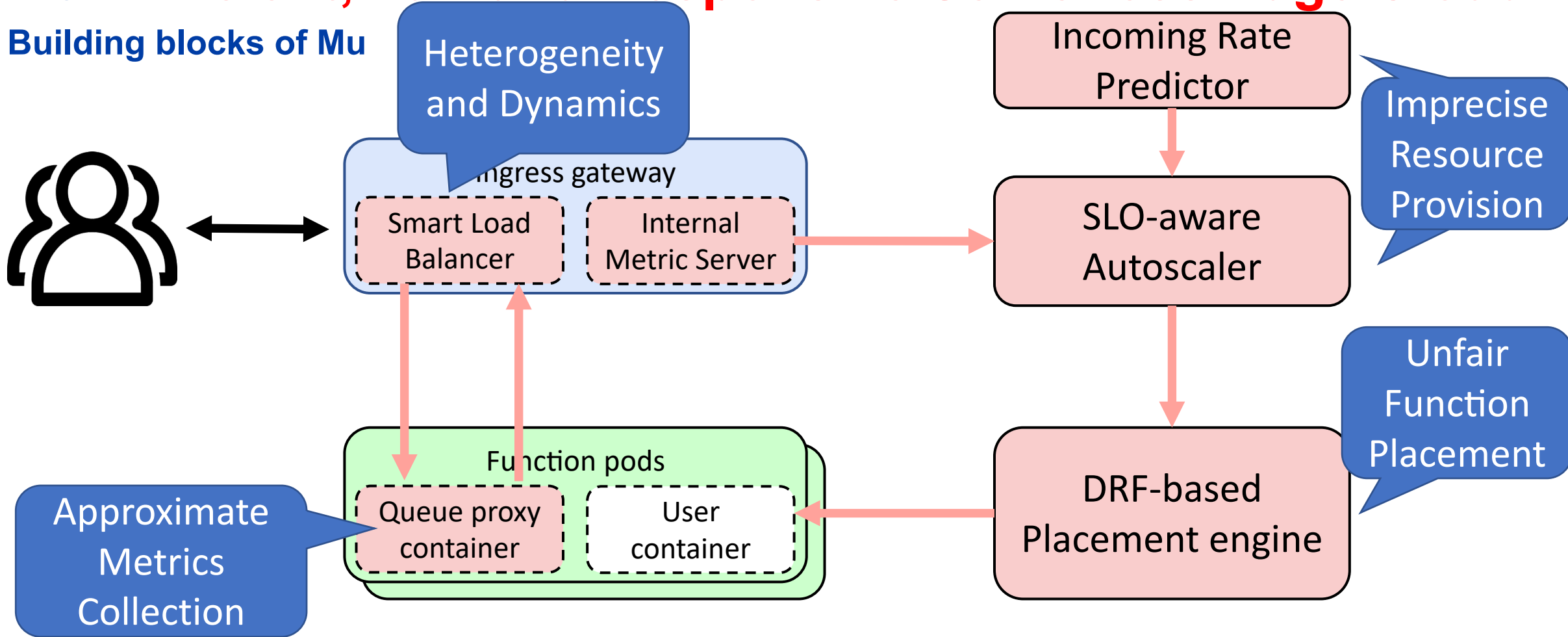
Existing design relies on approximate metrics collection to address scaling

- ☹️ an inaccurate view of system status
- ☹️ negative impact on resource provision, load balancing...

**Need a precise, lightweight
and scalable metric
collection mechanism**

Mu: Efficient, Fair and Responsive Serverless Edge Cloud

Building blocks of Mu



SLO-aware Autoscaler

Idea:

😊 From user's perspective, providing SLO is more meaningful rather than

inter (RPS)

✓ Single metric-based autoscaling (RPS, Concurrency)

✓ Existing autoscaling design depends on user inputs

How:

1 Users only provide the target SLO of their function

2 Provision resources by factoring in both the incoming request rate and the queue length

3 Ensure SLOs by factoring in the average request execution time

✓ Avoid over-allocation of resources to ensure performance with just the right amount of resources.

Incoming Rate Predictor

Our Goal:

🤔 Provision adequate resource in case of traffic bursts

Idea:

😊 Mu uses a simple online linear regression model to predict the incoming rate based on previous observations

😊 Mu uses multi-armed bandit to improve accuracy

Features:

1. lightweight and fast
2. Accurate prediction - dynamically select the model with minimum error

Placement Engine

2-stage heuristic algorithm

- **Resource fairness** between serverless functions
 - Function selection based on Dominant Resource Fairness (DRF)
- **Resource efficiency** between nodes
 - Node selection based on scoring
 - Alignment [1], WorstFit [2], and BestFit [2].
 - reduce the resource fragmentation, minimize unfairness

*Call two stages iteratively until there are no resources left or all functions are placed

Consider fairness and efficiency together

1. R. Grandl, G. Ananthanarayanan, S. Kandula, S. Rao, and A. Akella. 2014. Multi-resource packing for cluster schedulers. ACM SIGCOMM Computer Communication Review 44, 4 (2014), 455–466.

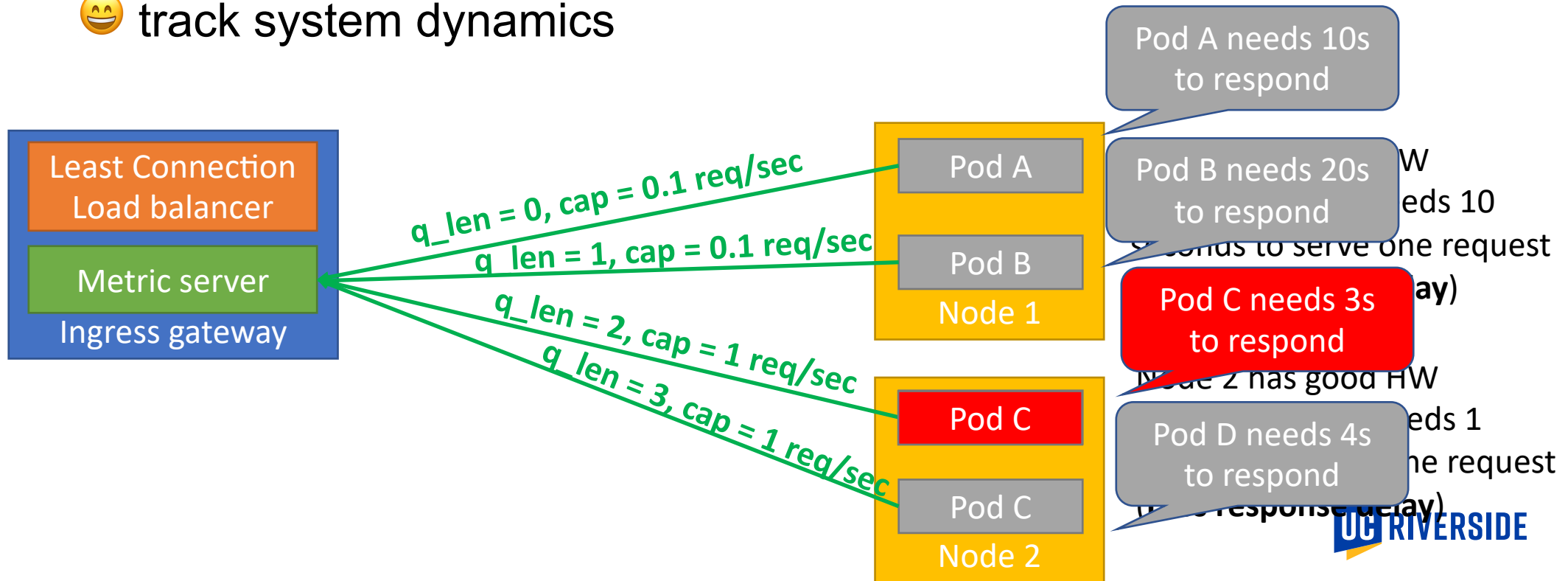
2. C. A. Psomasand, J. Schwartz. Beyond beyond dominant resource fairness: Indivisible resource allocation in clusters. Tech Report Berkeley, Tech. Rep. (2013).

Load Balancer

Our Goal: 🤔 to be aware of resource heterogeneity and system dynamics

- 😊 use extra metrics to estimate response time of each pod
 - 😊 differentiate “fast” and “slow” pods in the system
- 😊 use “piggybacked” metrics of each pods instead of “two random choices”
 - 😊 track system dynamics

request



Metric collection

Our goal:

- Precise metric collection that can reflect latest system state and achieve good scalability

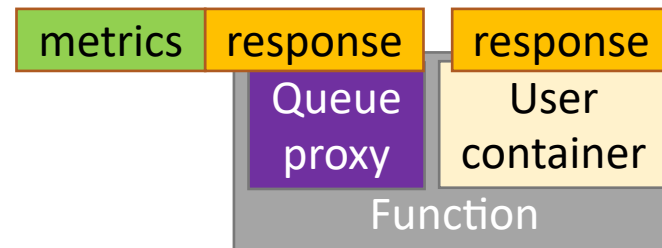
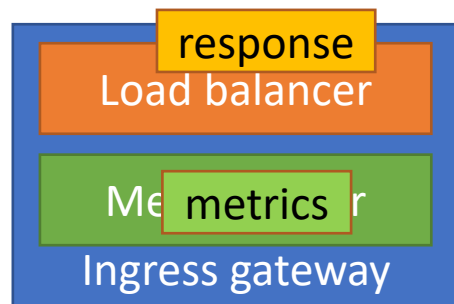
Precise

Dynamics

Heterogeneity

Scalable

- The response is instantaneous providing metrics based for the latest system state
- Piggybacking also help us accommodate heterogeneity and system dynamics
- Piggybacking eliminates the explicit need of metric collection



Summary of overall evaluation

Experiment setup

- Three different system configurations

Mu

Default Knative with RPS autoscaling (RPS)

Default Knative with Concurrency autoscaling (CC)

- Scaling target is fair among different scaling policies
- System is slightly overloaded
- Evaluated metrics
 - Latency and Fairness
 - Pod allocation (efficiency)
 - SLO performance

Azure workloads

Parameter/Specification	Values	
Invocation Range	W-1	41-230 rps
	W-2	69-182 rps
Average invocations	W-1	154 rps
	W-2	146 rps
Container Concurrency	4	
Grace Flag (Mu only)	16	
Execution time	500ms	
Maximum pod capacity	48	
CPU and Mem. per pod	7 cores, 30GB	
Target	RPS	8
	CC	40
SLO	5 seconds	

Latency and Fairness

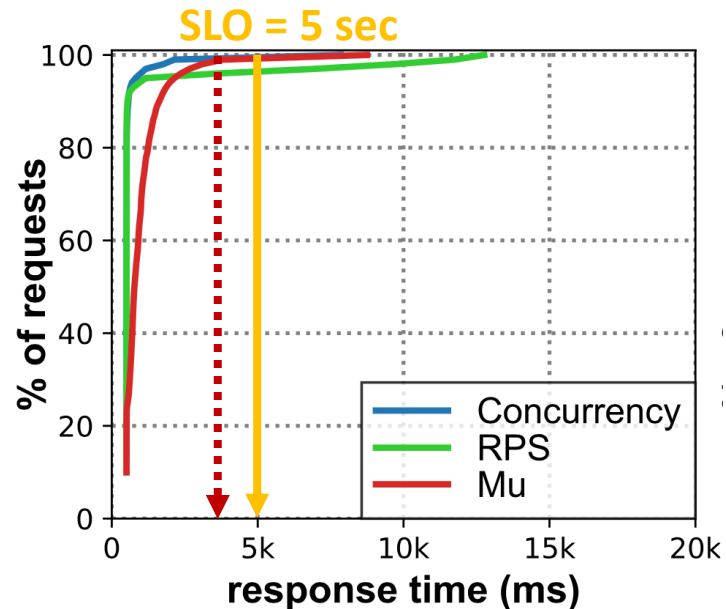
Response time CDF for **Mu** and standard Knative approaches (**CC** and **RPS**)

Mu has good control over response time

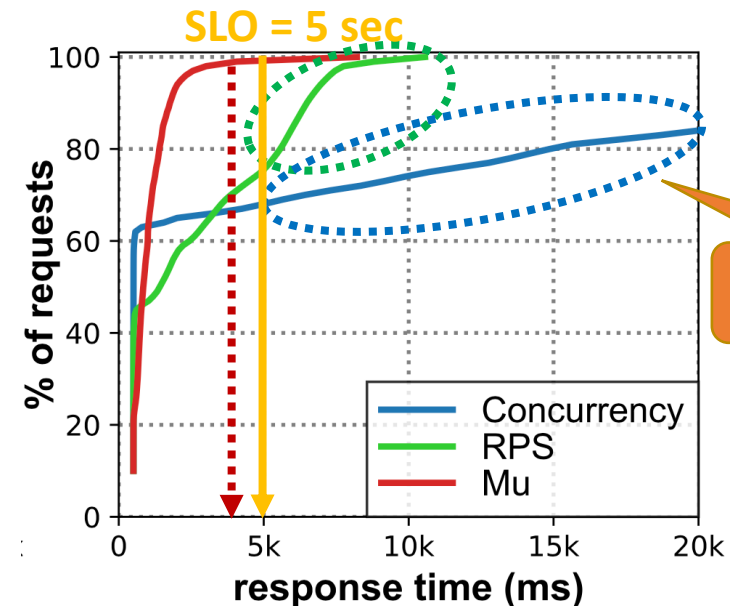
- Mu limits the tail latency with SLO of 5 seconds **for both workloads**
- Standard Knative approaches result in much larger response time tail

Mu achieves better fairness

- Mu treats **Workload-1** and **Workload-2** equally
- Standard Knative approaches unfairly treat **Workload-2**



Workload-1

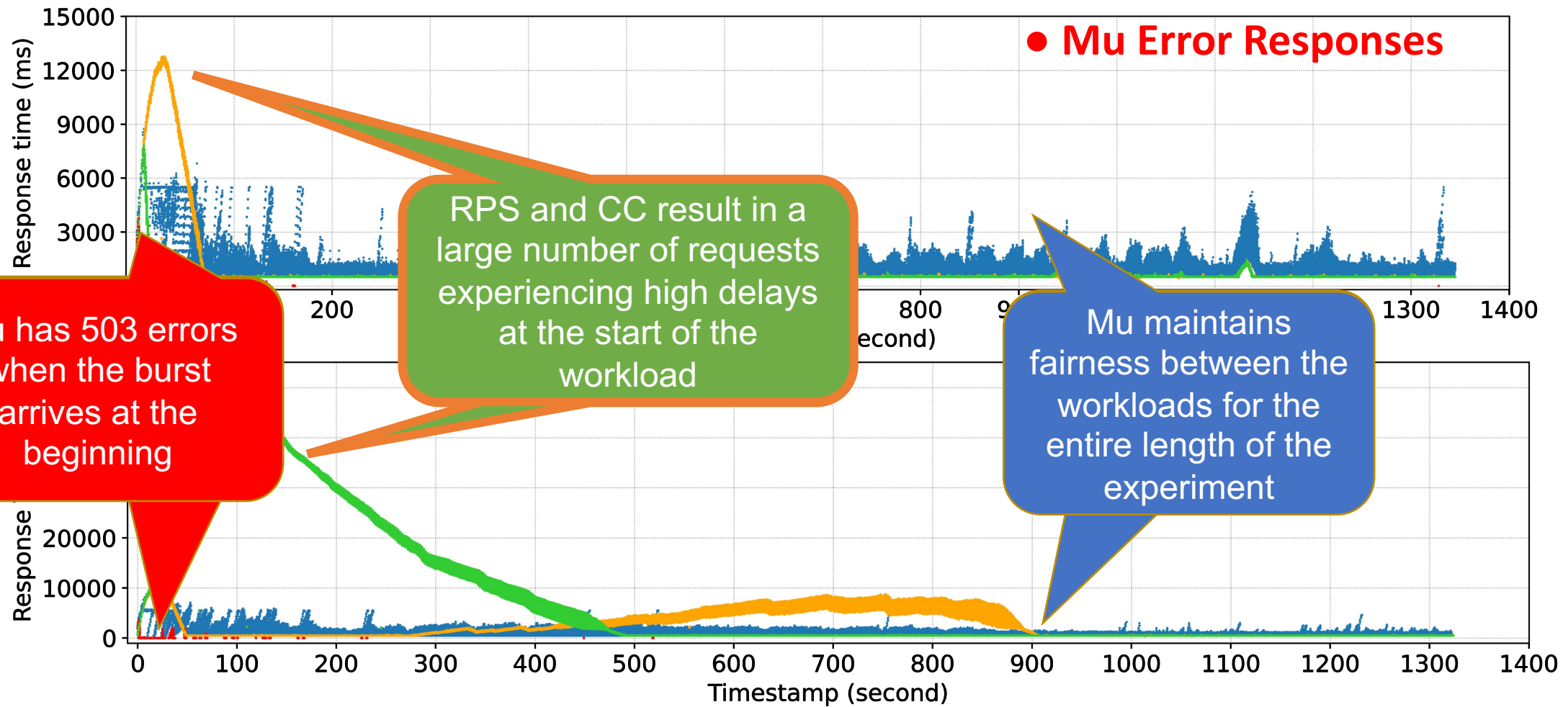


Workload-2

Latency and Fairness

Time series of Response Time for Mu and standard Knative approaches

- RPS Successful Responses
- CC Successful Responses
- Mu Successful Responses
- Mu Error Responses



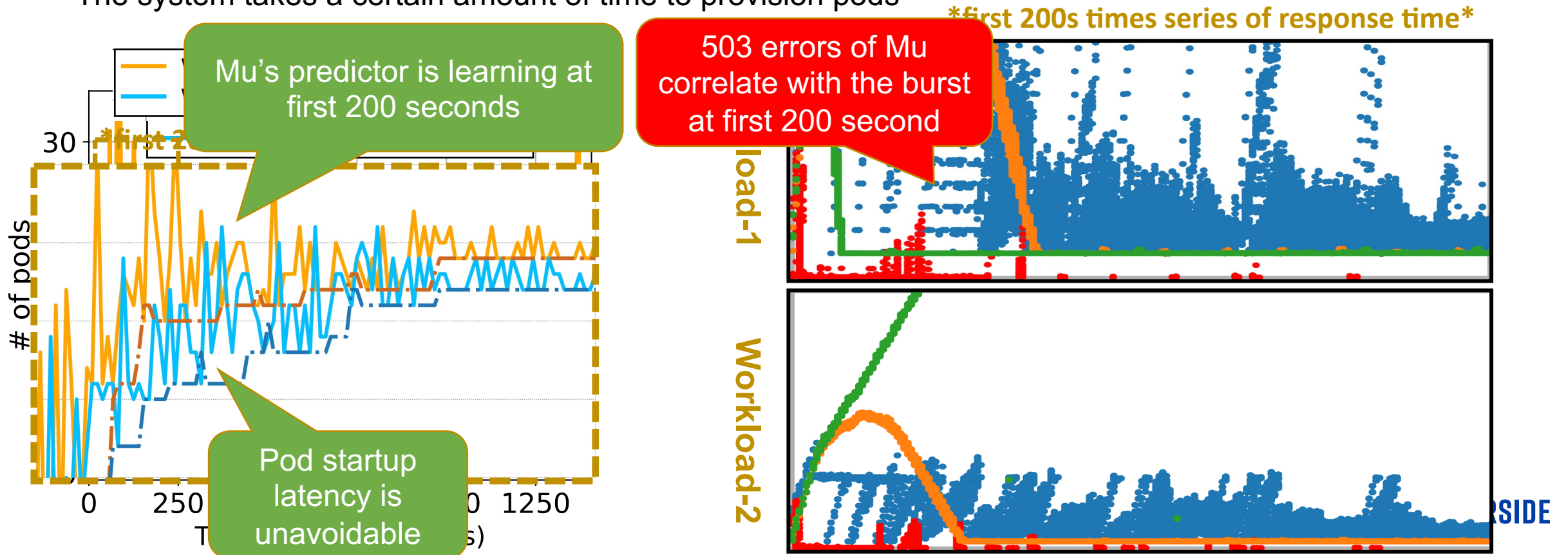
Workload-1
Workload-2

SLO Performance

Mu: Correlation between 503 errors and occurrence of bursts

Most of the 503 errors occur when the burst arrives at the beginning (see first 200 seconds)

- When the predictor has not yet learned the characteristics of the workload
- The system takes a certain amount of time to provision pods



SLO Performance Summary

Mu achieves better SLO than standard Knative approaches

- 96.8% requests served within SLO compared to RPS (86.2%) and CC (84.2%)

Mu uses SLO-aware admission control and returns 503 errors

- This avoids the build up of a large queue with the arrival of a *burst* of requests
- RPS and CC choose to buffer the burst of requests -> SLO violation

503 errors result in limited negative impact and bring more SLO benefit compared to a large queue

- Better SLO performance
- Low tail latency

			503 error / total requests	Requests served within SLO
CC	Workload-1	3805	503 / 3805	86.5%
	Workload-2	1073	503 / 1073	53.0%
	Workload-3	1757	503 / 1757	71.1%
RPS	Workload-1	8808	0 / 8808	100.0%
	Workload-2	209905	0 / 209905	100.0%
Concurrency	Workload-1	2141	0 / 220126	99.6%
	Workload-2	49526	0 / 209905	68.0%

503 errors returned by Mu impacts relatively small number (<5%) of requests,

CC and RPS build up a large queue resulting in very long latencies

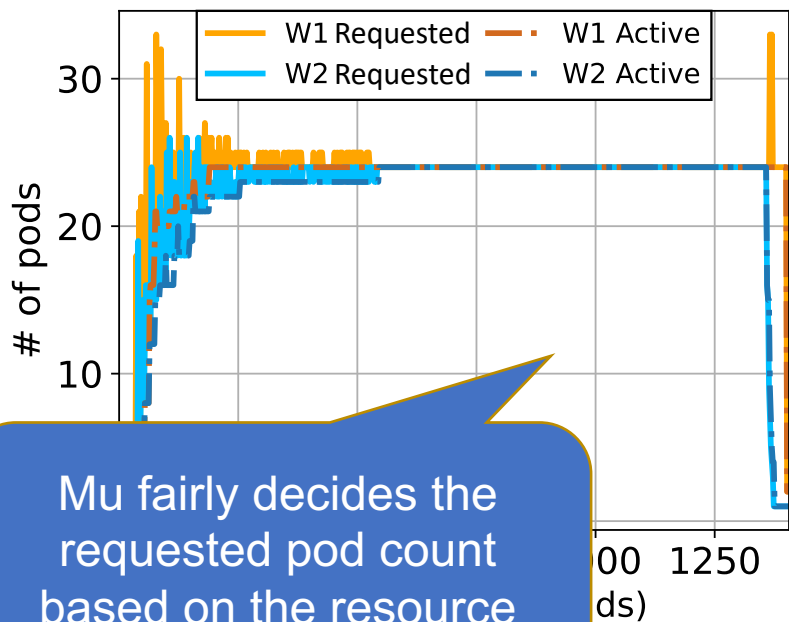
25-30% of total requests failed to meet SLO

Pod Allocation

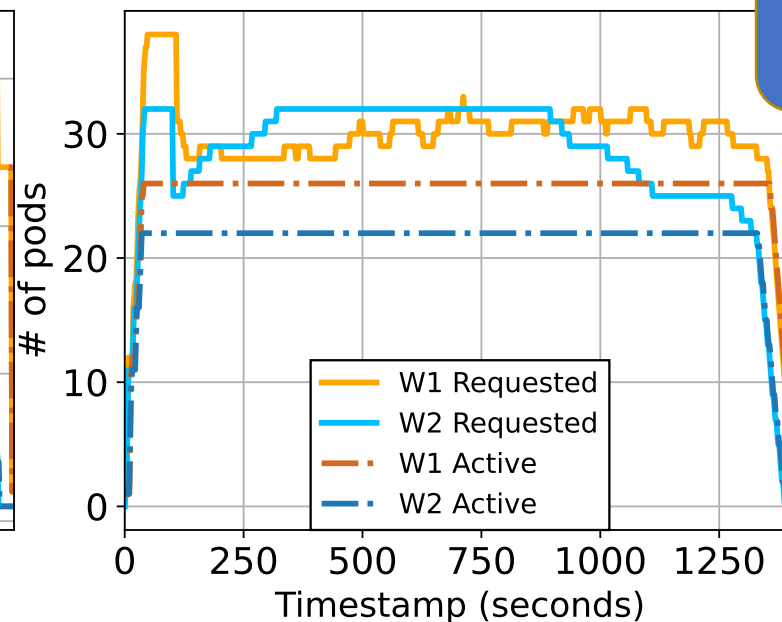
Time series of Pod counts for Mu and standard Knative approaches

Mu uses less pods than standard Knative approaches

- On average RPS and Concurrency use 18% and 17% more pods than Mu
- Mu tends to request fewer pods since its goal is to meet SLOs, not necessarily minimize response times, and its predictor helps it judge the workload.

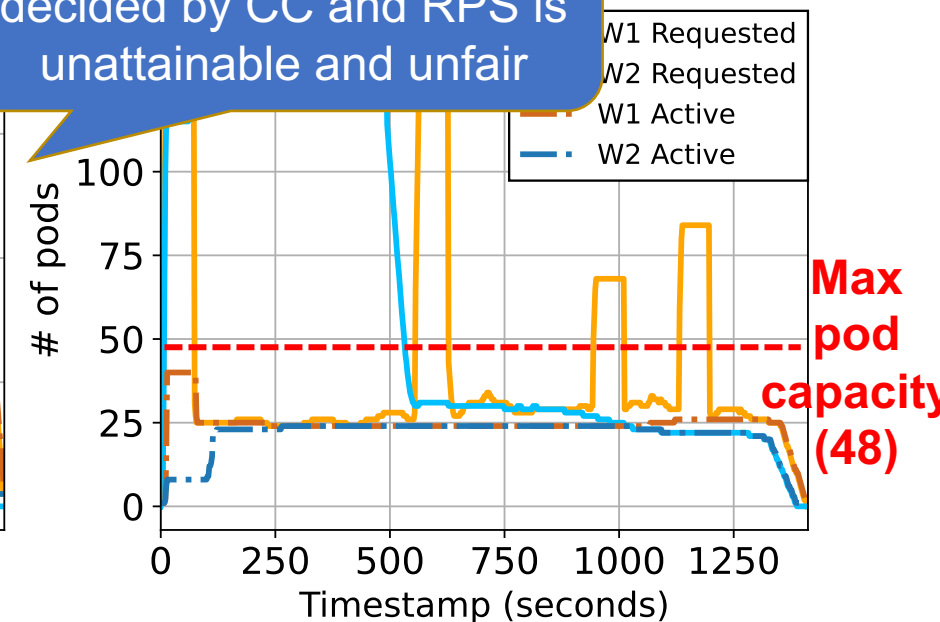


Mu fairly decides the requested pod count based on the resource availability.



RPS-based Knative

Requested pod count decided by CC and RPS is unattainable and unfair



Concurrency-based Knative

Max pod capacity (48)

Conclusion

Mu achieves better

SLO-aware Autoscaler

Incoming rate predictor

01 latency performance

DRF-based placement engine

02 resource fairness

Smart load balancer

03 resource efficiency

Piggybacked metrics

04 SLO performance

compared to Knative

An Efficient, Fair and Responsive Serverless Framework for Resource-Constrained Edge Clouds