

Escra: Event-driven, Sub-second Container Resource Allocation

Greg Cusack, Maziyar Nazari, Sepideh Goodarzy, Erika Hunhoff, Prerit Oberai, Eric Keller,
Eric Rozner, Richard Han

ICDCS - July 11, 2022
Bologna, Italy

vmware[®]





Containerized Infrastructure

- Light-weight
- Rich orchestration systems
- Development workflow integration
- Strong resource isolation

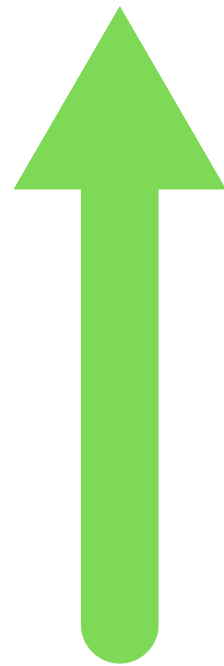


Developers must make tradeoffs!

Developers must make tradeoffs

Prioritize Performance

Overallocate compute resources



Performance



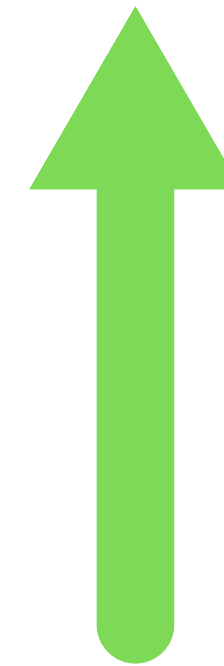
Efficiency

Prioritize Resource Efficiency

Underallocate compute resources



Performance



Efficiency



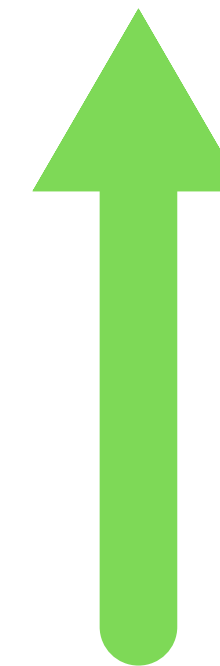
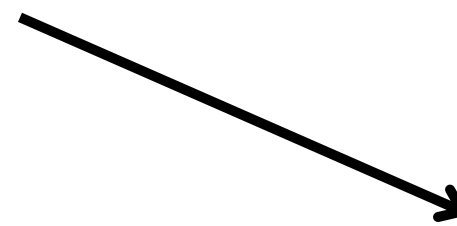
Developers must make tradeoffs

Prioritize Resource Efficiency

Underallocate compute resources

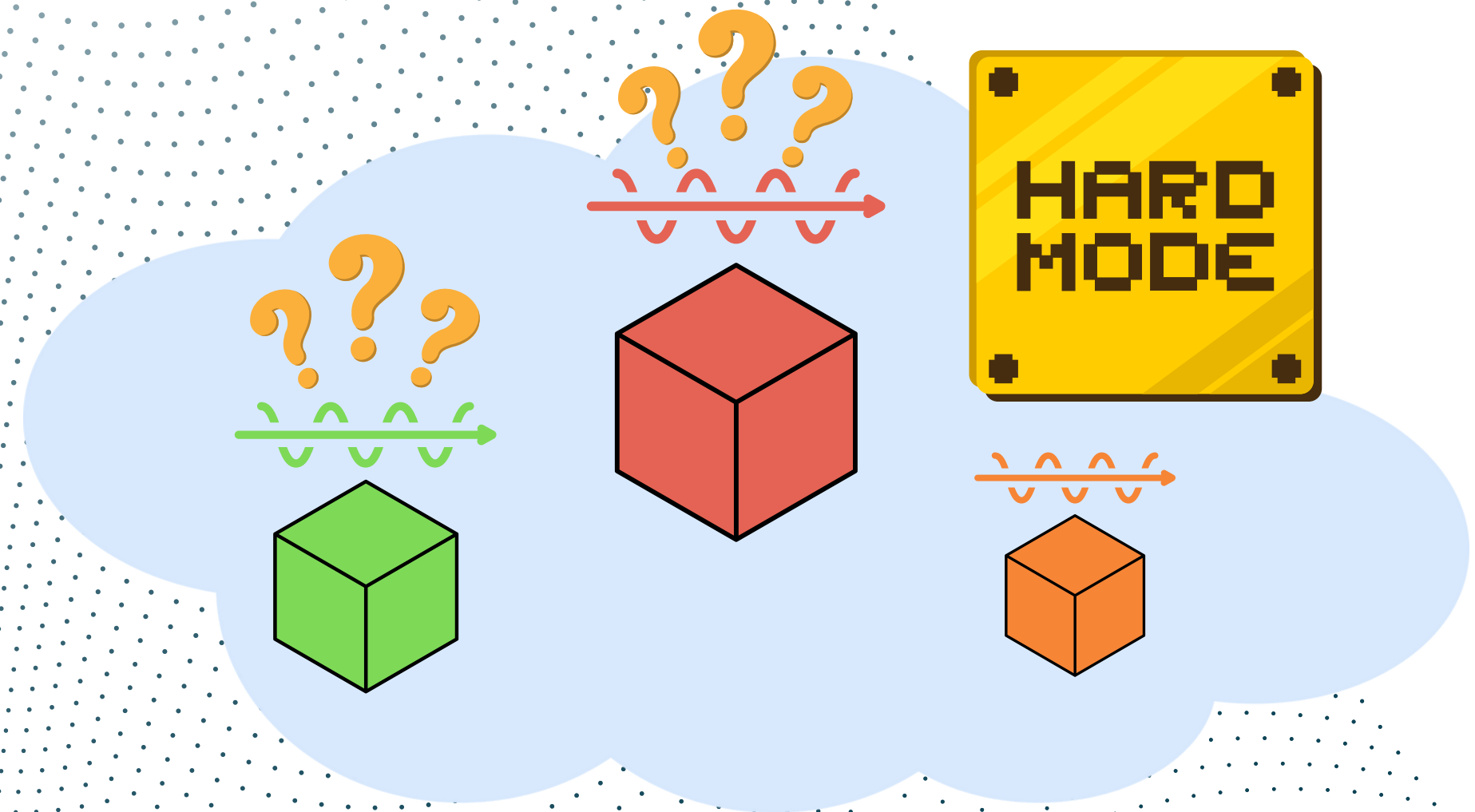
Out of Memory Events (OOMs)

CPU Throttles



Performance

Efficiency



Setting Container Limits is Difficult!

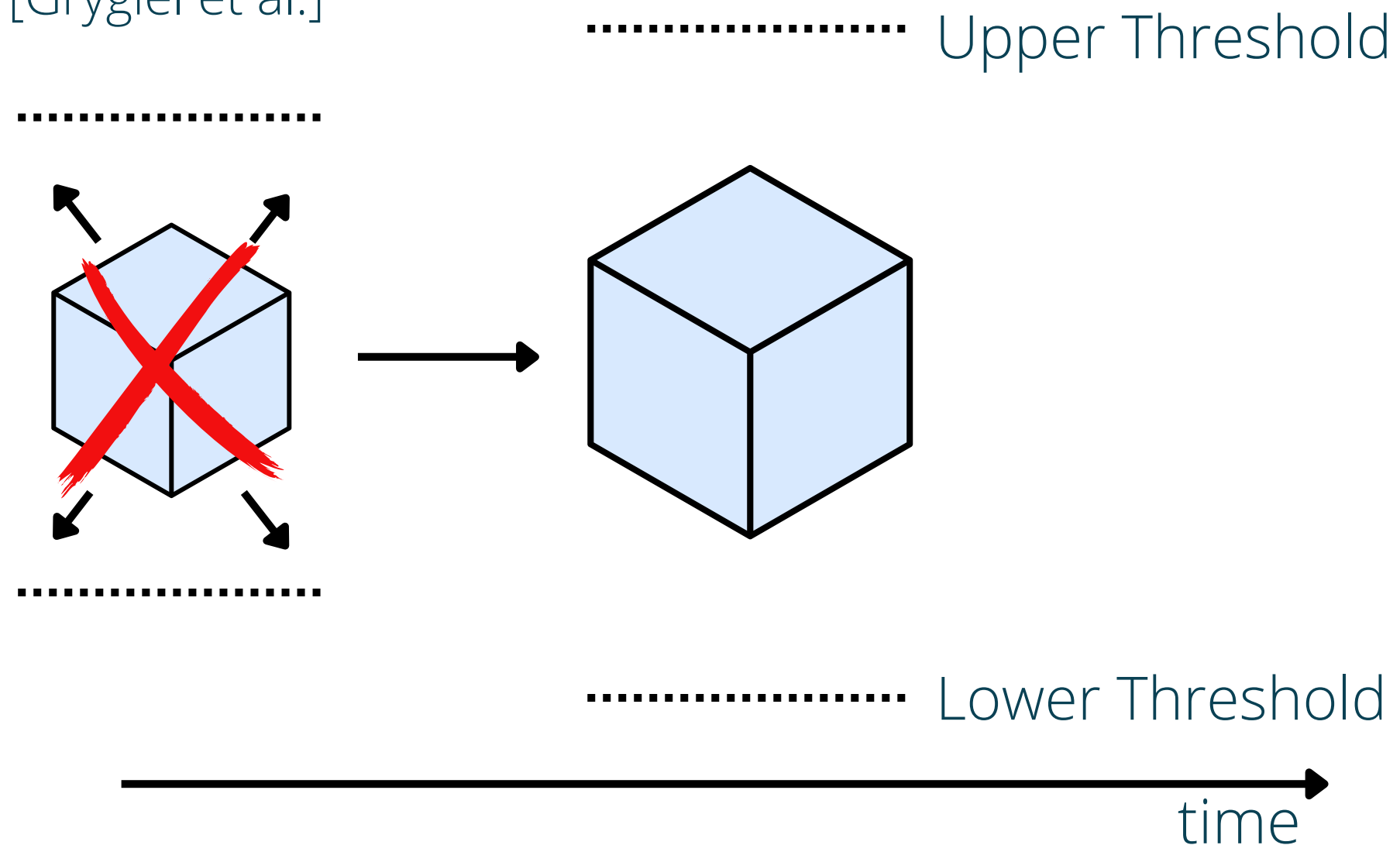
- Requires representative workload
- Resource needs vary over time
- Tools aggregate usage over coarse-grained timescales
- **Result: Over provisioned containers**

Recent Work

Threshold Scaling

VPA

[Grygiel et al.]

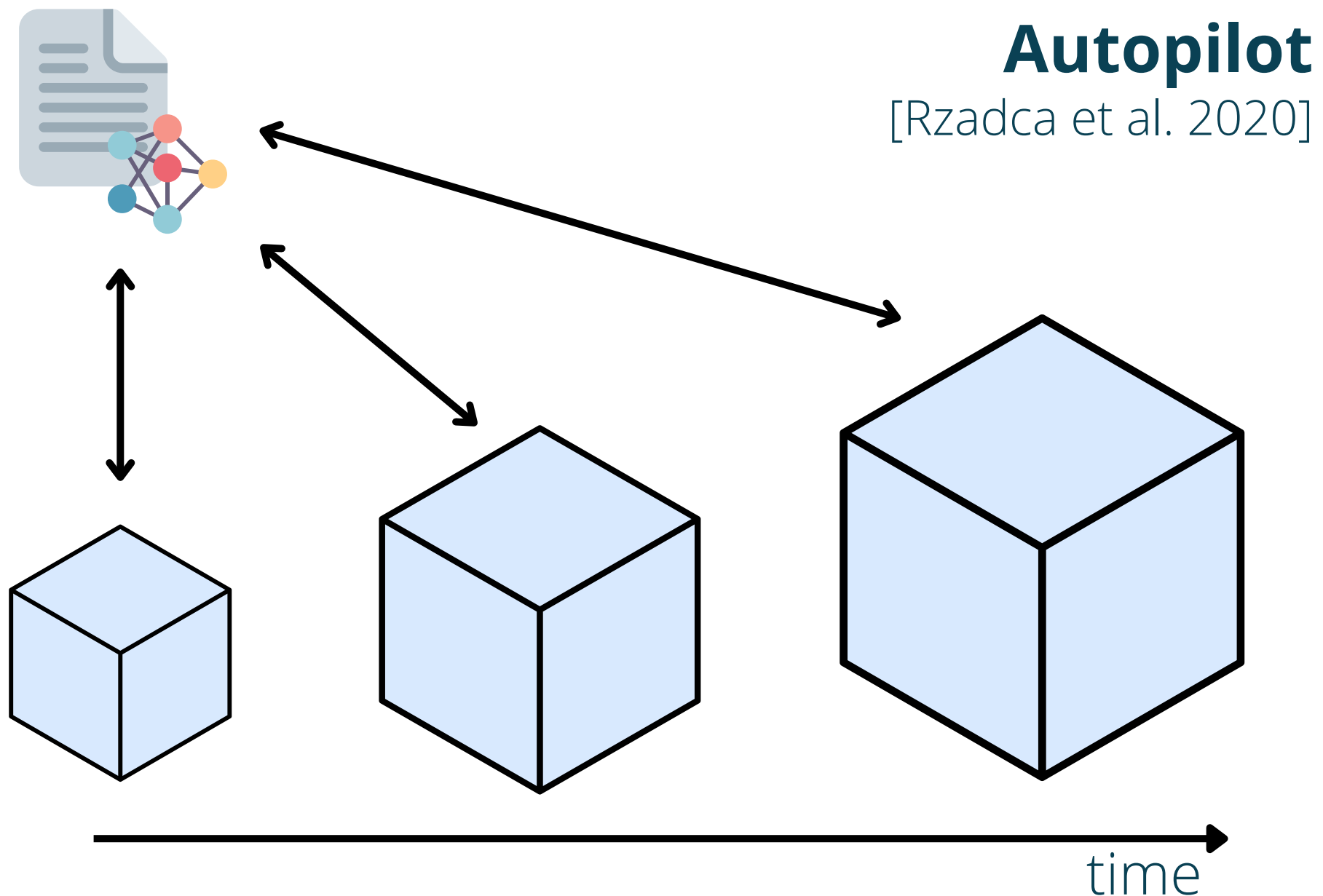


Shortcomings

- Container restart to scale
- Scale once per minute
- High slack

Recent Work

ML-based Scaling



Shortcomings

- Cannot respond to quick changes in workload
- Must allocate to maximum prediction over next 5-minute interval
- Unable to correct inaccurate predictions
- Not suitable for short-lived containers (e.g. serverless)

Key Insight

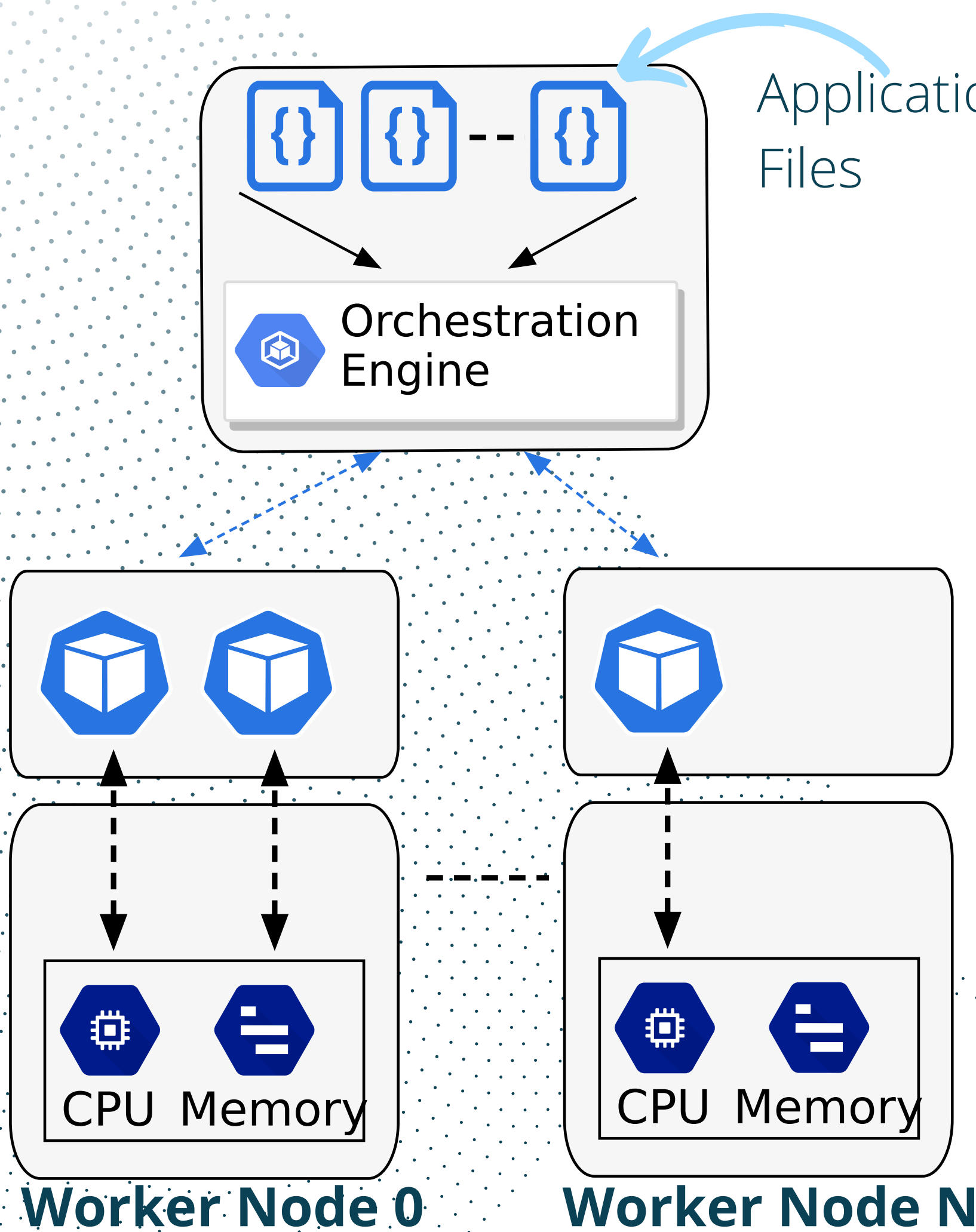
If we can make resource allocation fine grained, we don't need to predict; we can react to workload changes.





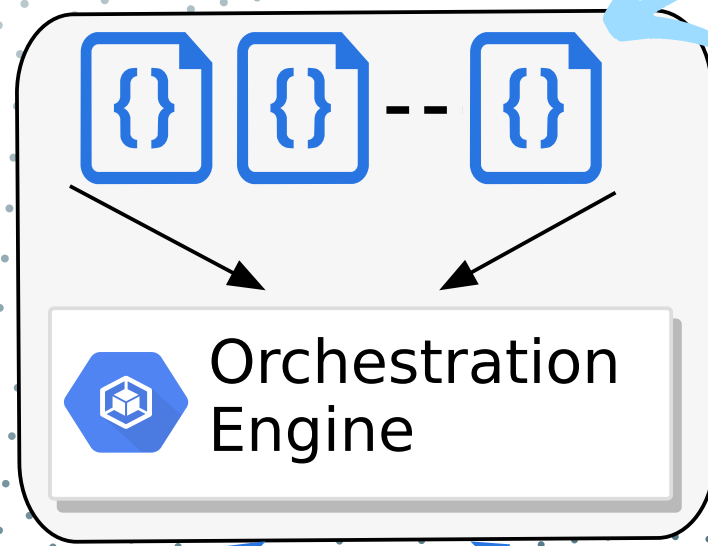
Escra

- Distributed resource allocation
- Sub-second interval scaling within and across hosts
- OOM prevention and scaling
- Immediate CPU throttle response
- No performance penalty



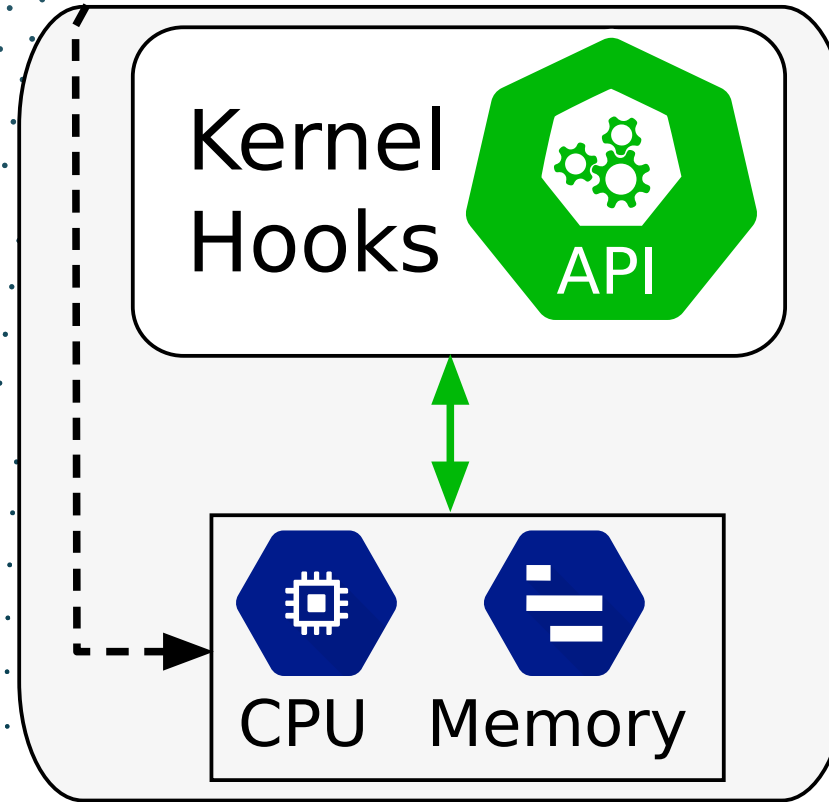
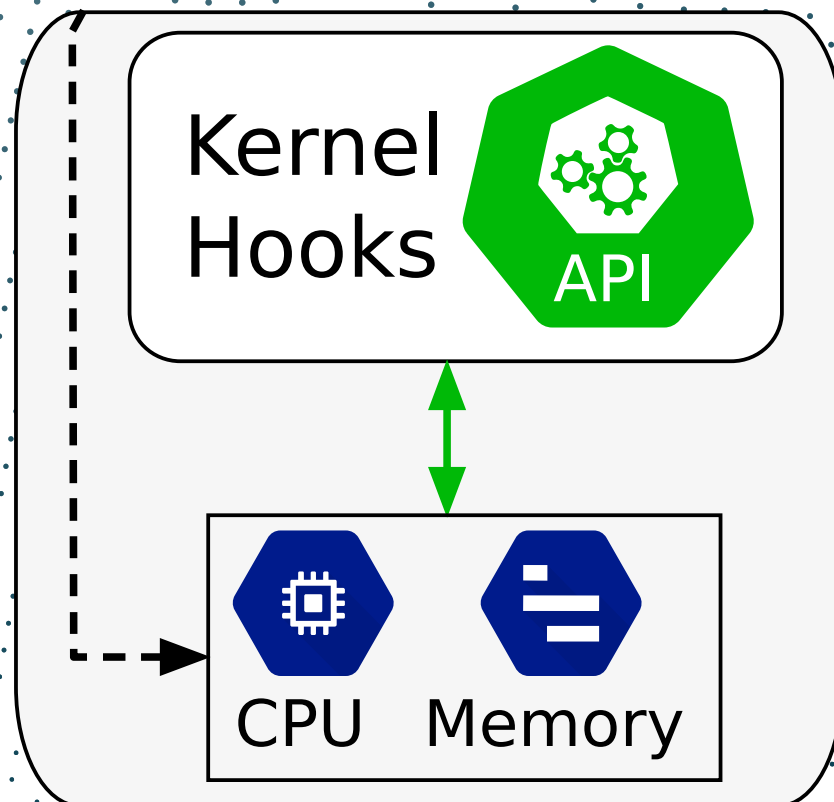
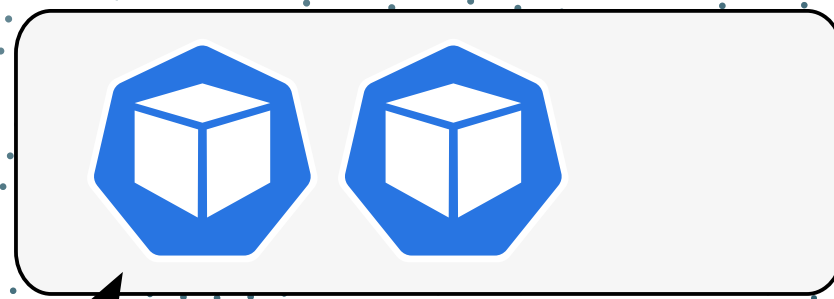
Simplified Cloud Orchestration Architecture

Application Description Files



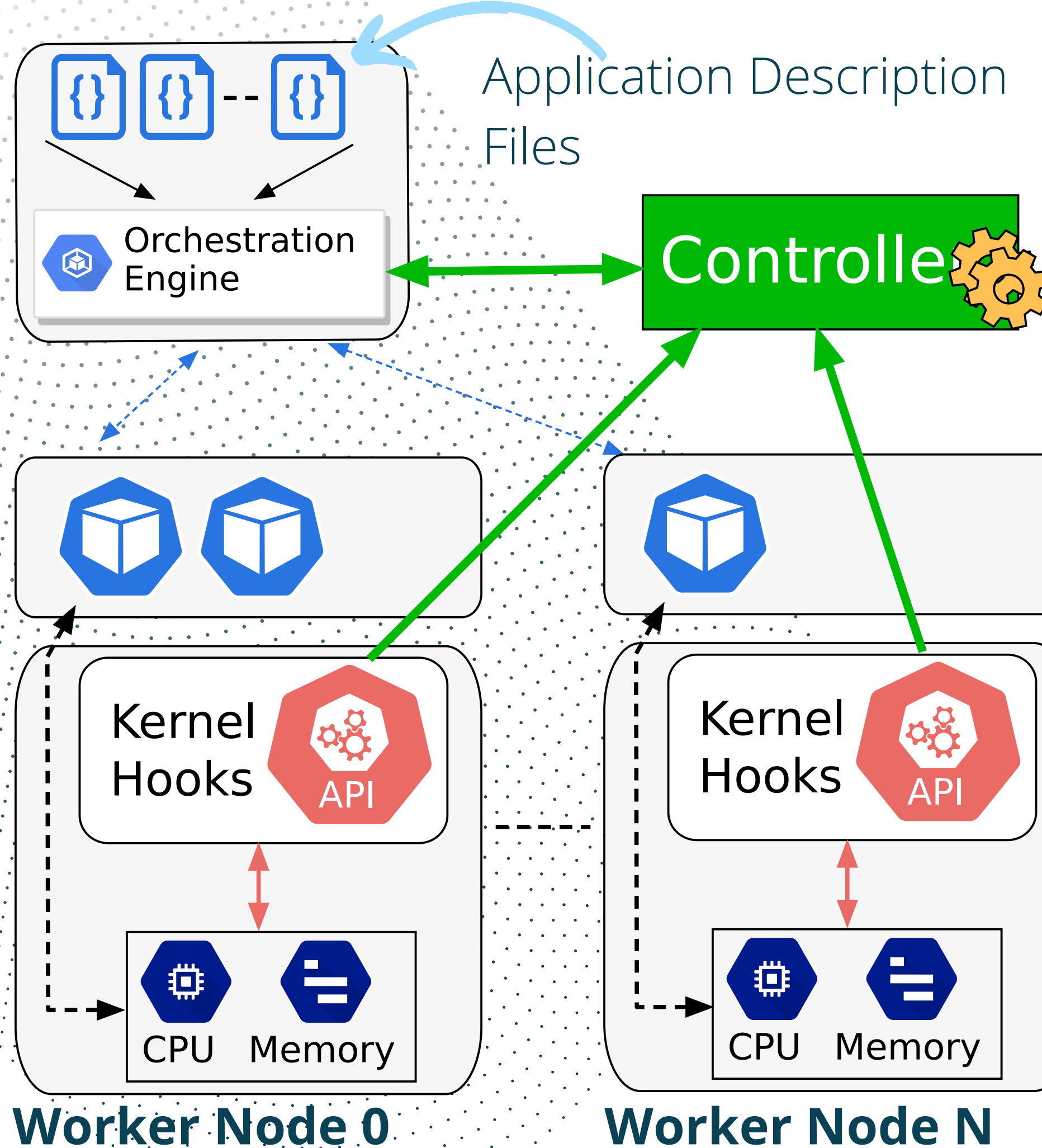
Escra Architecture

- Kernel Hooks



Worker Node 0

Worker Node N

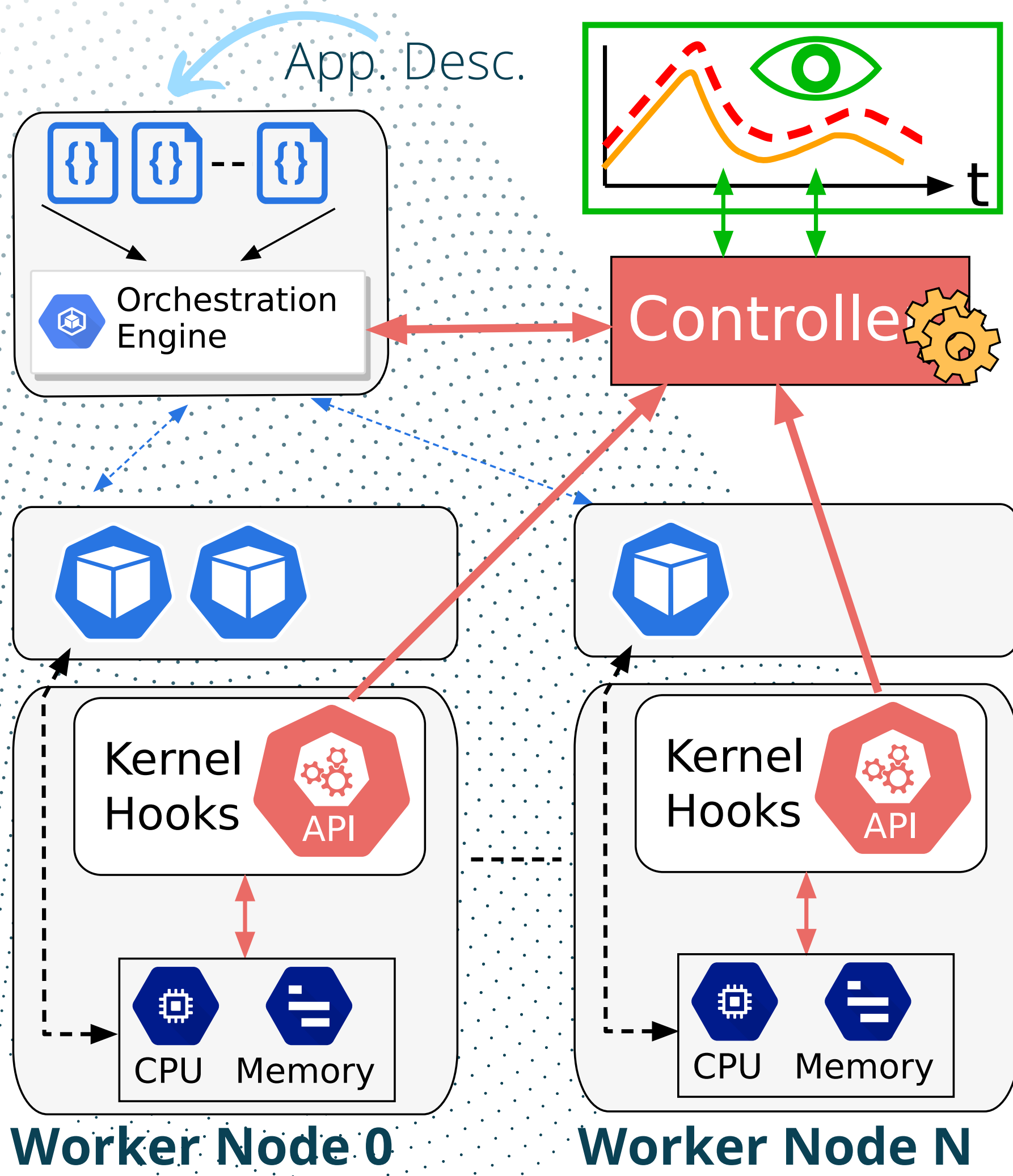


Escra Architecture

- Kernel Hooks
- **Controller**

Worker Node 0

Worker Node N



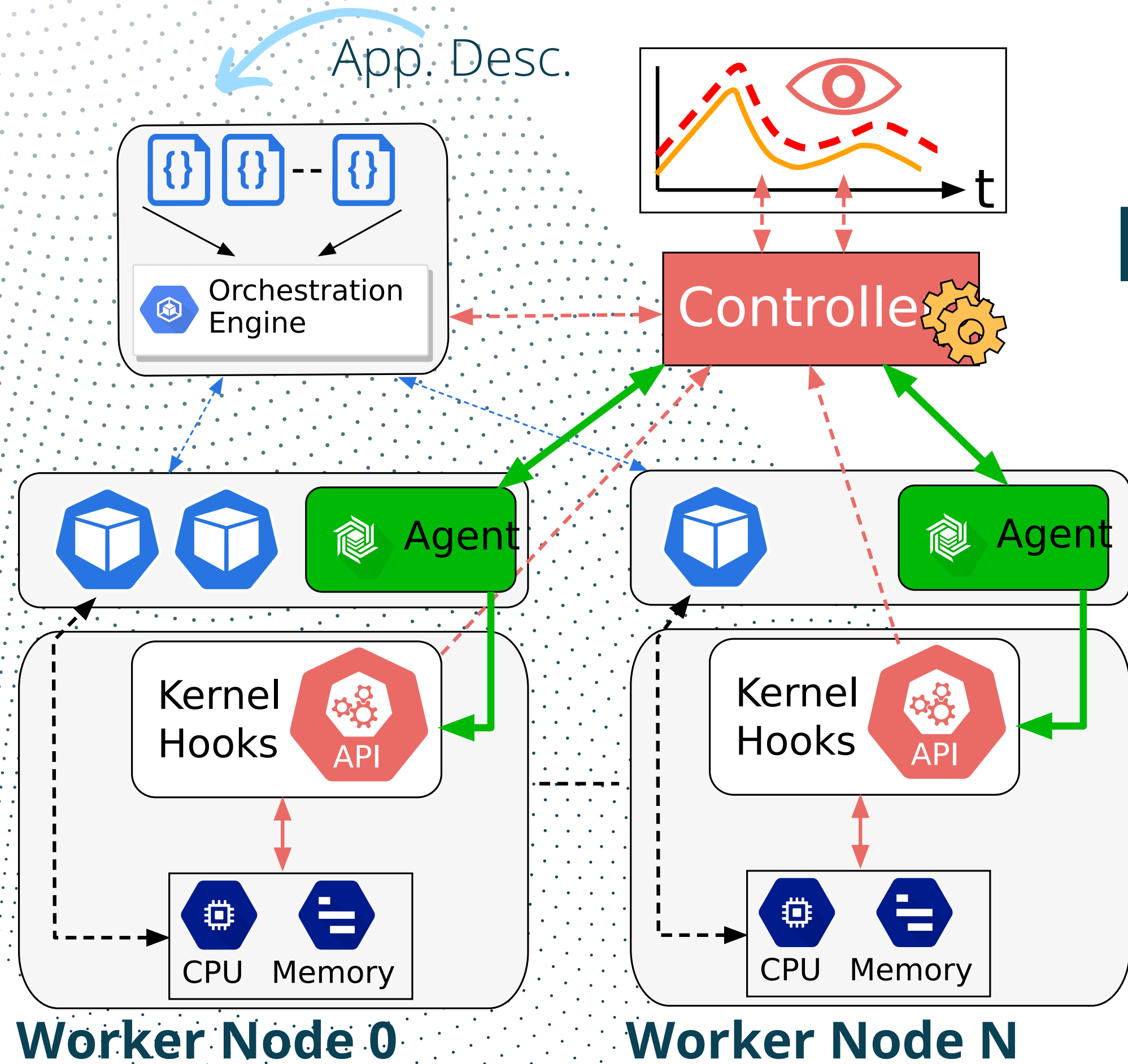
Limits
Usage

Escra Architecture

- Kernel Hooks
- Controller
- **Resource Allocator**

Escra Architecture

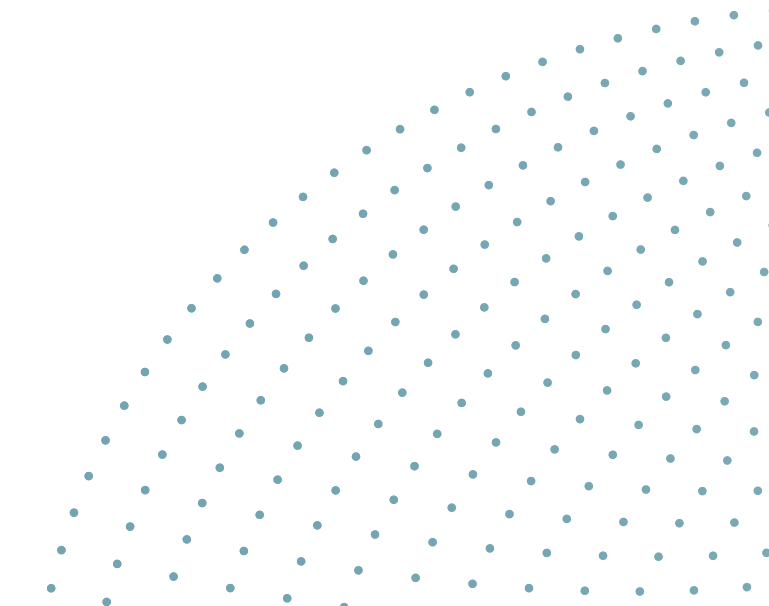
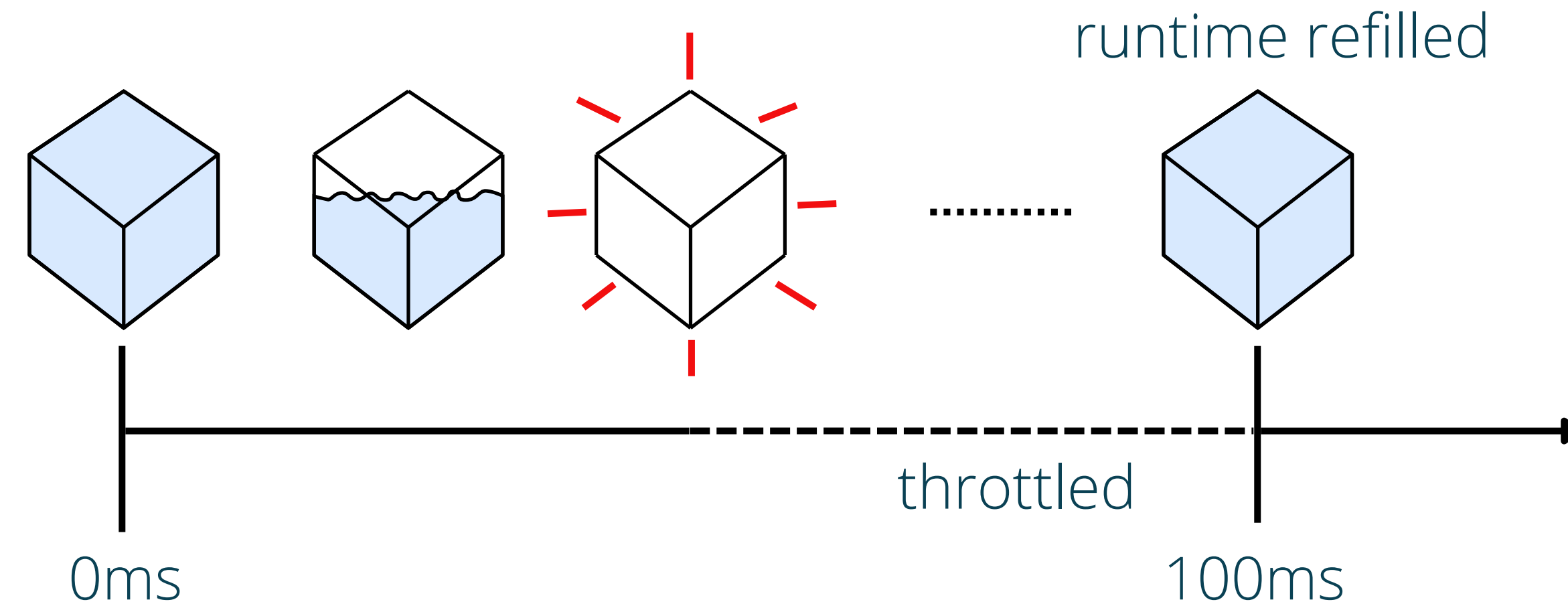
- Kernel Hooks
- Controller
- Resource Allocator
- **Agent**



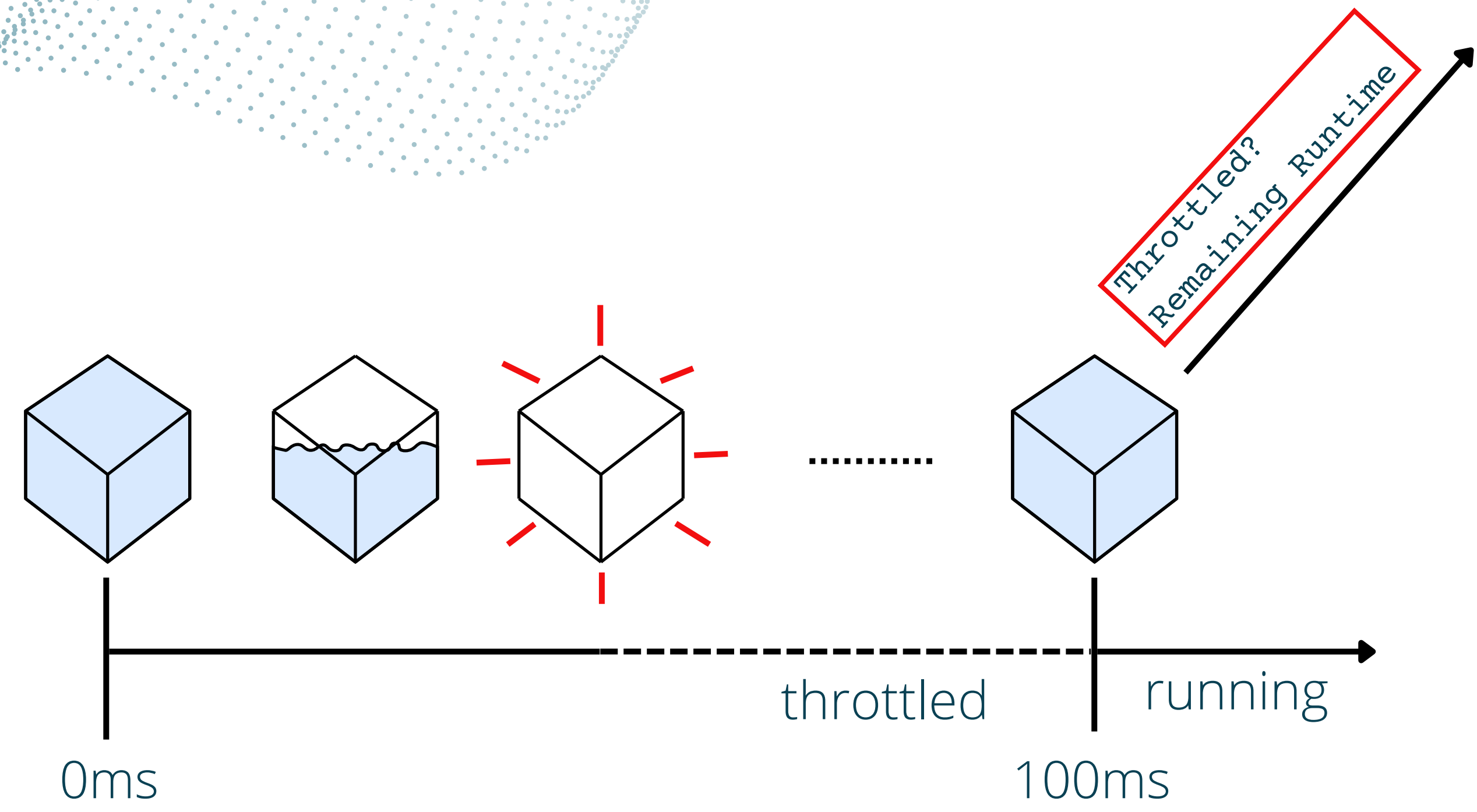
Worker Node 0

Worker Node N

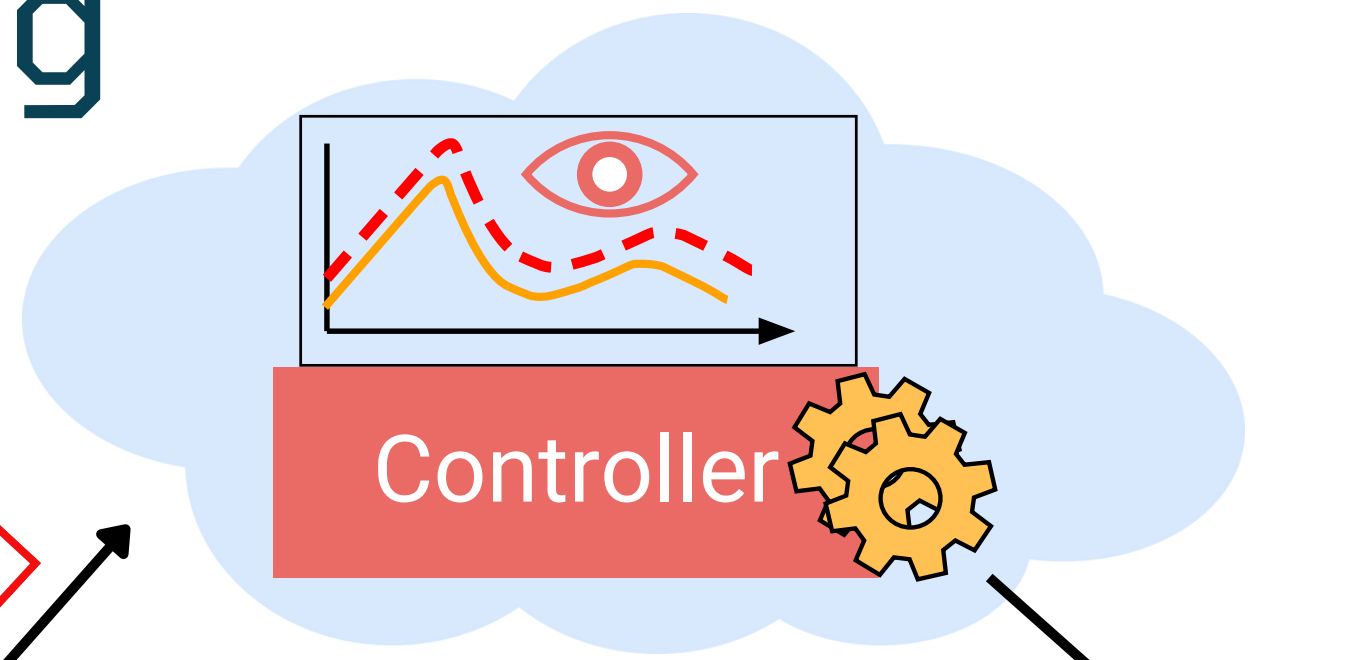
CPU Telemetry and Scaling



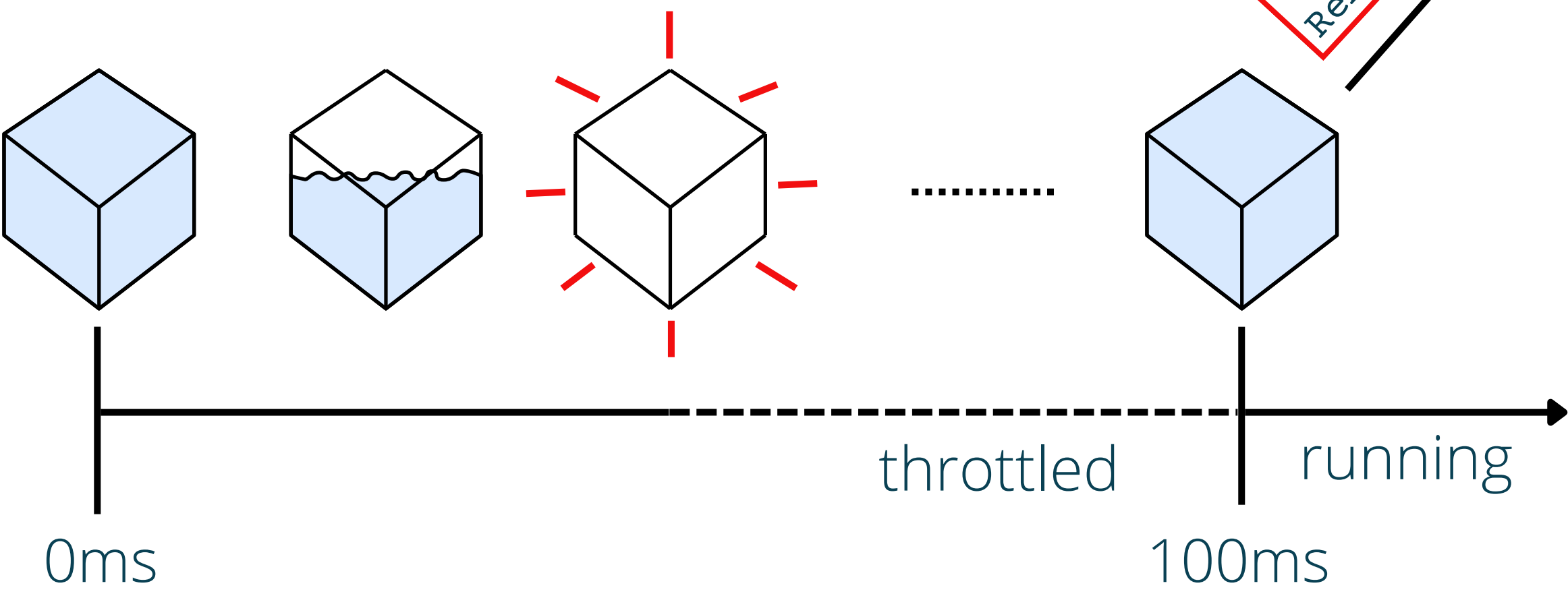
CPU Telemetry and Scaling



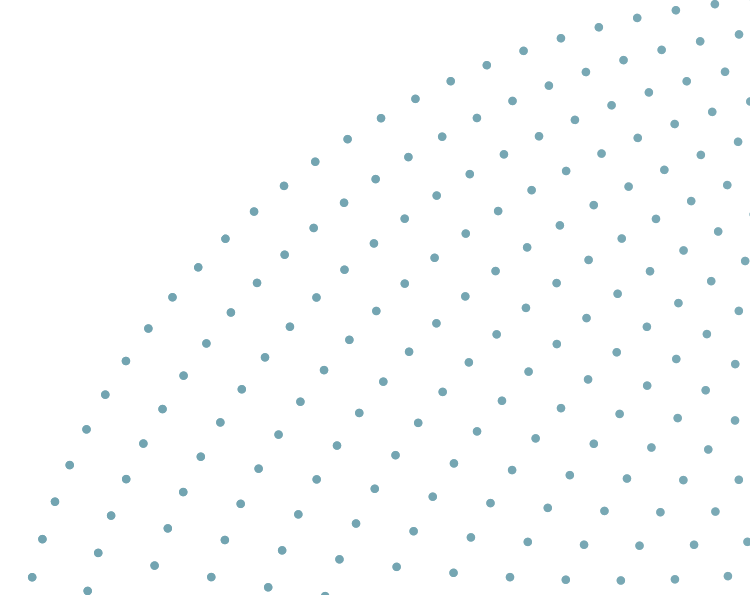
CPU Telemetry and Scaling



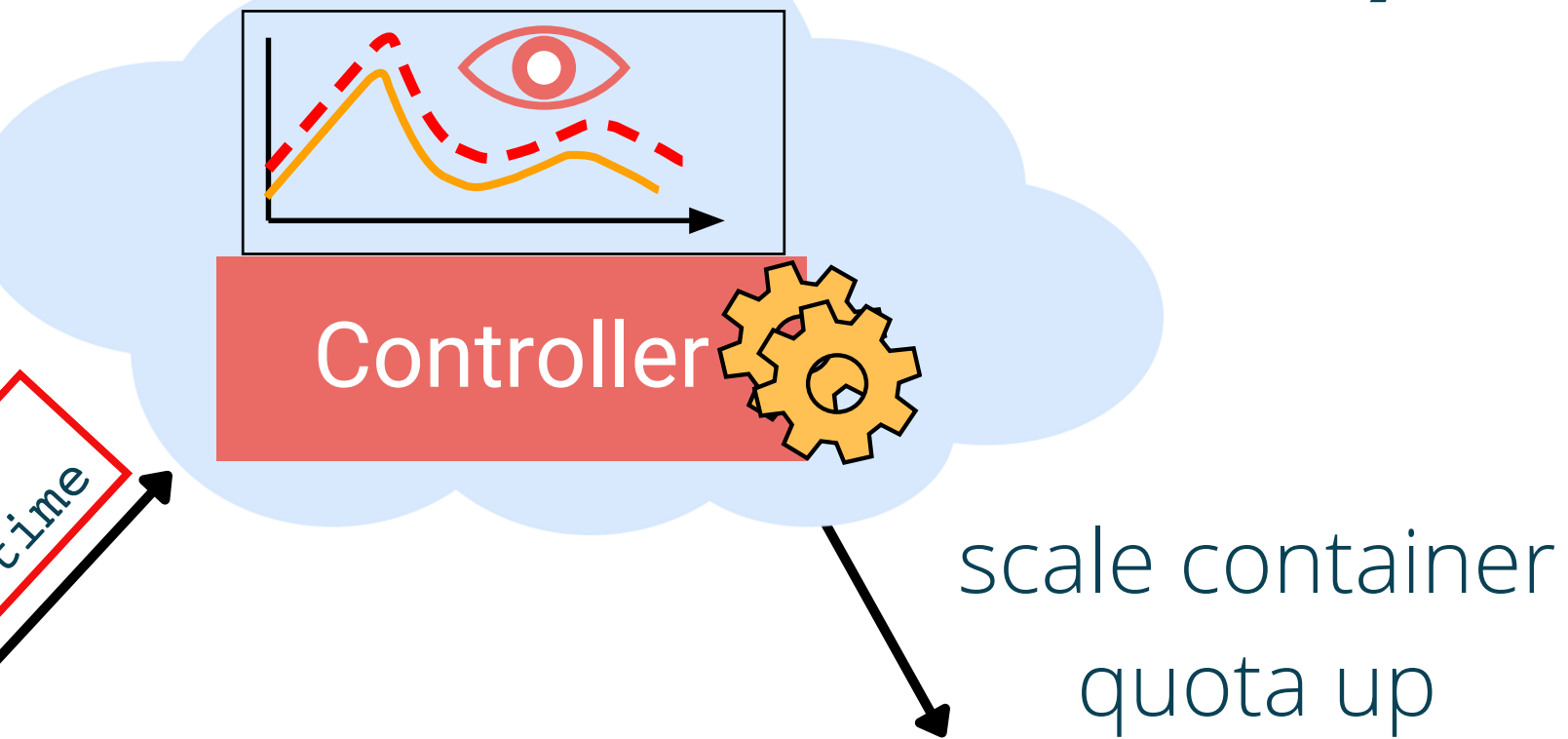
scale container
quota up



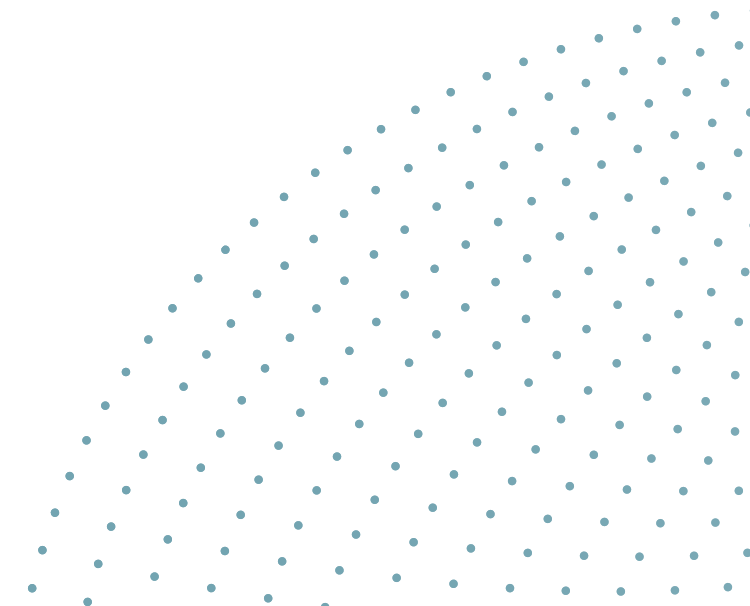
Throttled?
Remaining Runtime



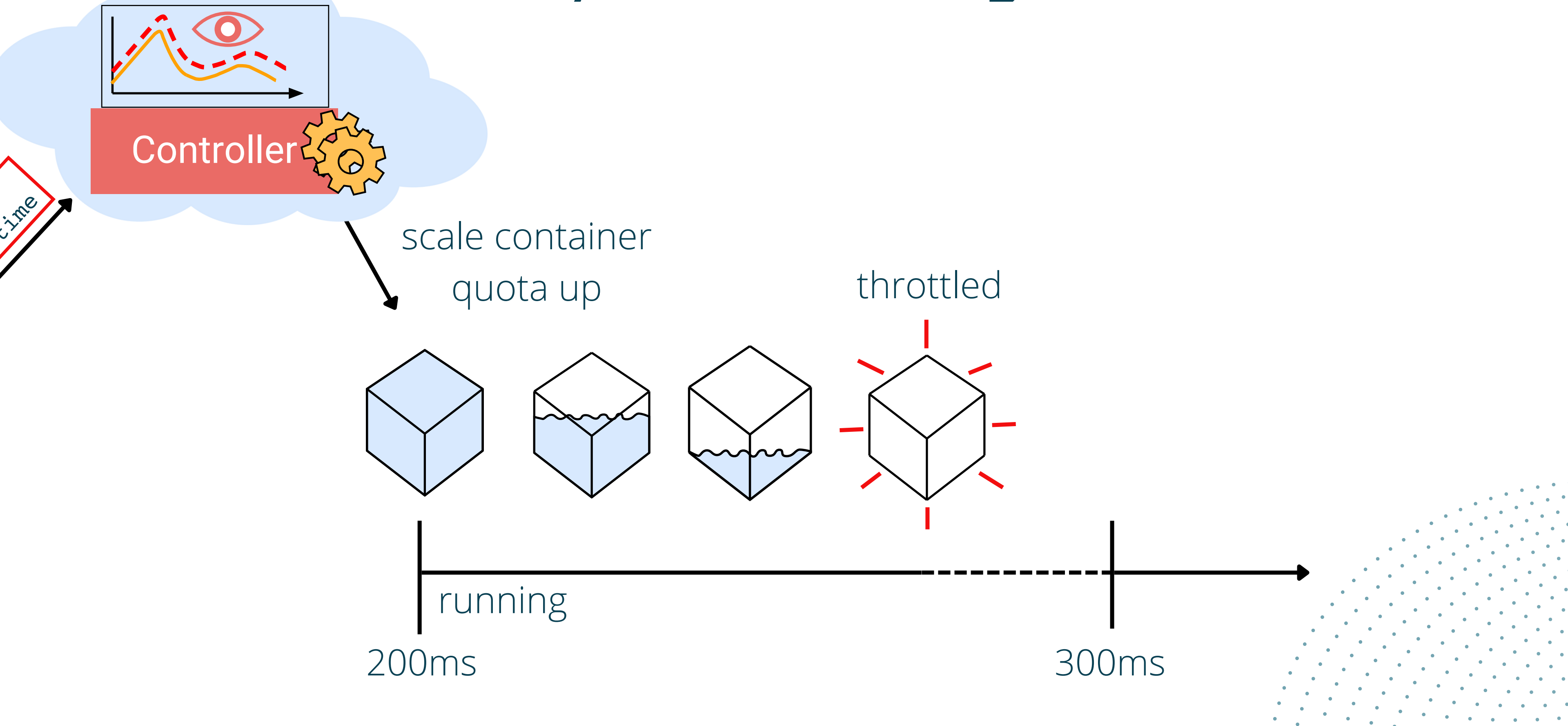
CPU Telemetry and Scaling



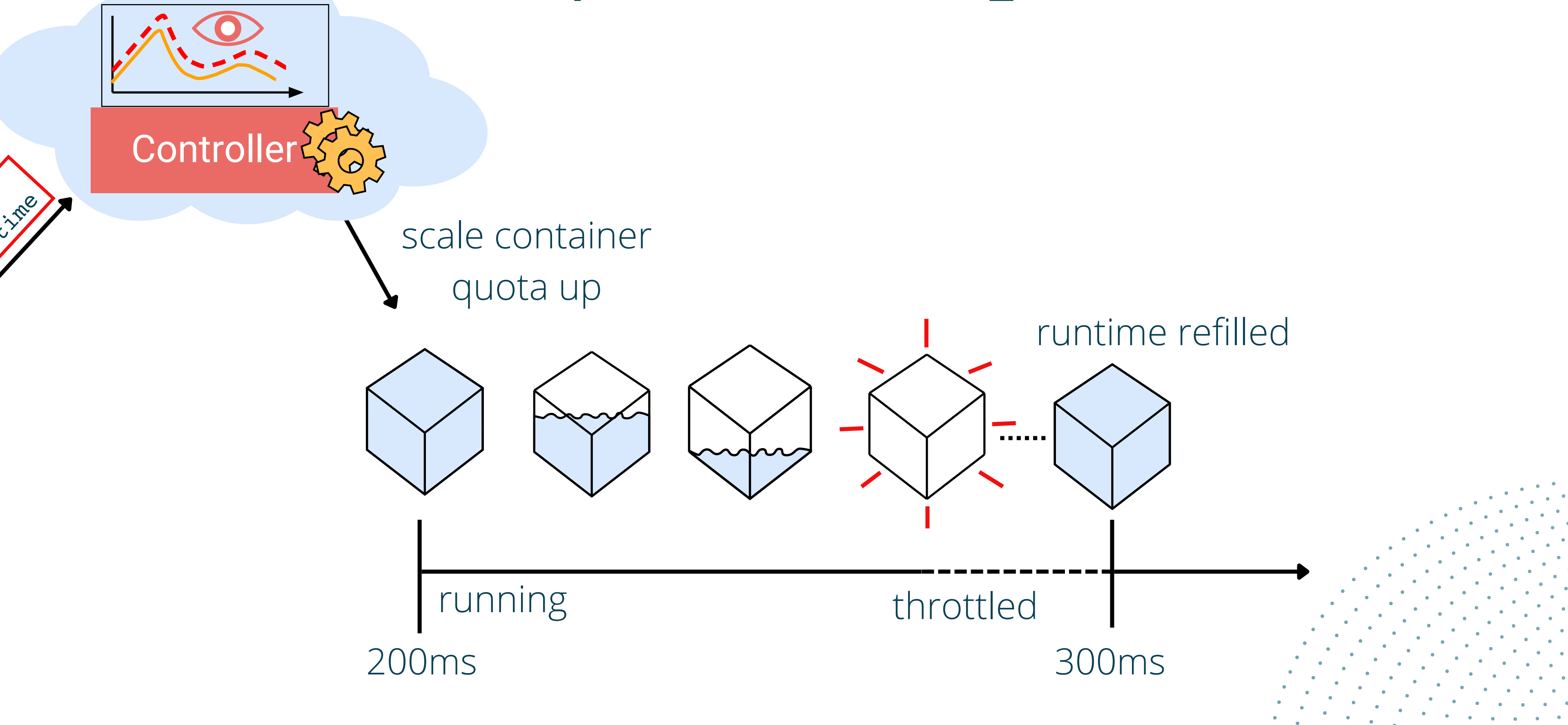
$$C(i)_q[t+1] = C(i)_q[t] + \frac{\sum_{t=0}^n C(i)_{th}[t]}{n} * \Upsilon(\Omega_l - \sum_{i=0}^{\lambda} C(i)_q[t])$$



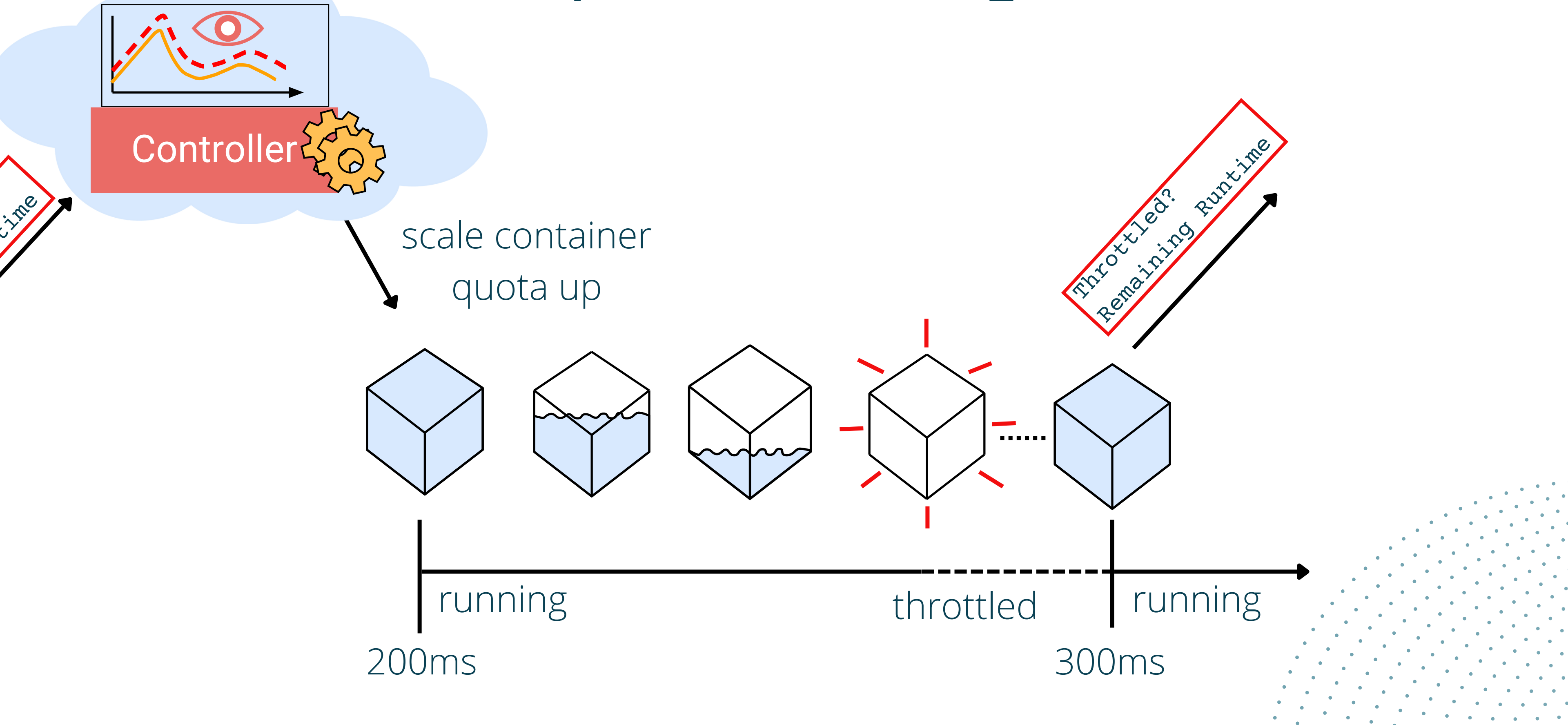
CPU Telemetry and Scaling



CPU Telemetry and Scaling

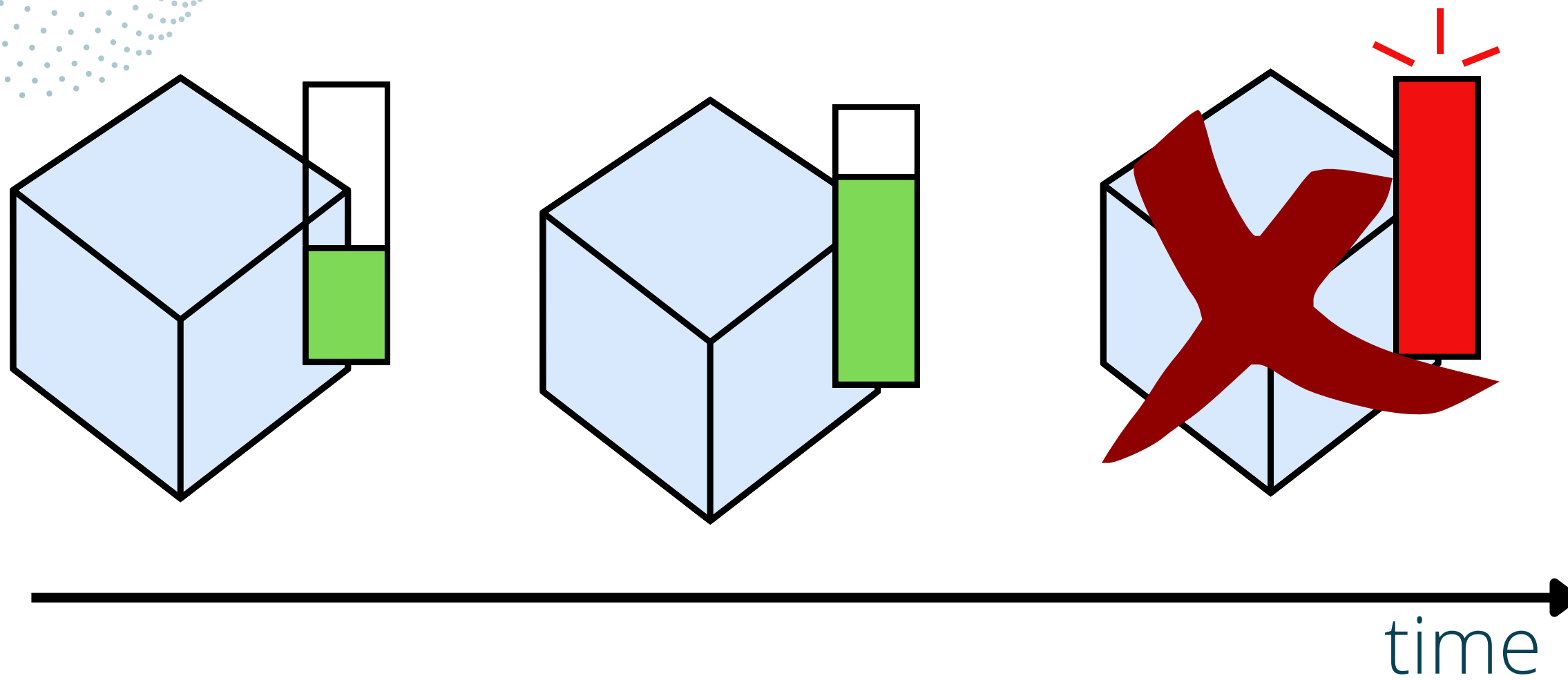


CPU Telemetry and Scaling



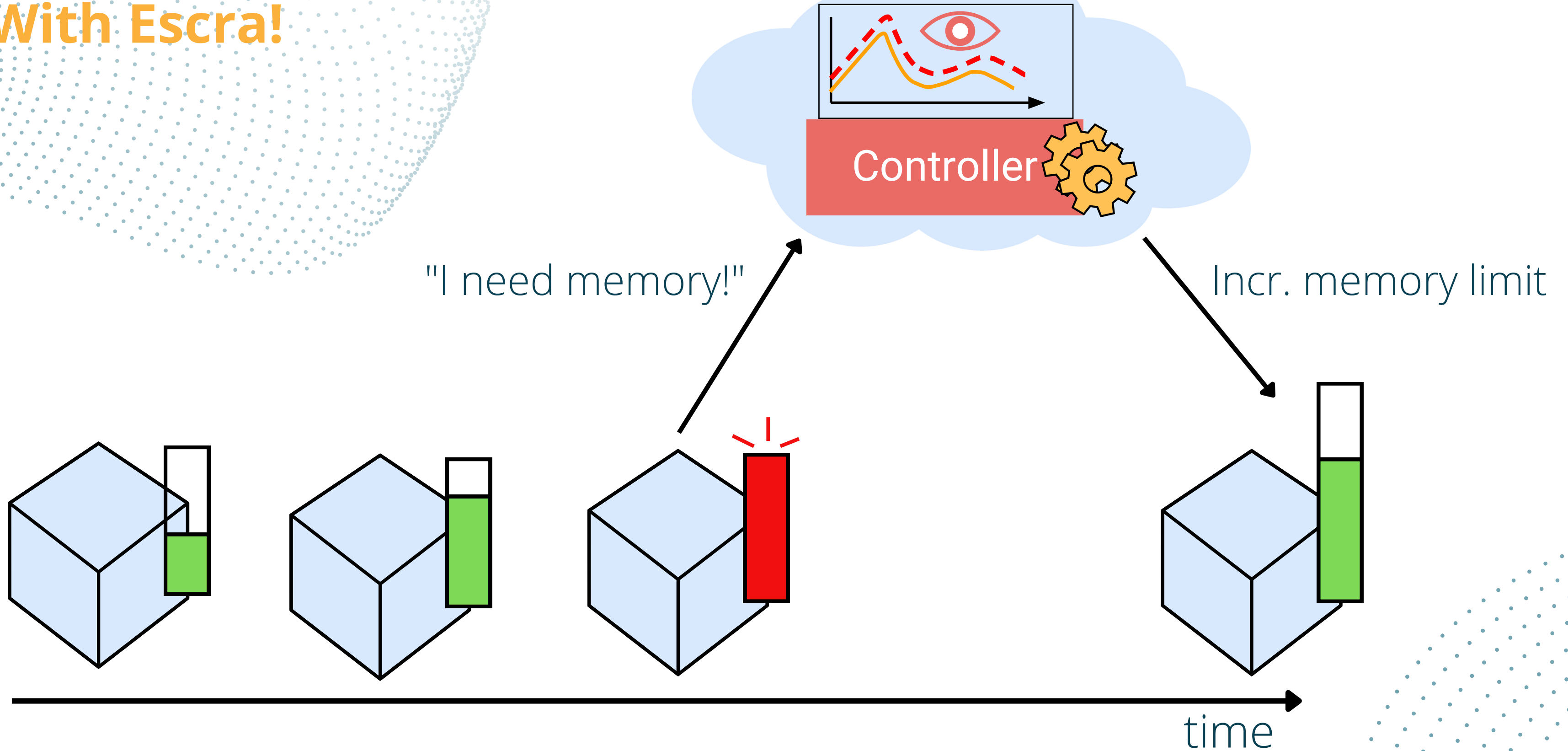
Event-Based Memory Scaling

Typical Scenario



Event-Based Memory Scaling

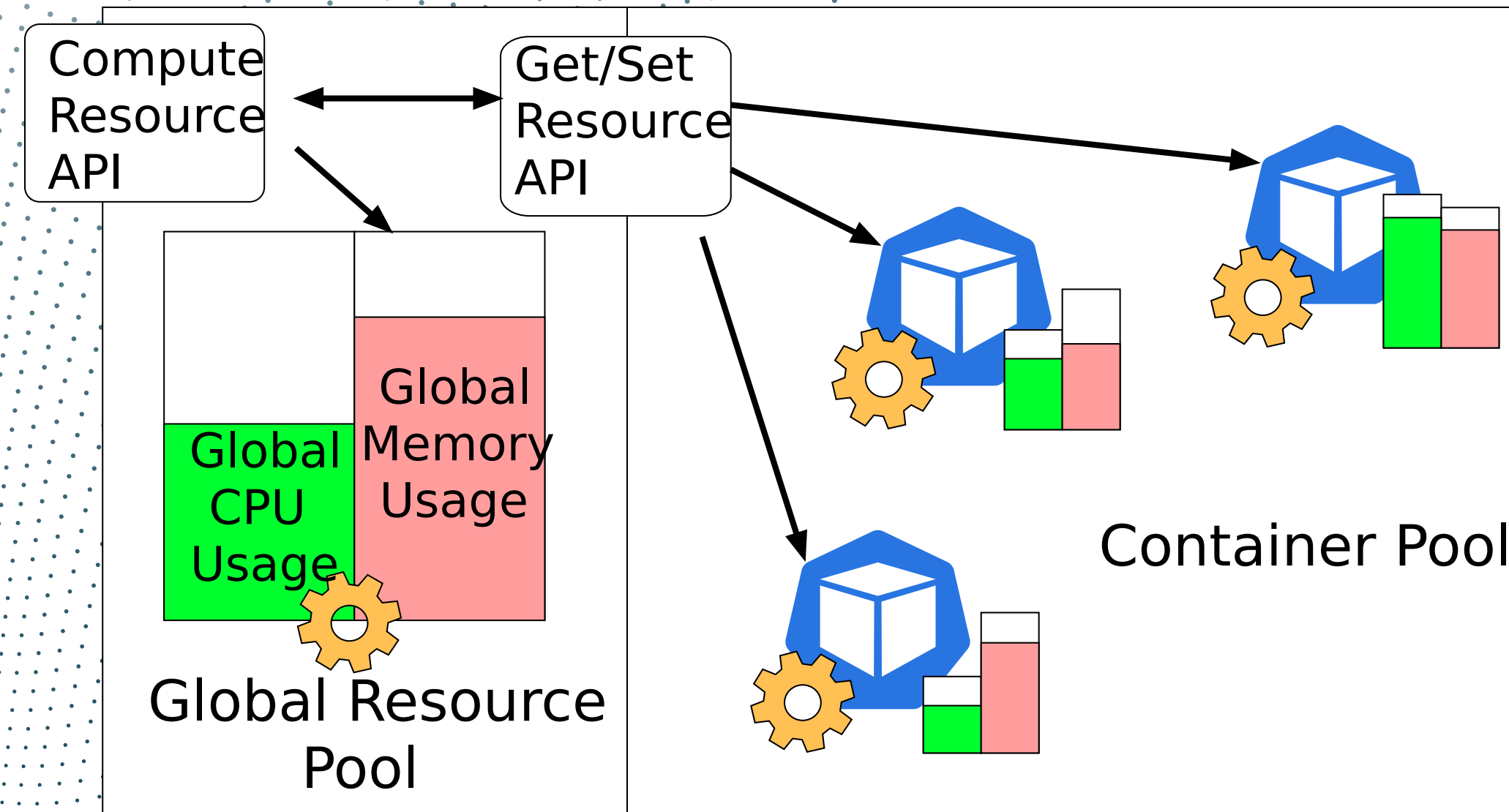
With Escra!



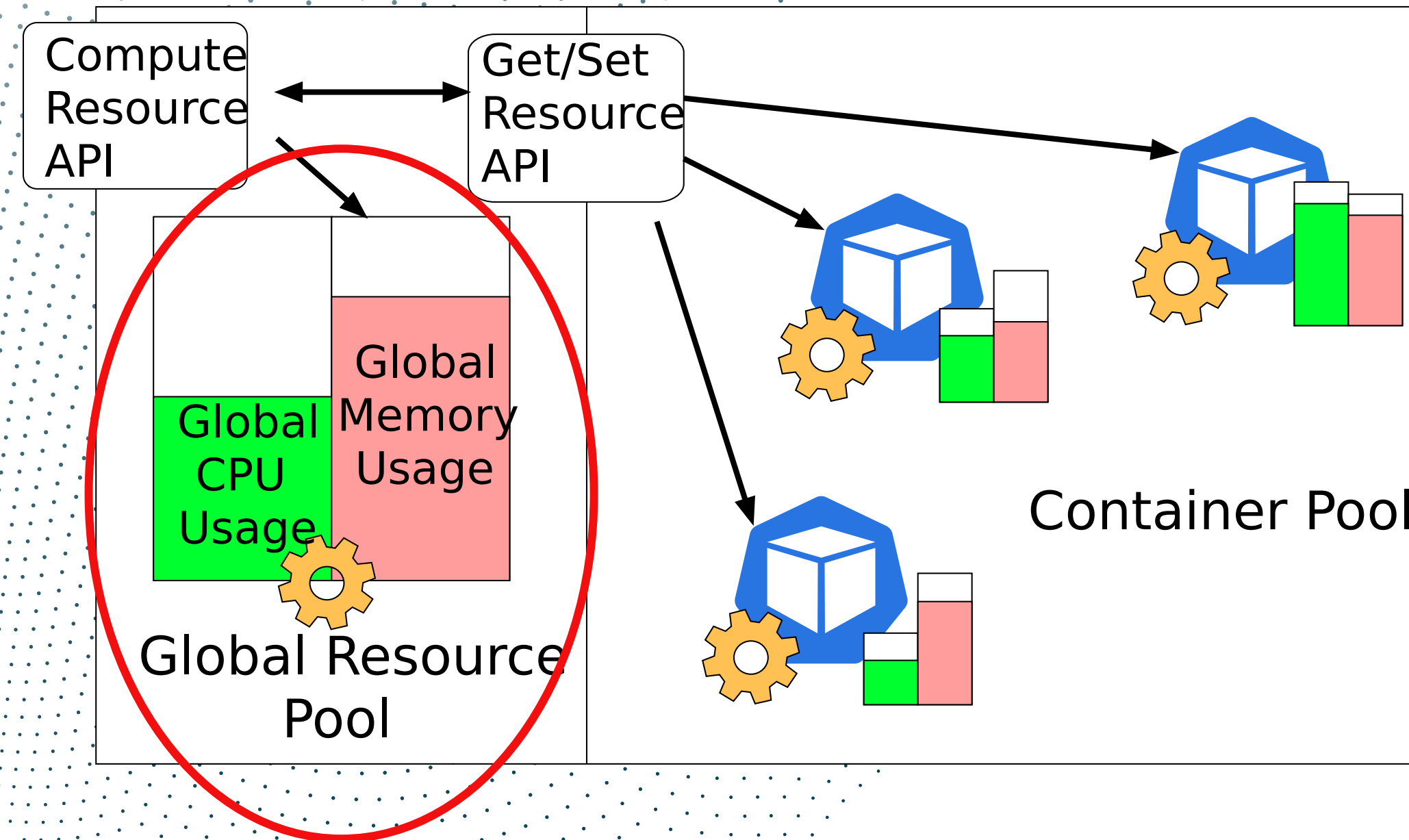


Distributed Container

Distributed Container

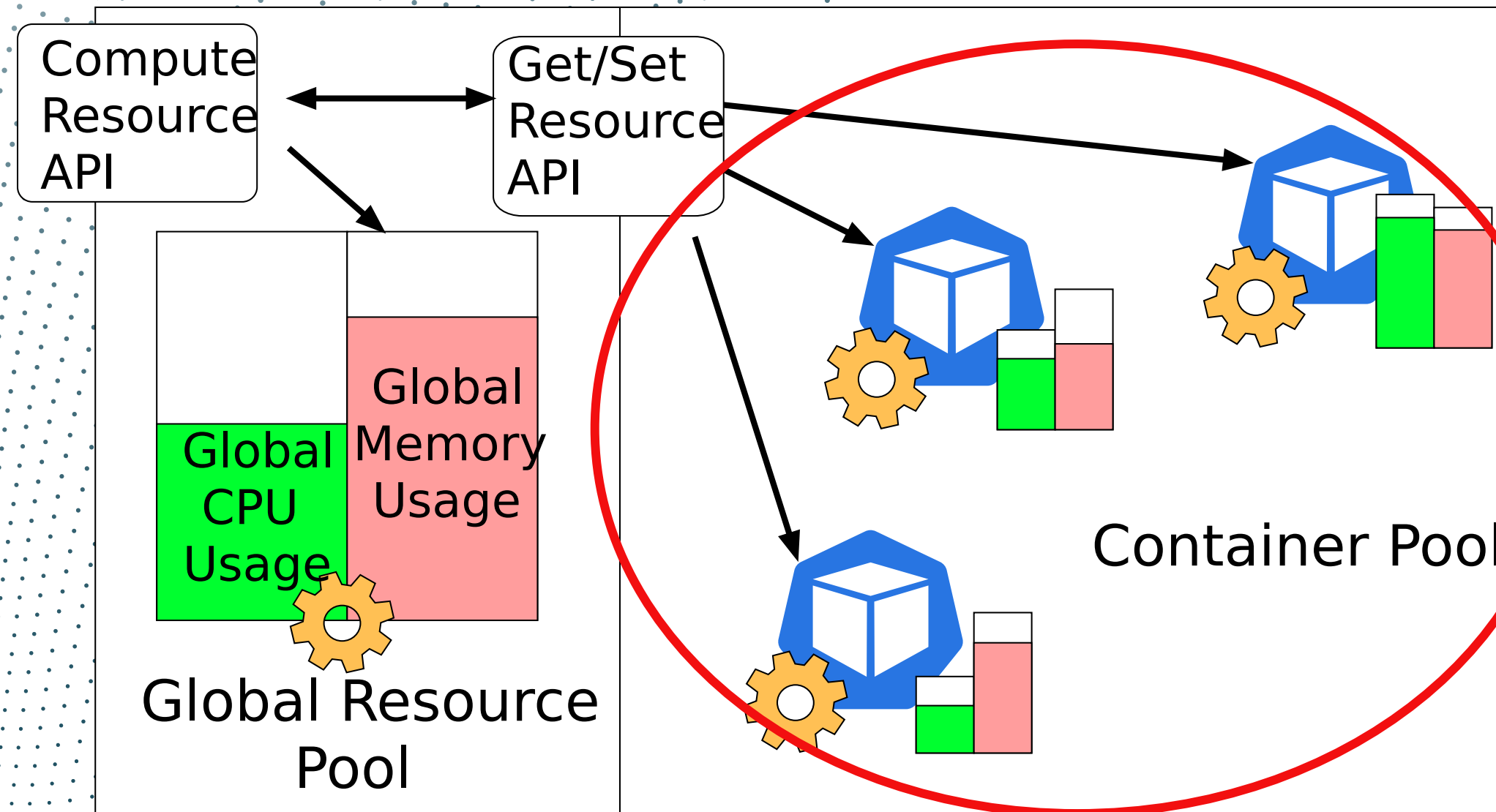


Distributed Container



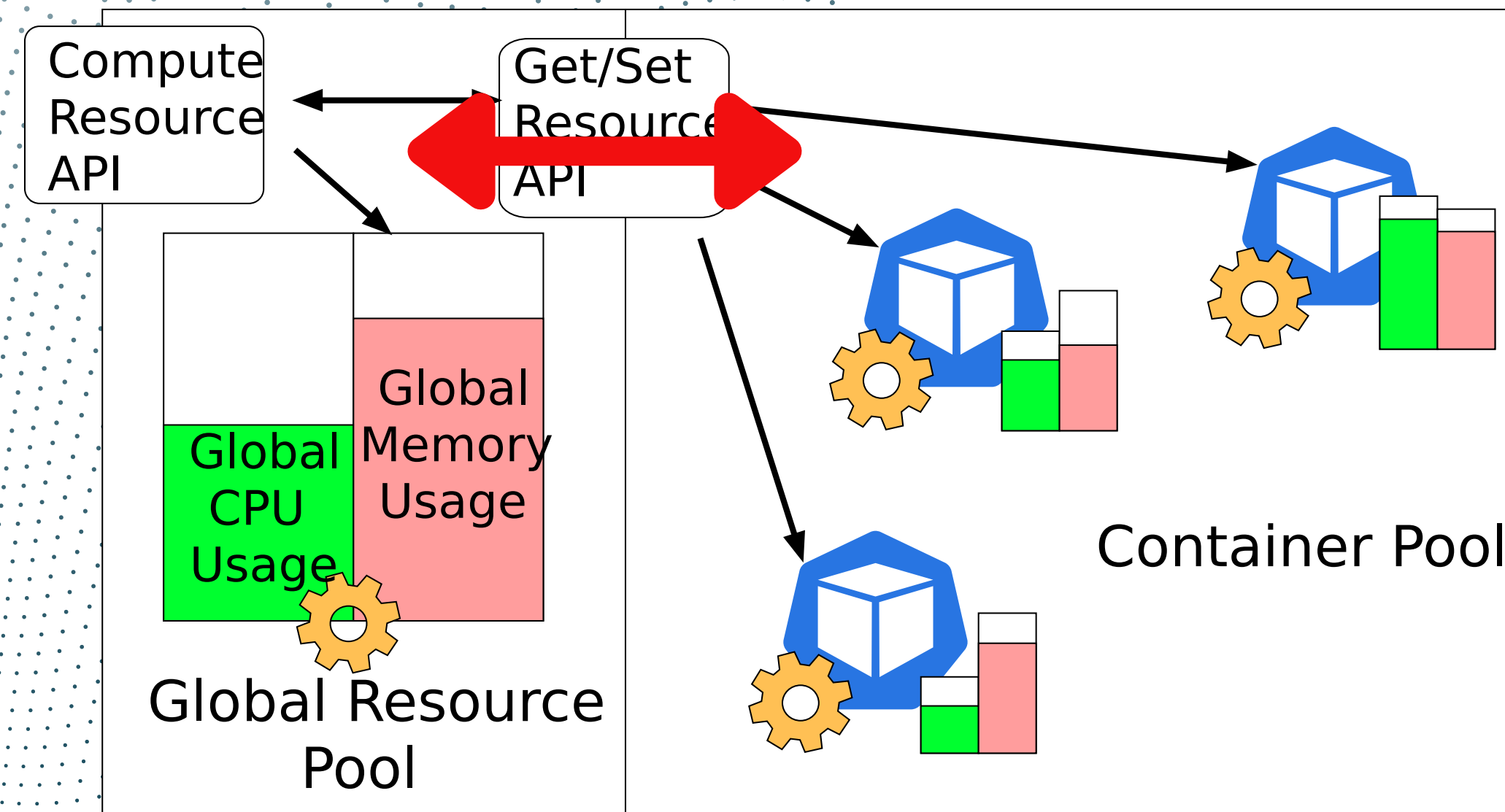
- Enforces per-application resources limits

Distributed Container



- Enforces per-application resources limits
- Per-container resource usage and limit tracking

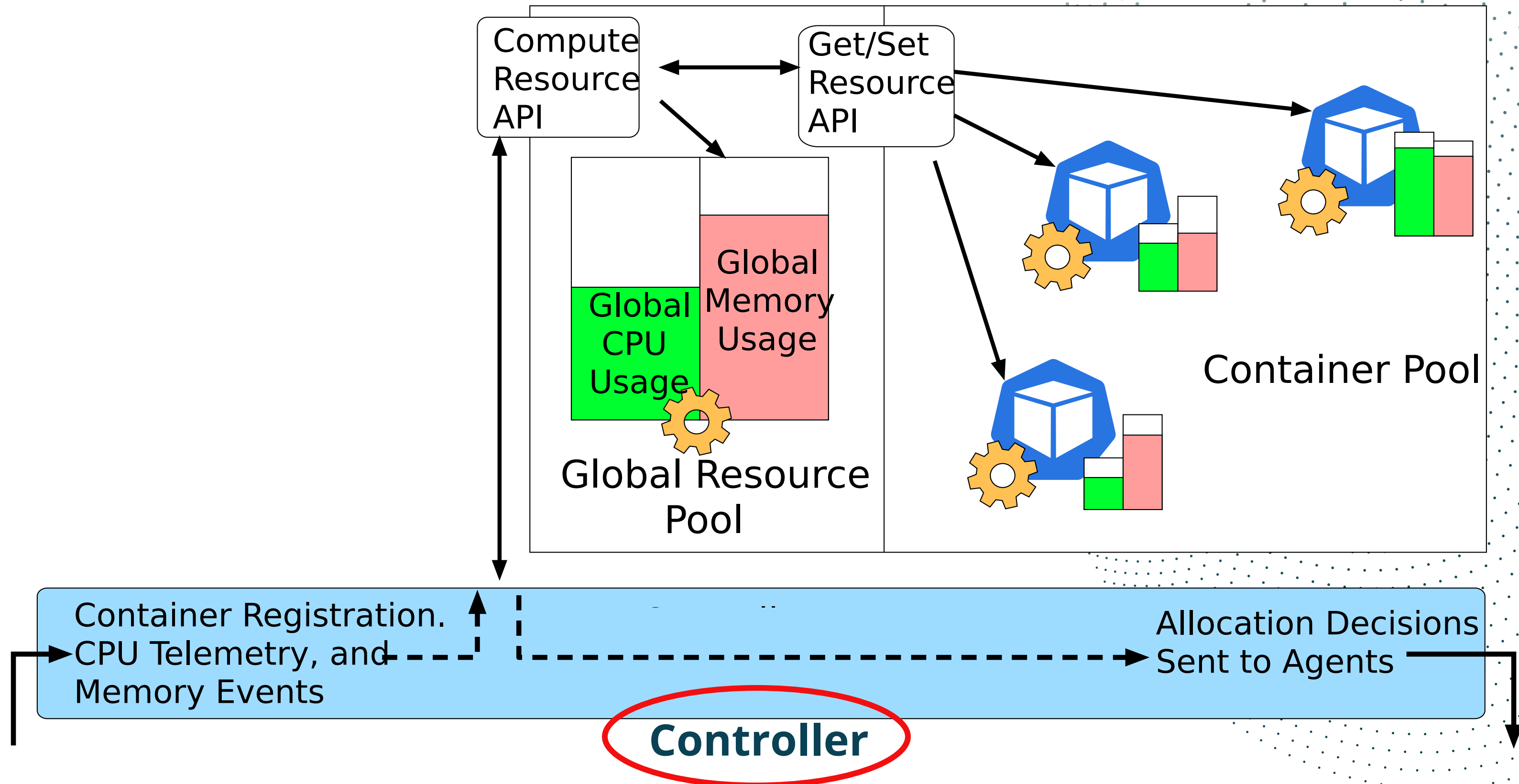
Distributed Container



- Enforces per-application resources limits
- Per-container resource usage and limit tracking
- Containers dynamically share compute resources at runtime

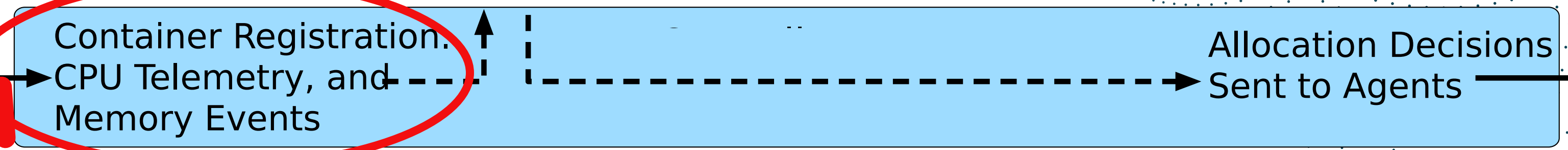
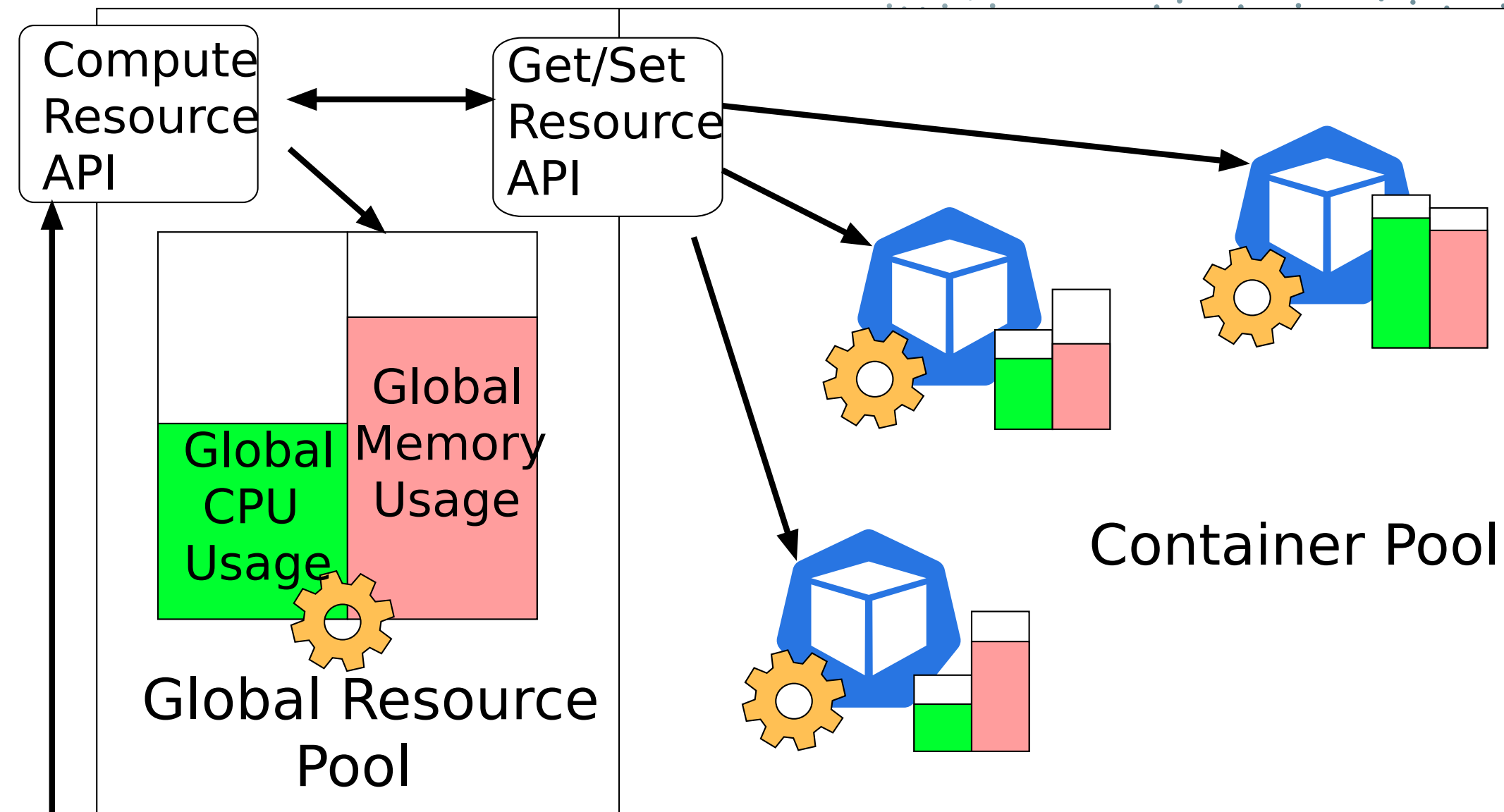
Escra Controller

Distributed Container



Escra Controller

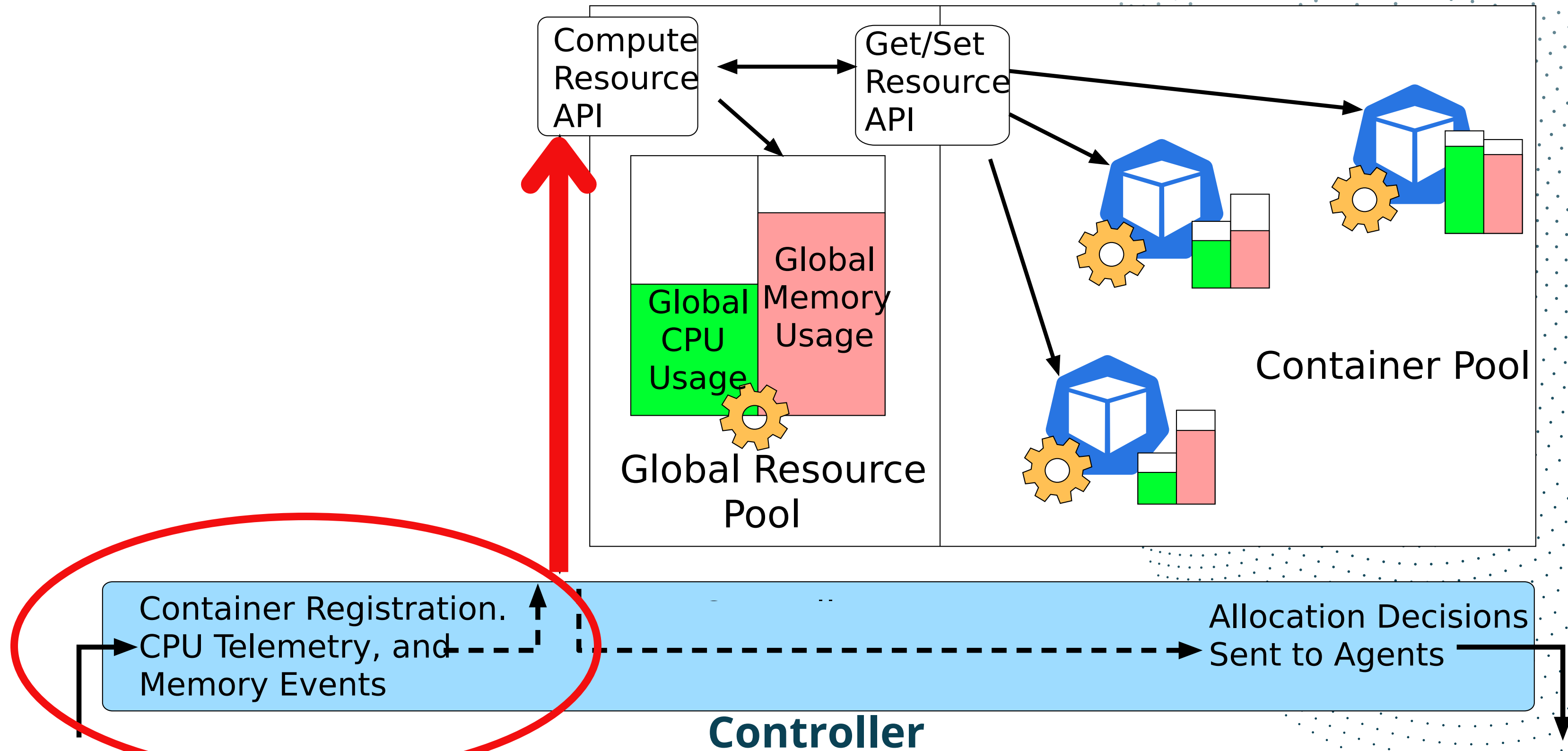
Distributed Container



Controller

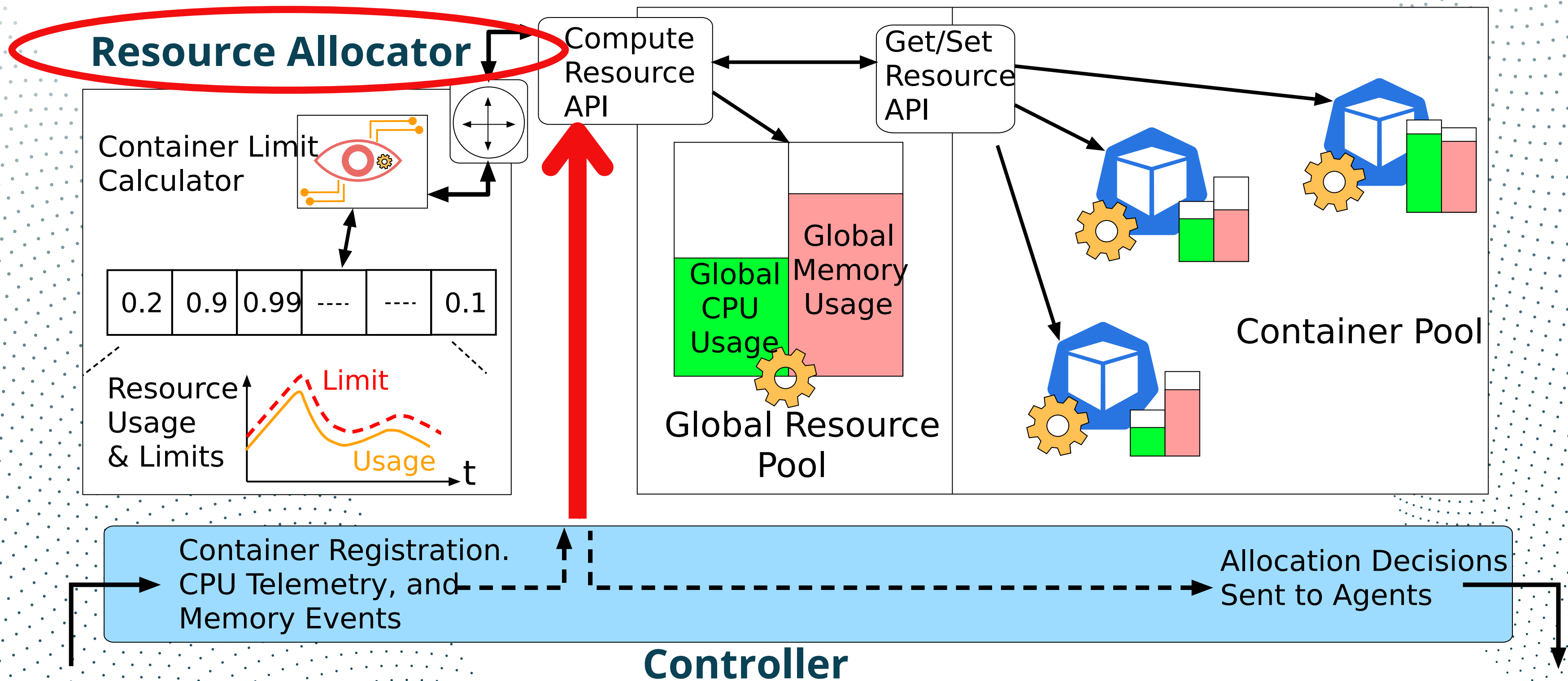
Escra Controller

Distributed Container



Escra - Control Node

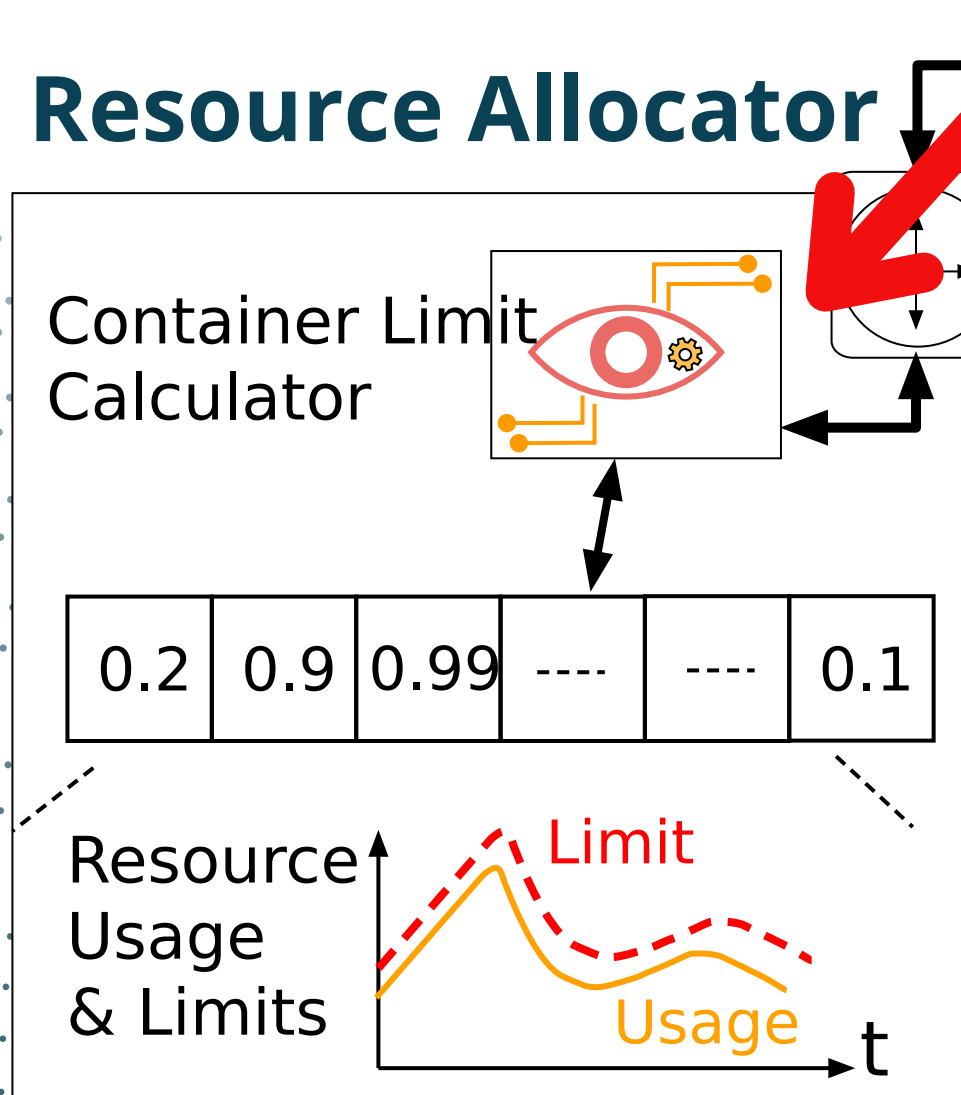
Distributed Container



Escra - Control Node

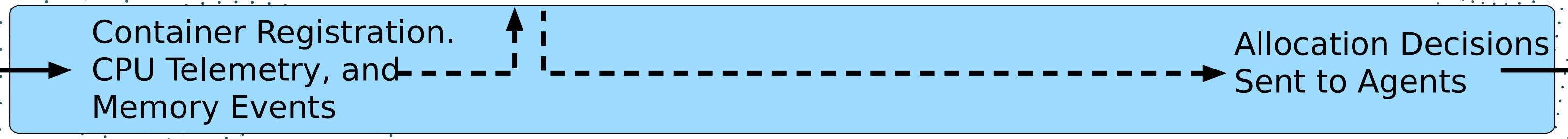
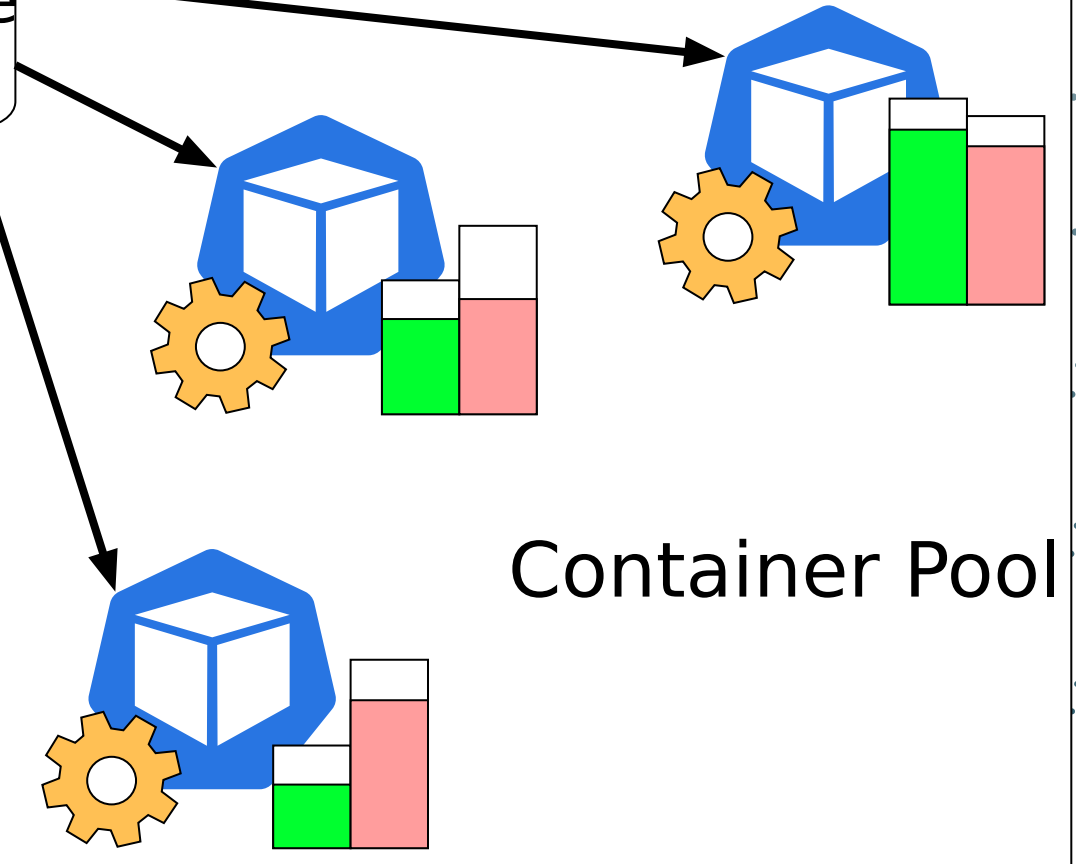
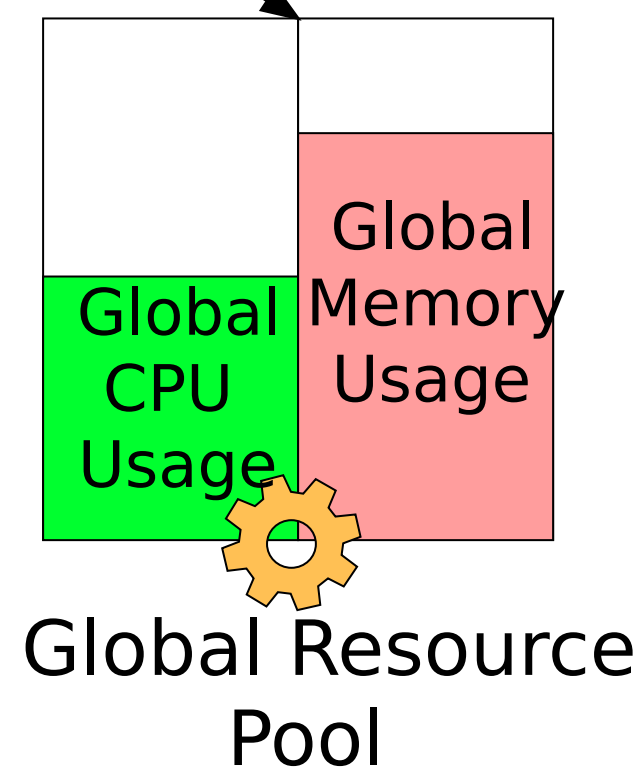
Distributed Container

Resource Allocator



Compute Resource API

Get/Set Resource API

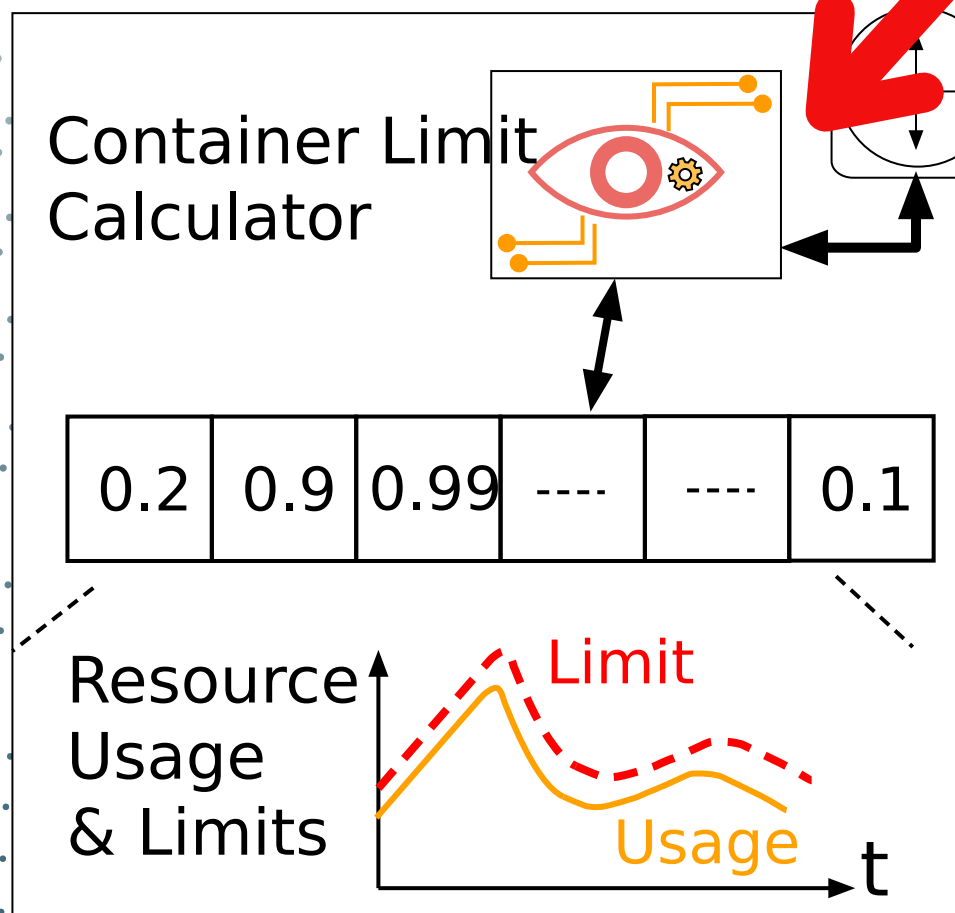


Controller

Escra - Control Node

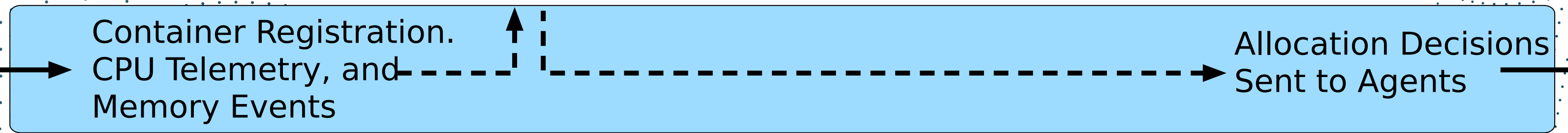
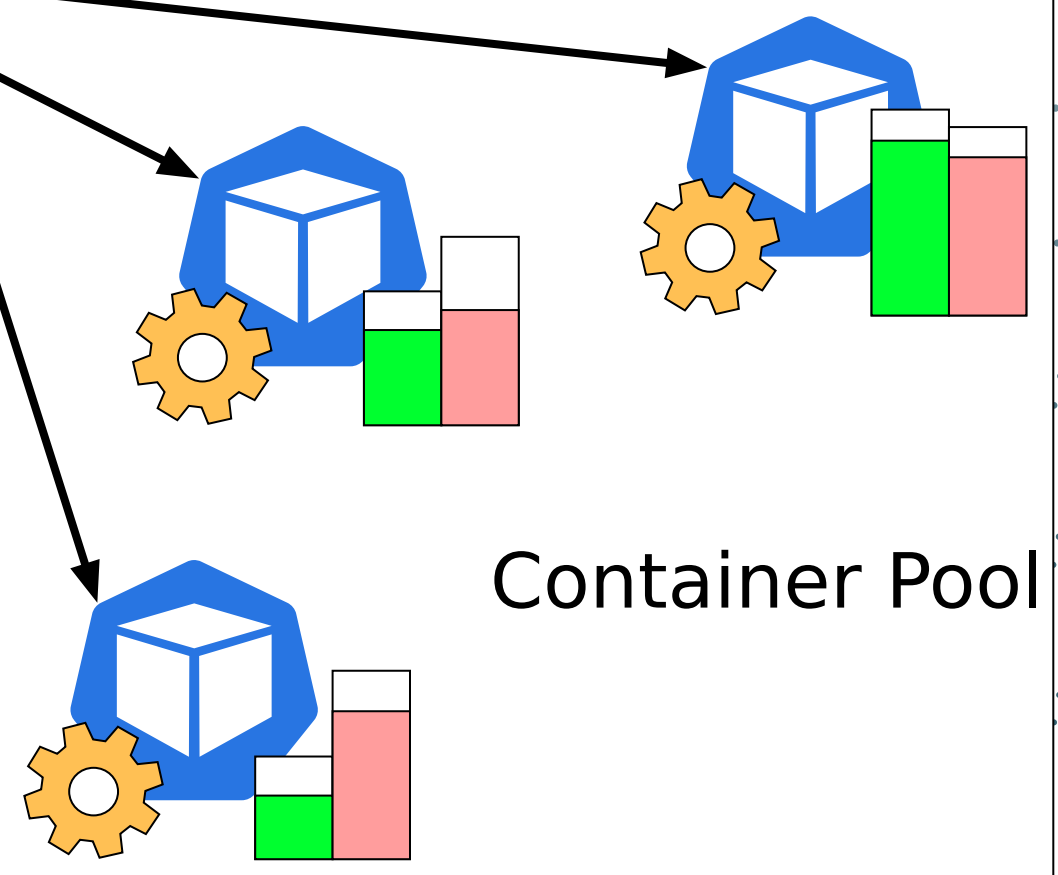
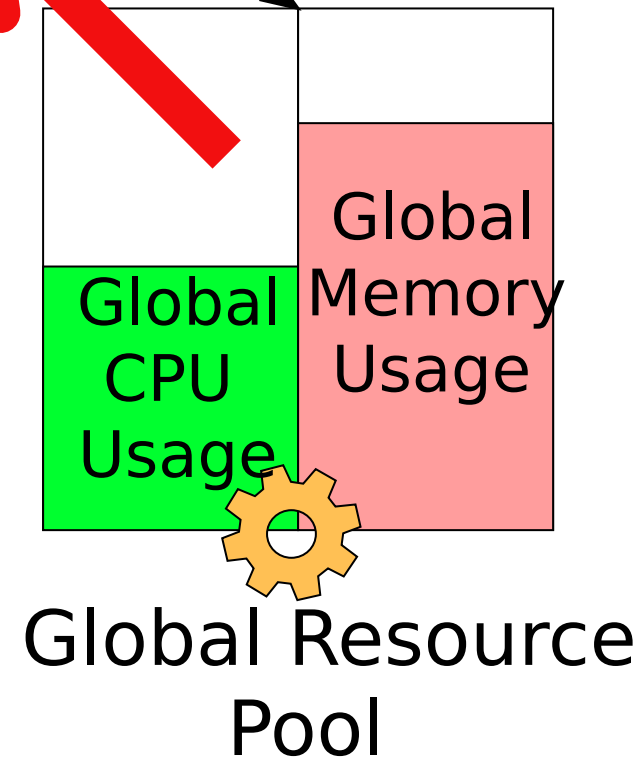
Distributed Container

Resource Allocator



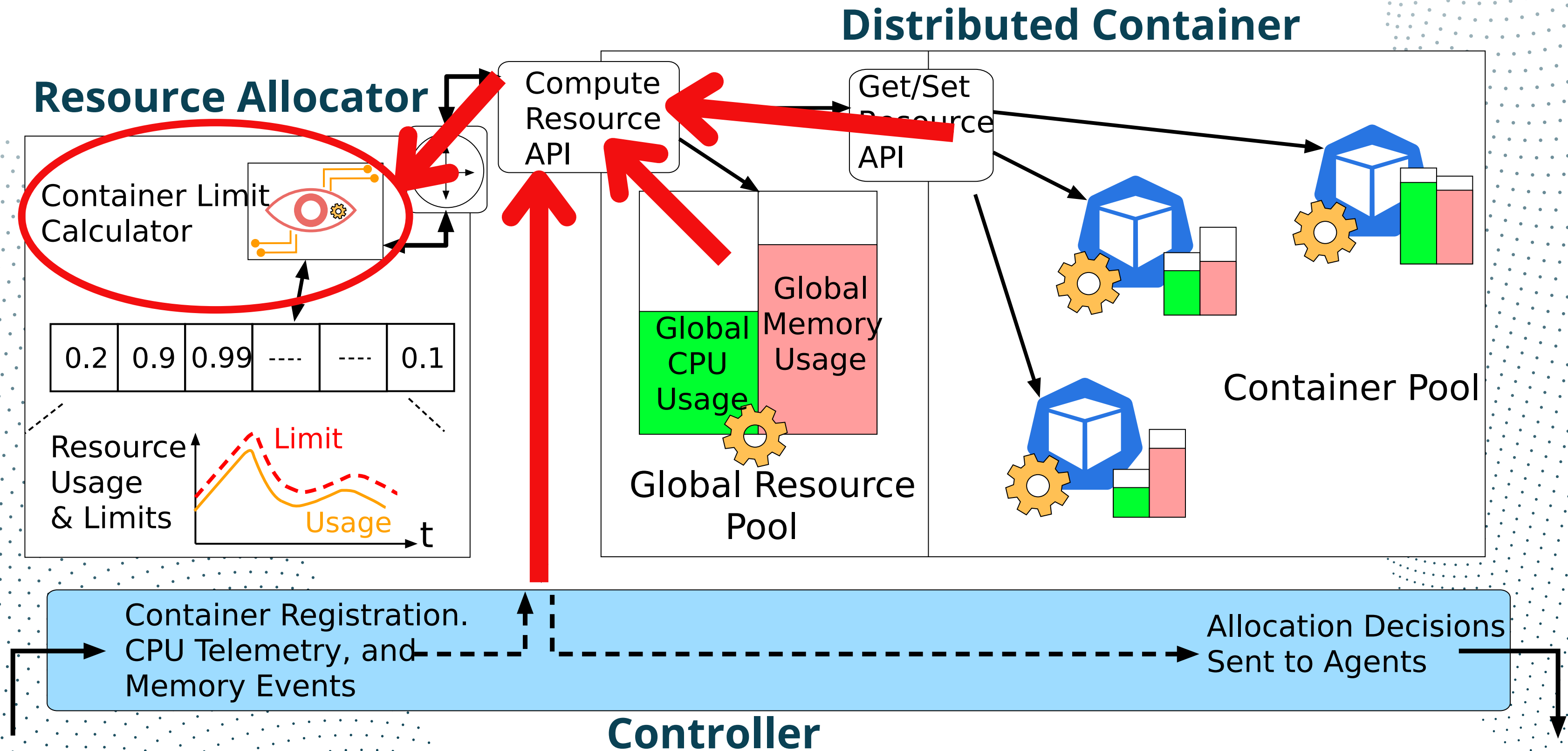
Compute Resource API

Get/Set Resource API



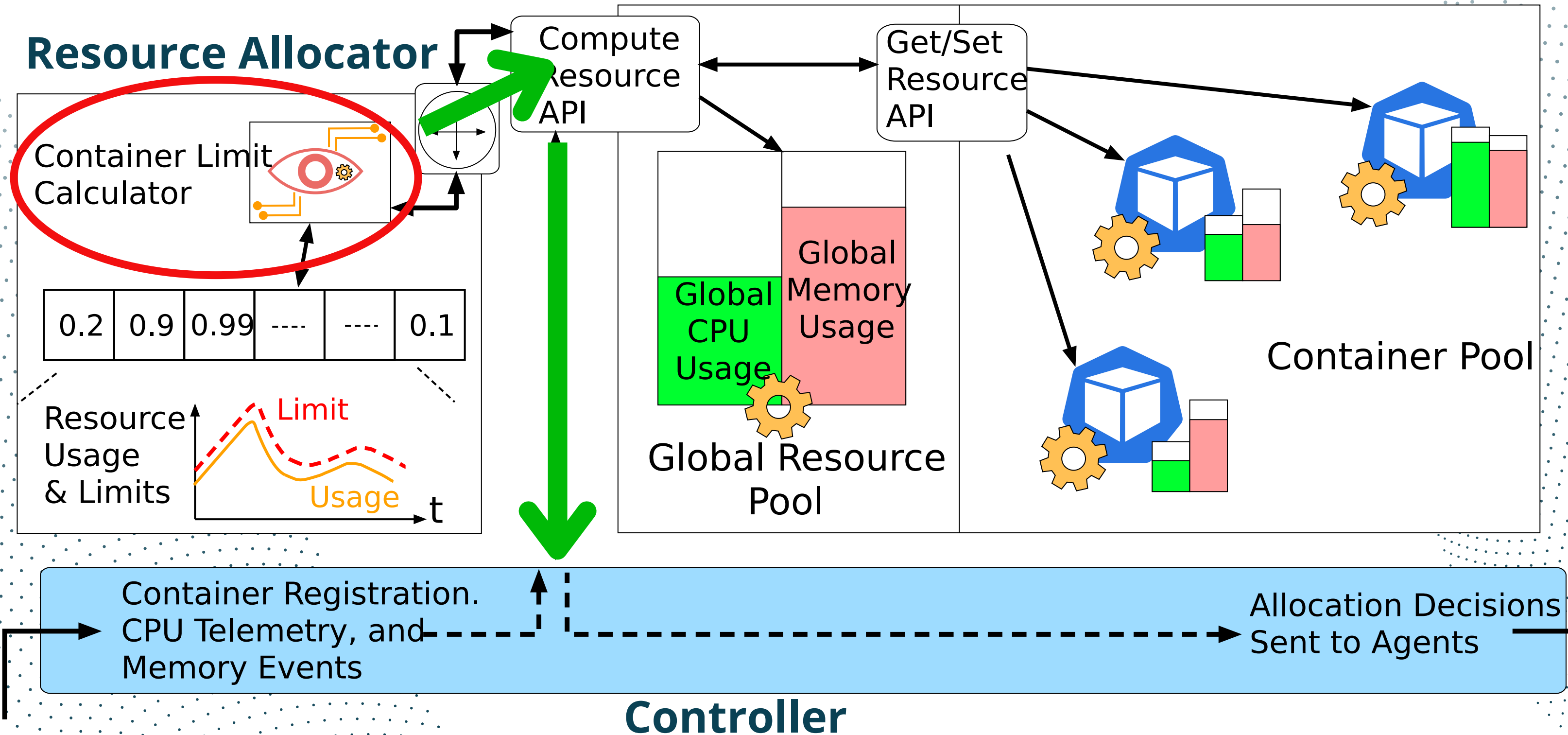
Controller

Escra - Control Node

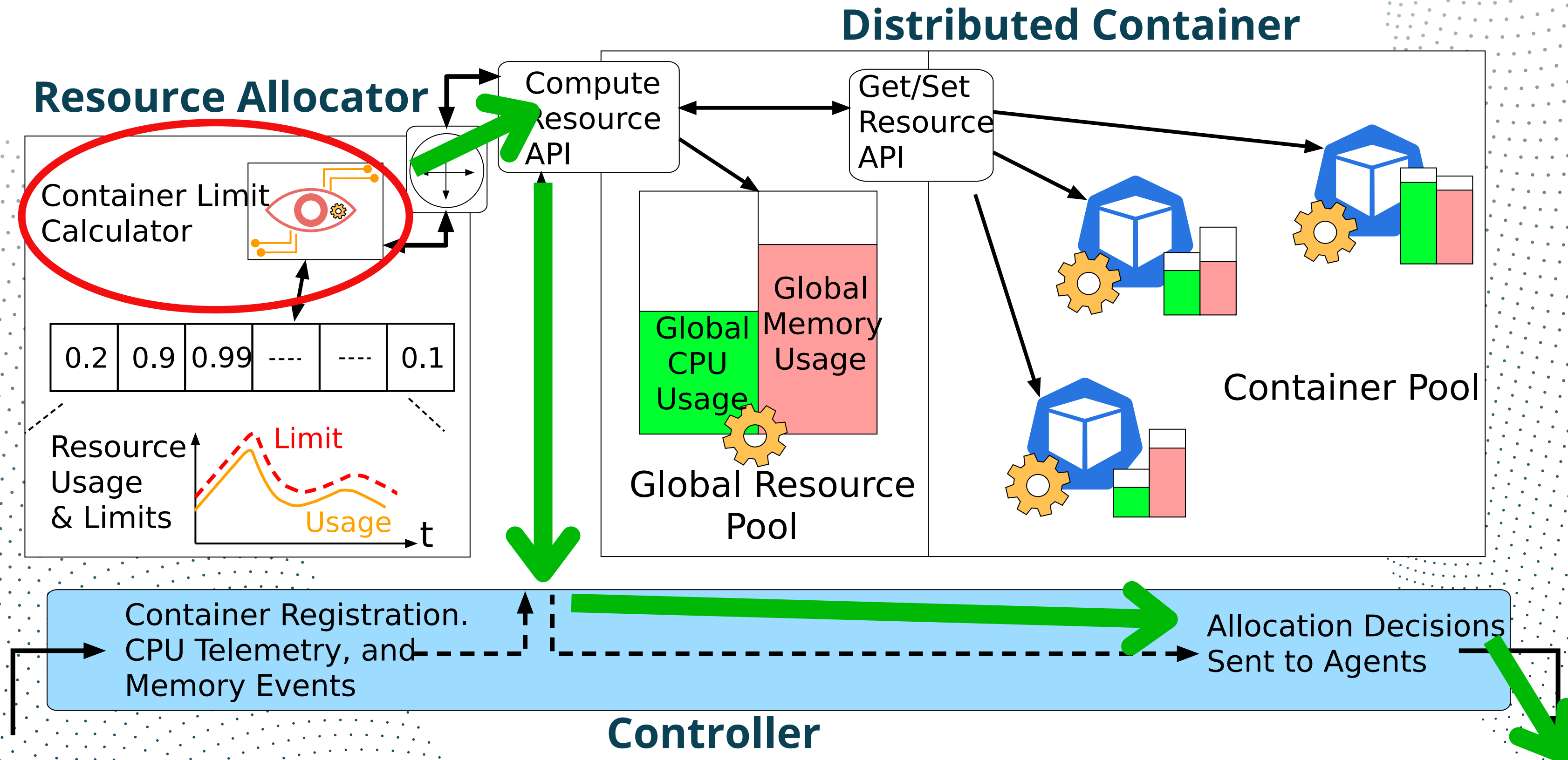


Escra - Control Node

Distributed Container



Escra - Control Node





Evaluation - Microservices



Evaluation - Microservices

Metrics

Absolute Slack

Container Limit -
Container Usage

Application Throughput

Successful requests/second

Application 99.9%ile

Latency

99.9%ile end-to-end latency

Applications & Workloads

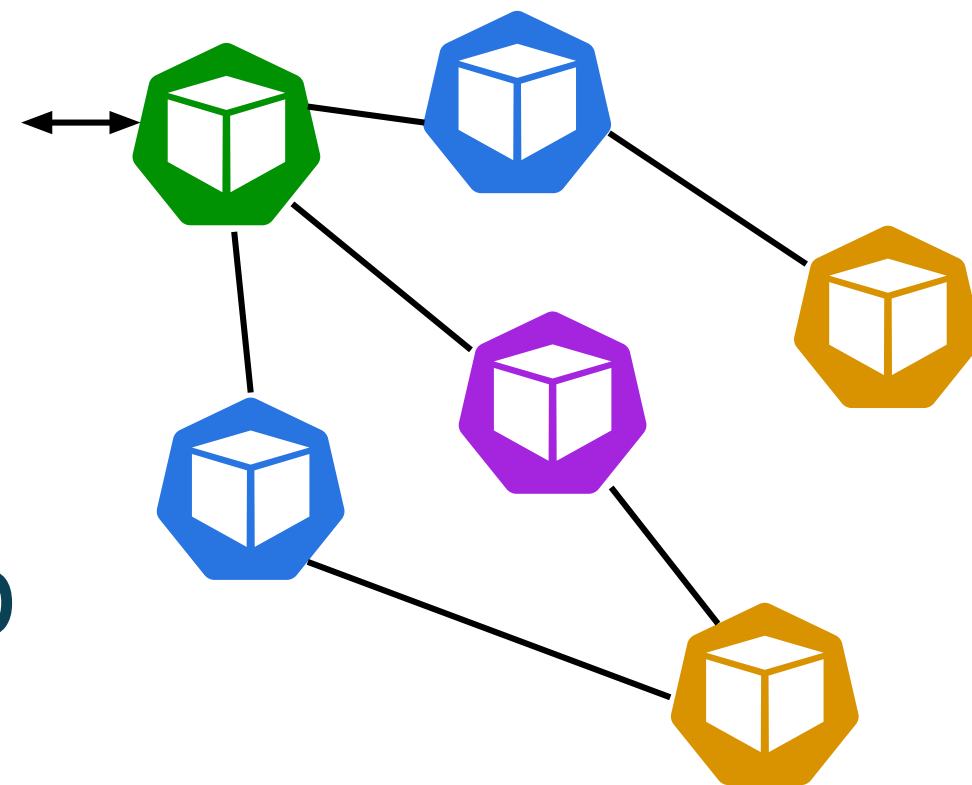
Applications

- TrainTicket

- TeaStore

- HipsterShop

- MediaMicroservice

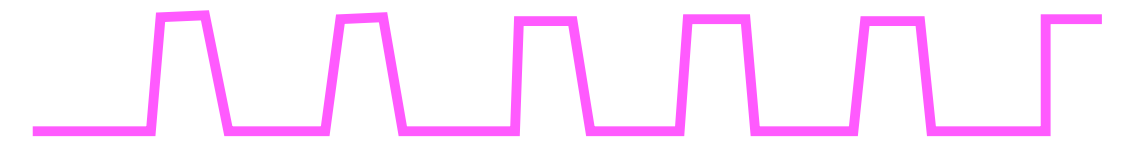


Workloads

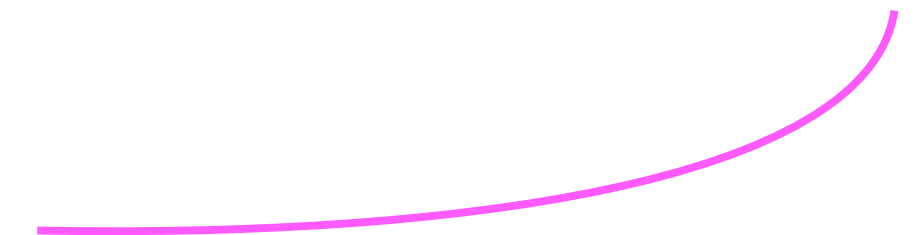
- Alibaba



- Burst



- Exponential

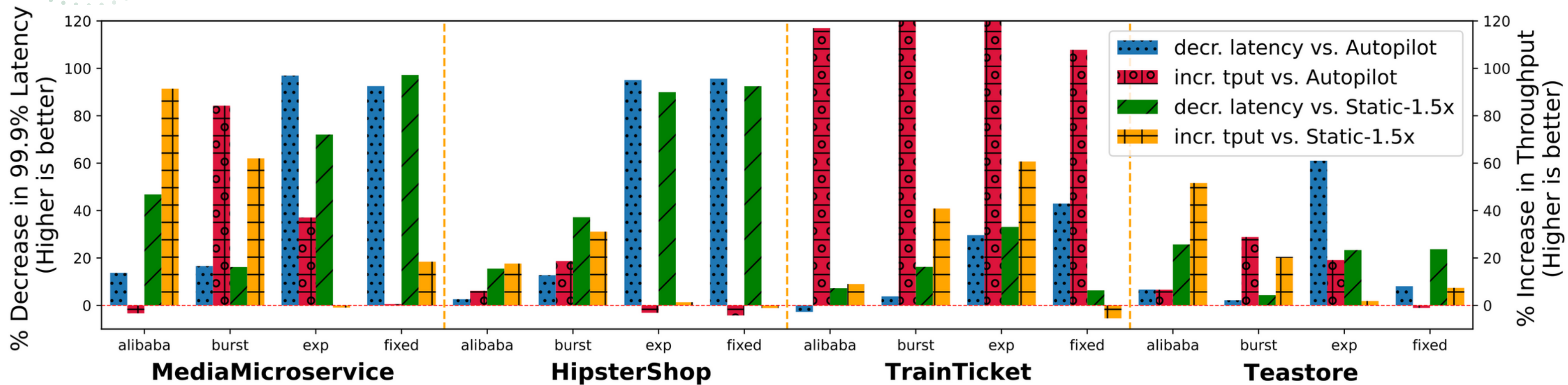


- Fixed



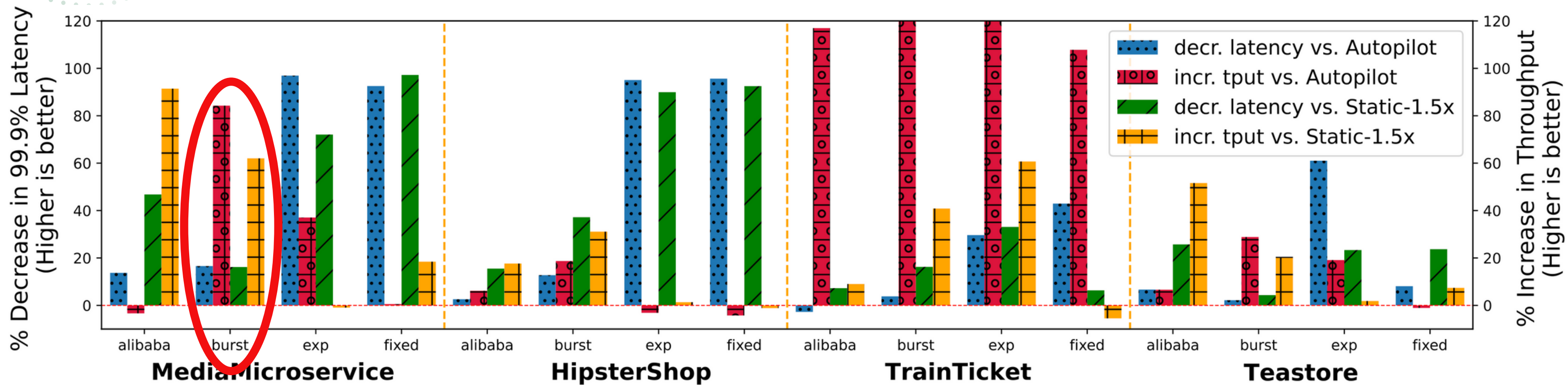
Change in Latency and Throughput

Microservices



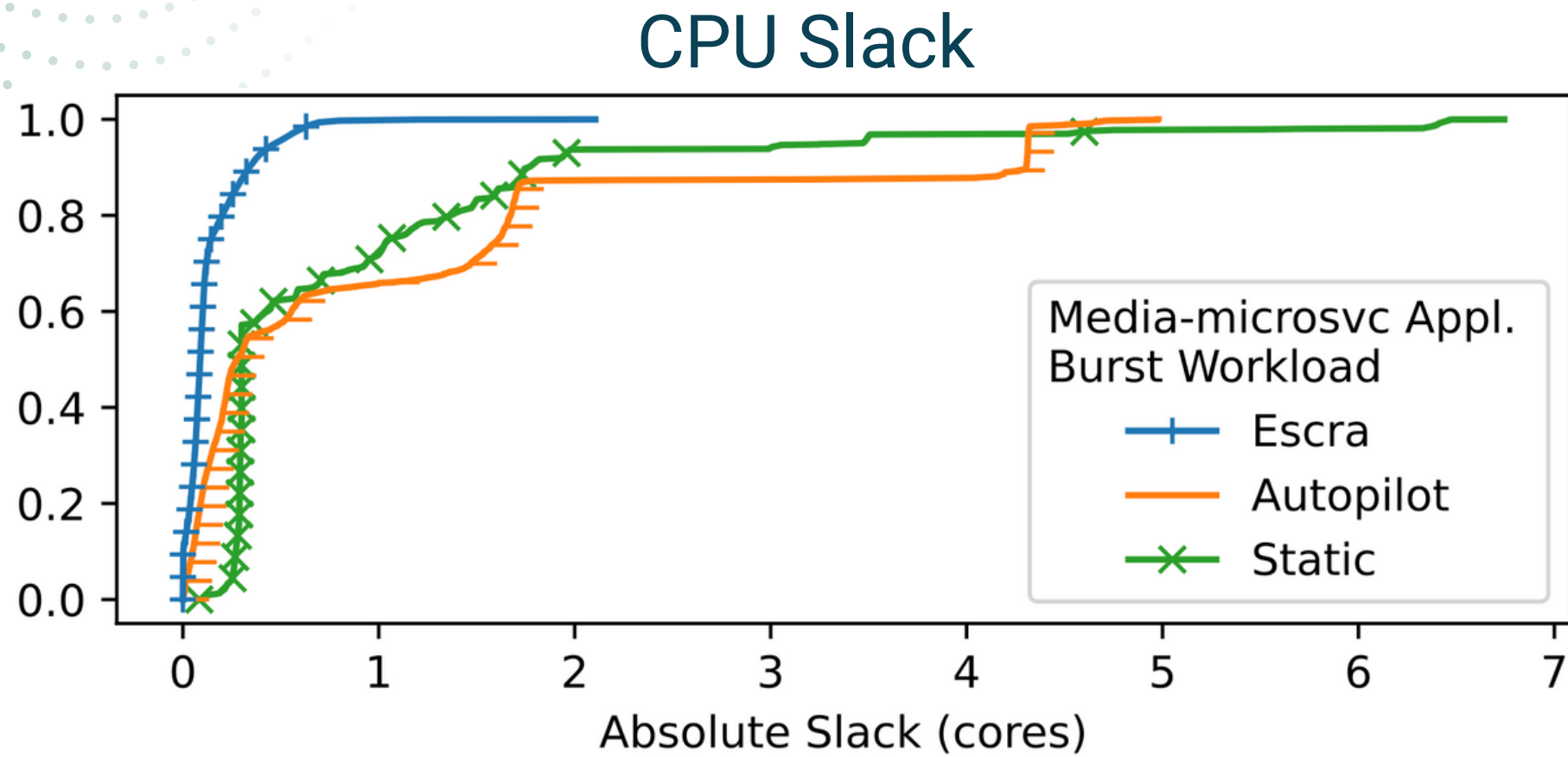
Change in Latency and Throughput

Microservices



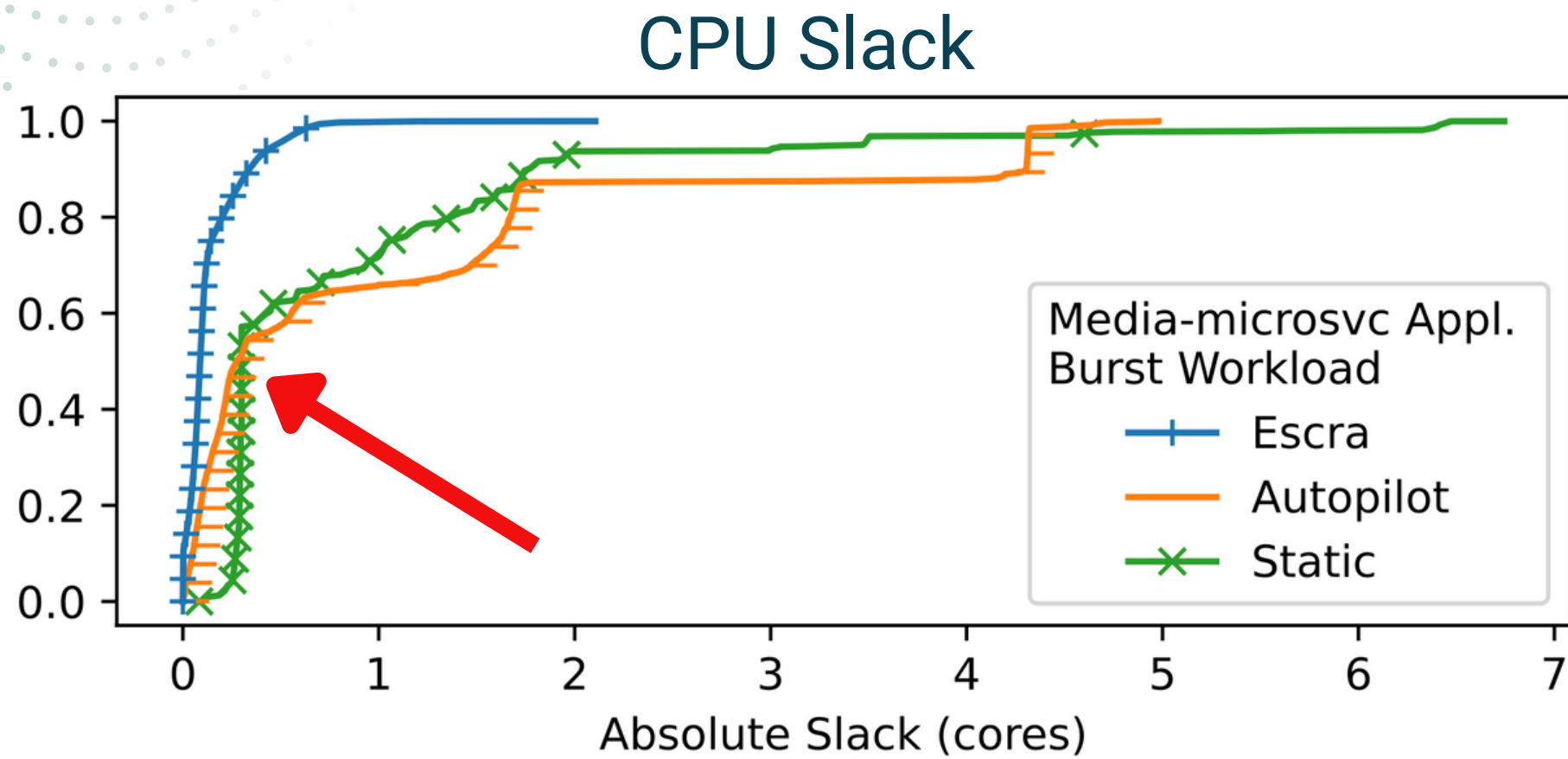
Escra vs. Autopilot

Microservices - MediaMicroservice/Burst Workload



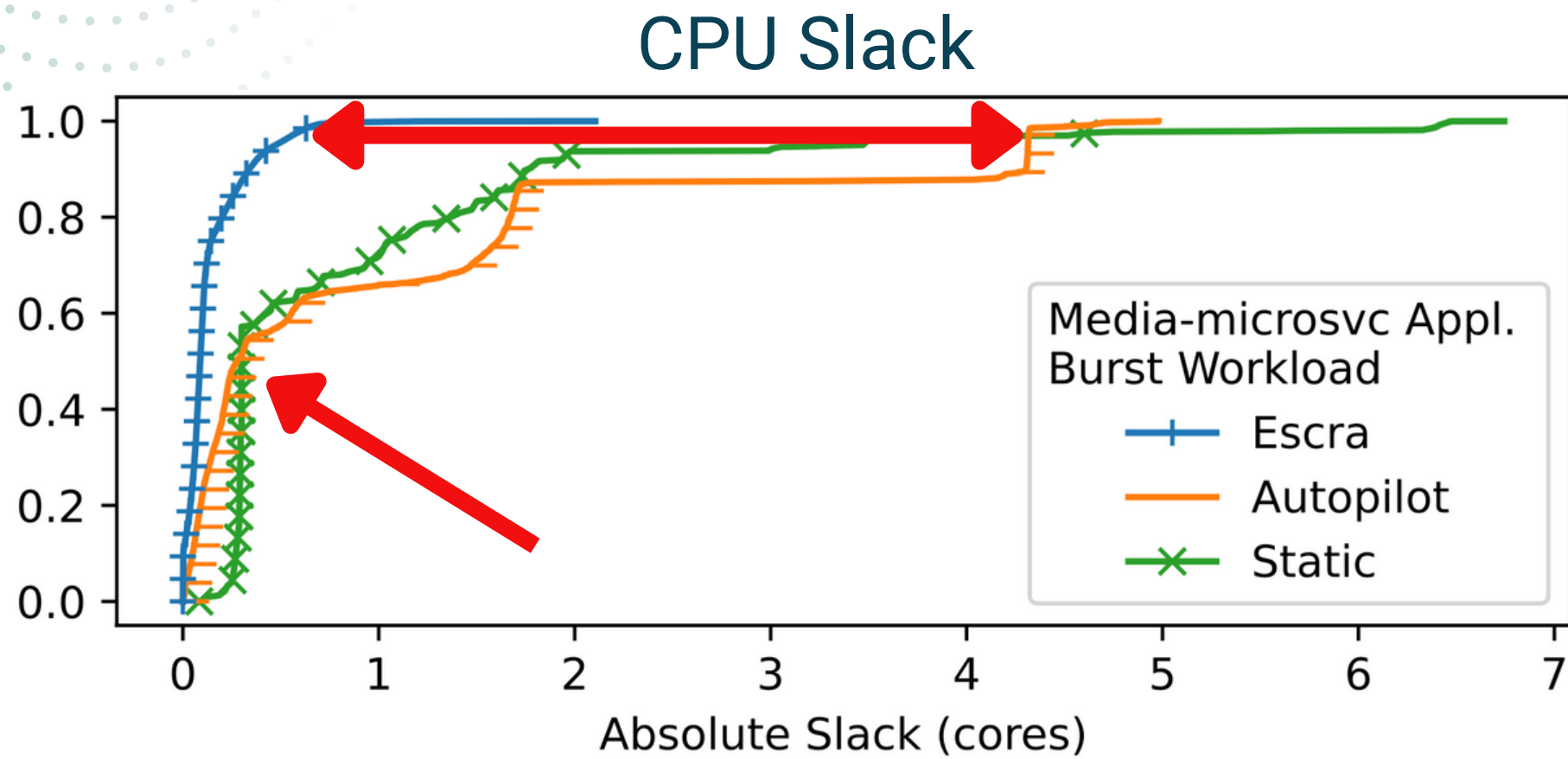
Escra vs. Autopilot

Microservices - MediaMicroservice/Burst Workload



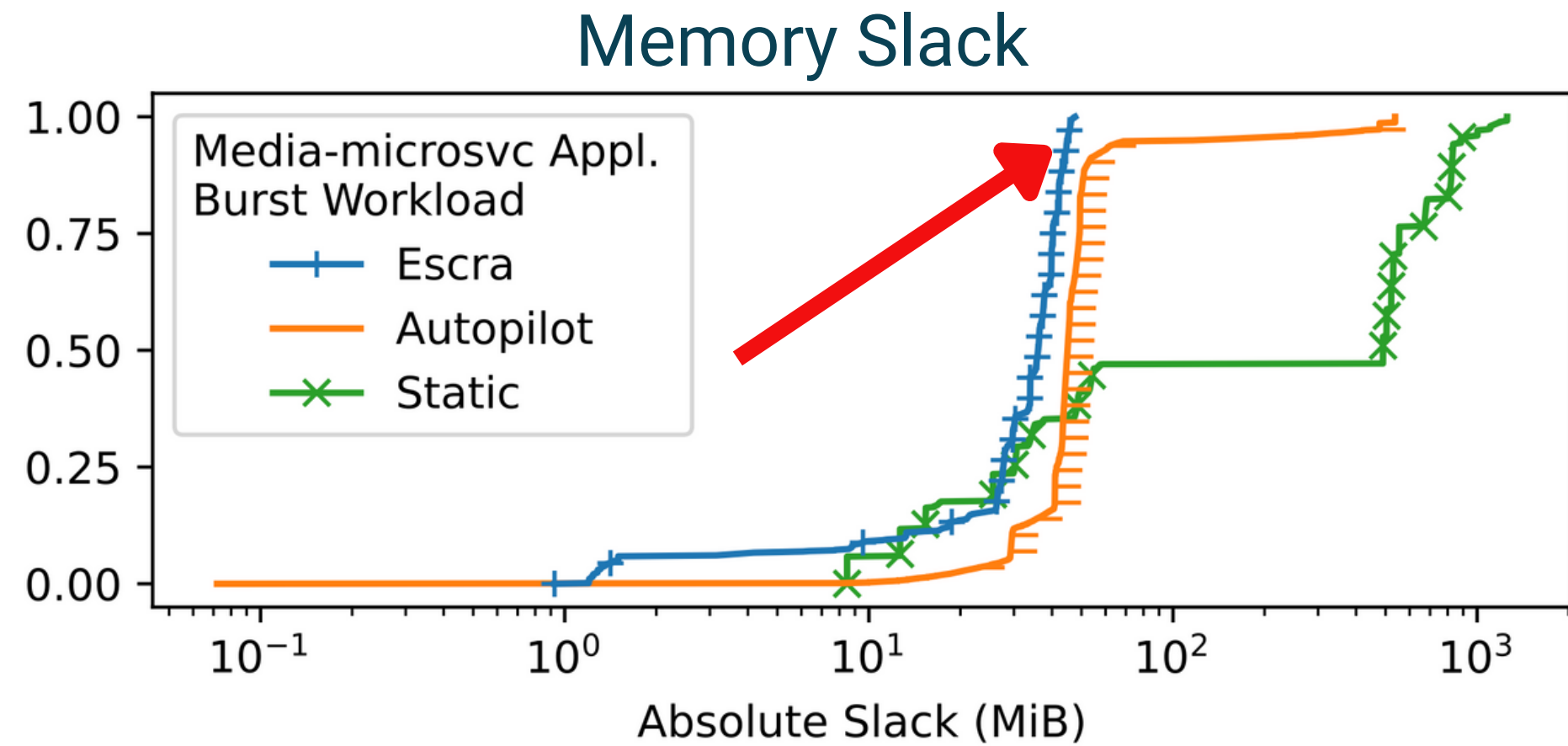
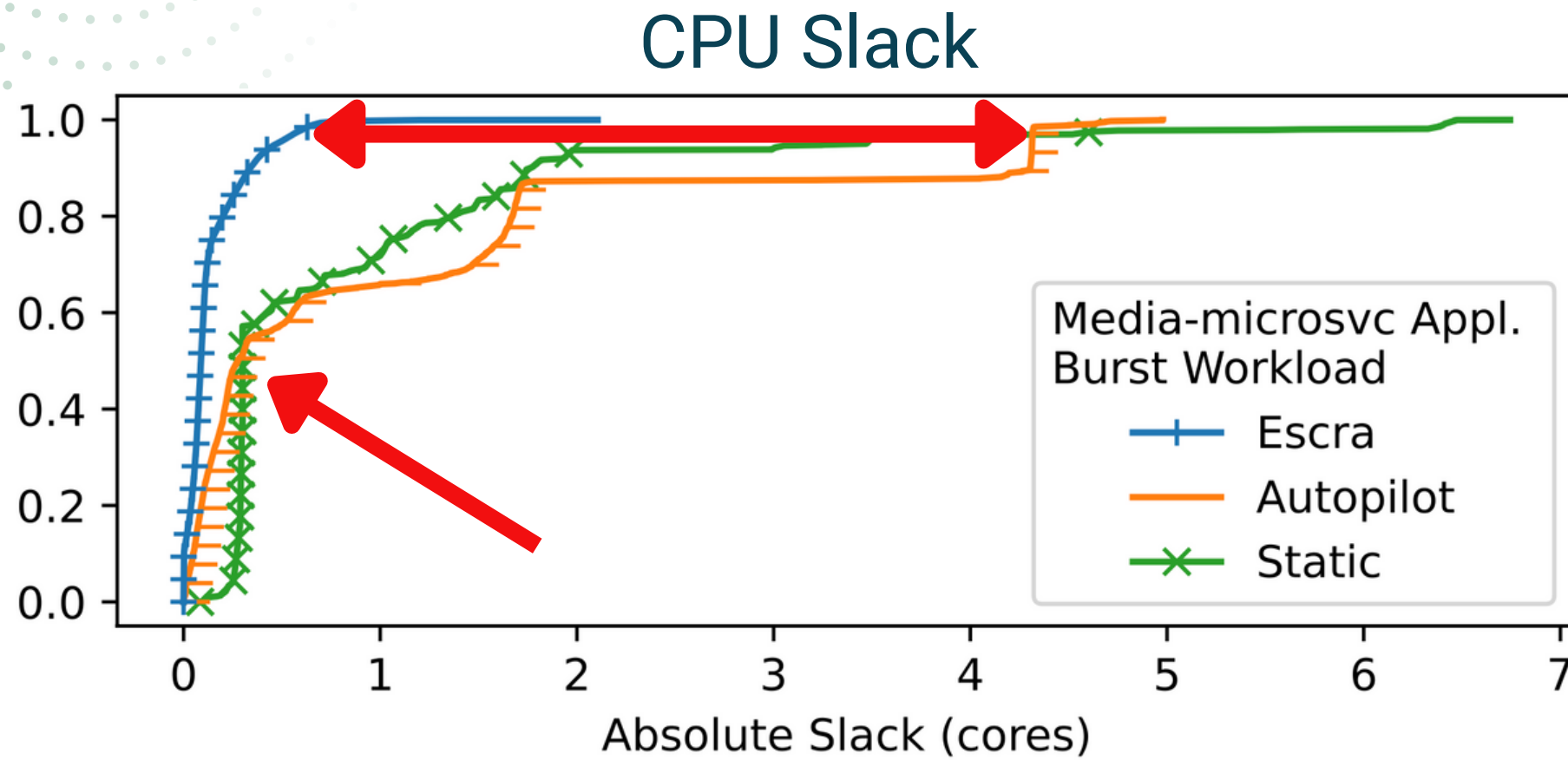
Escra vs. Autopilot

Microservices - MediaMicroservice/Burst Workload



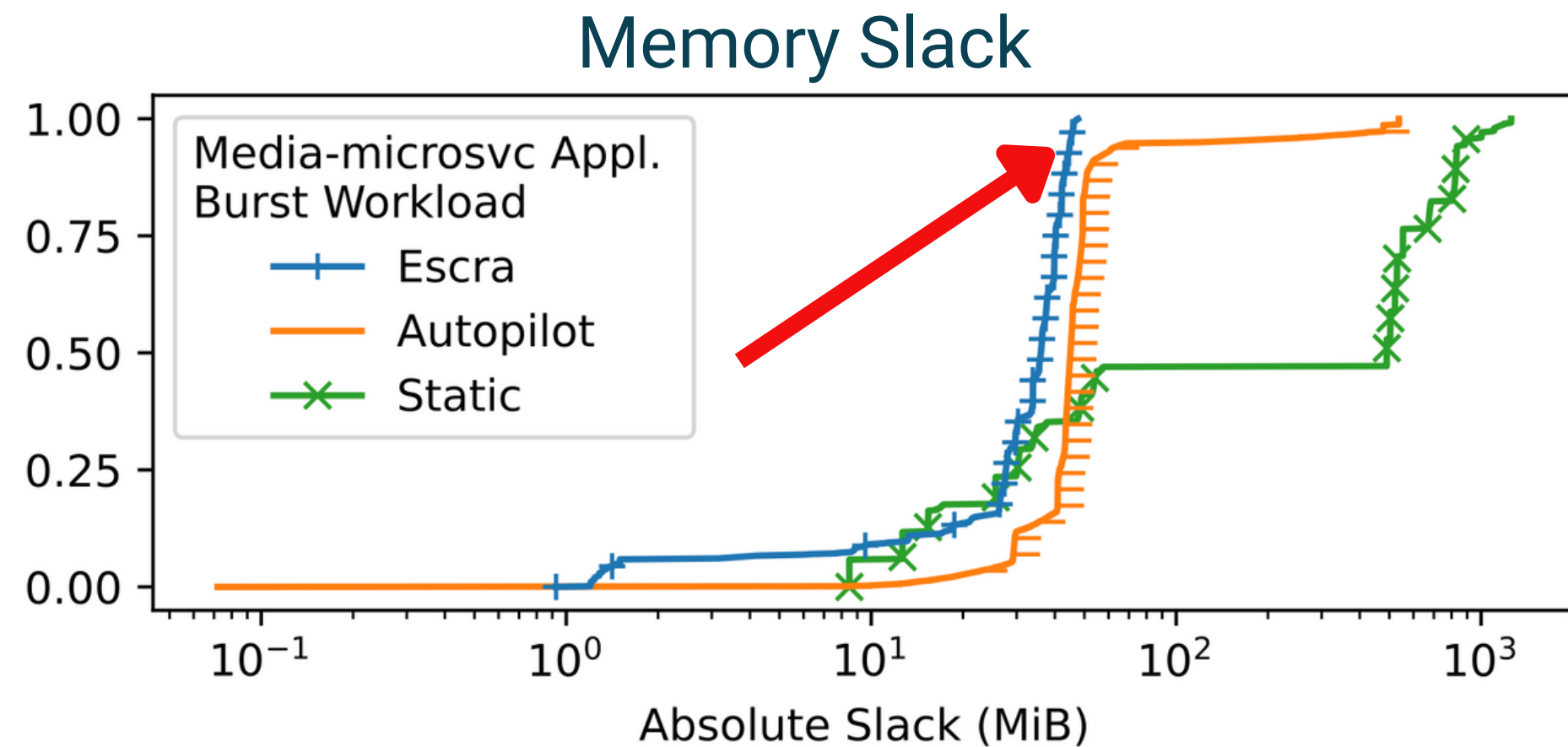
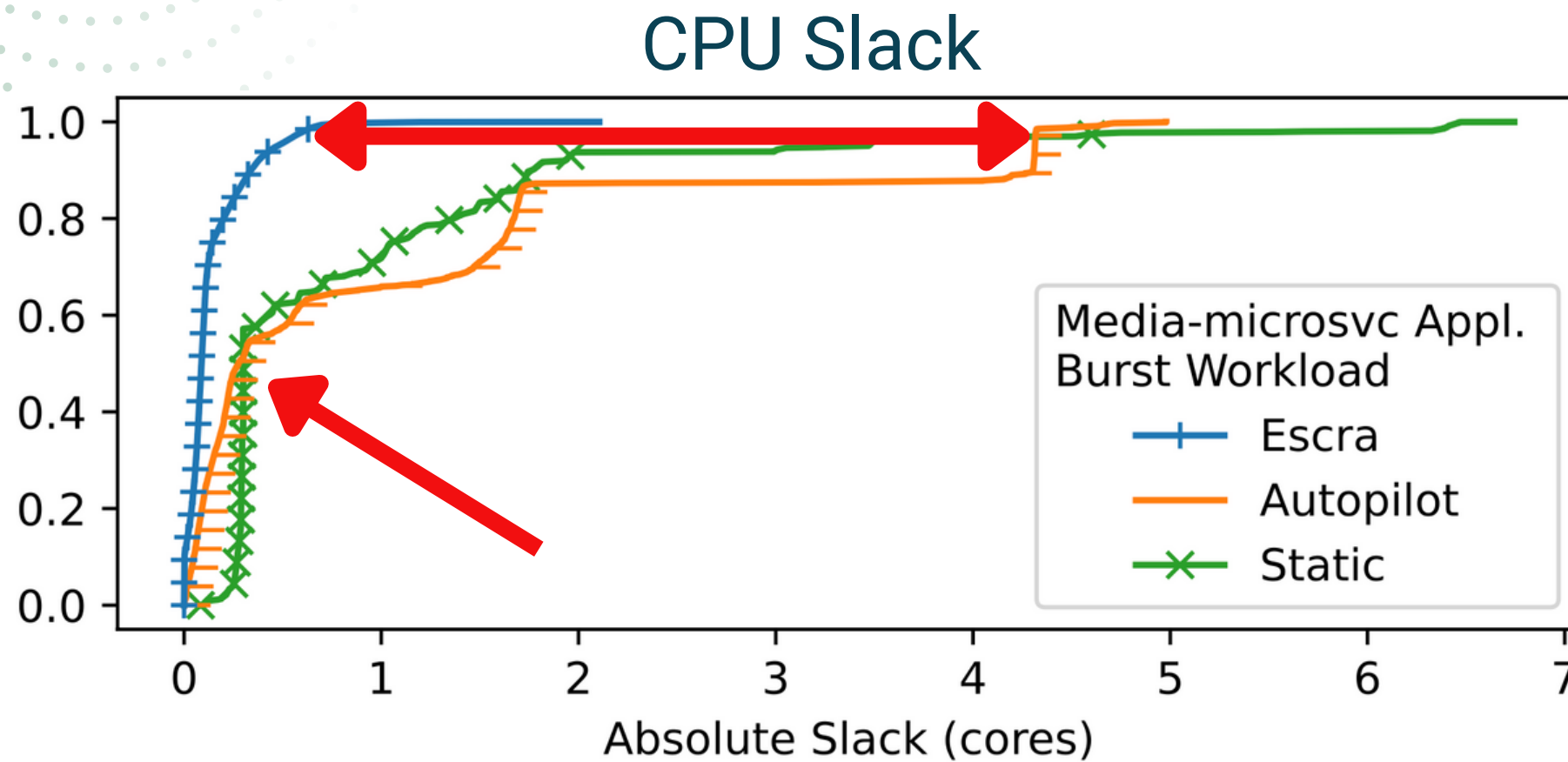
Escra vs. Autopilot

Microservices - MediaMicroservice/Burst Workload



Escra vs. Autopilot

Microservices - MediaMicroservice/Burst Workload



- **Latency Decrease: 16.6%**
- **Throughput Increase: 84.3%**



Evaluation - Serverless





Evaluation - Serverless

Metrics

Aggregate Limits

Σ container limits -
 Σ container usages

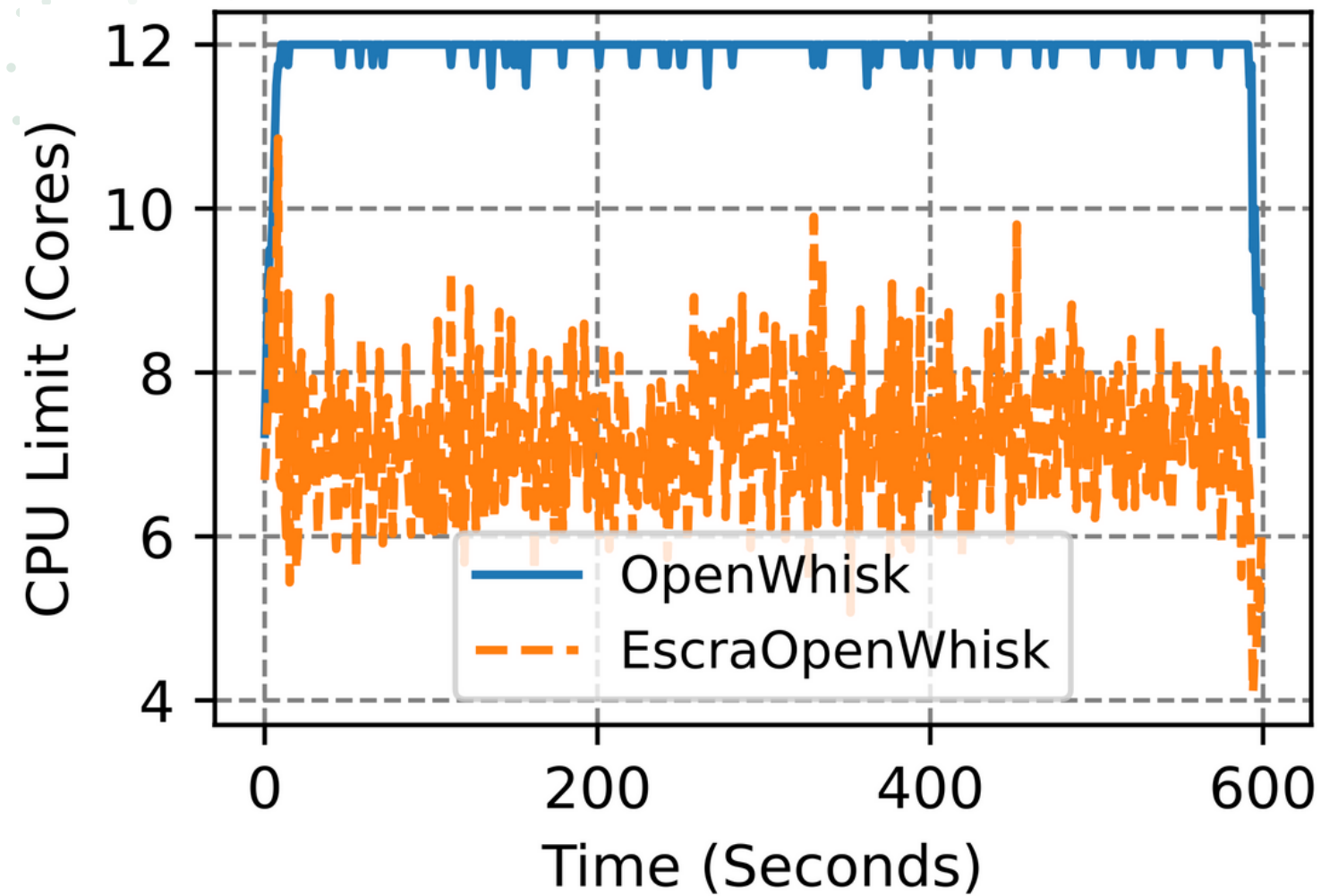
Application Latency

End-to-end latency per
request or job

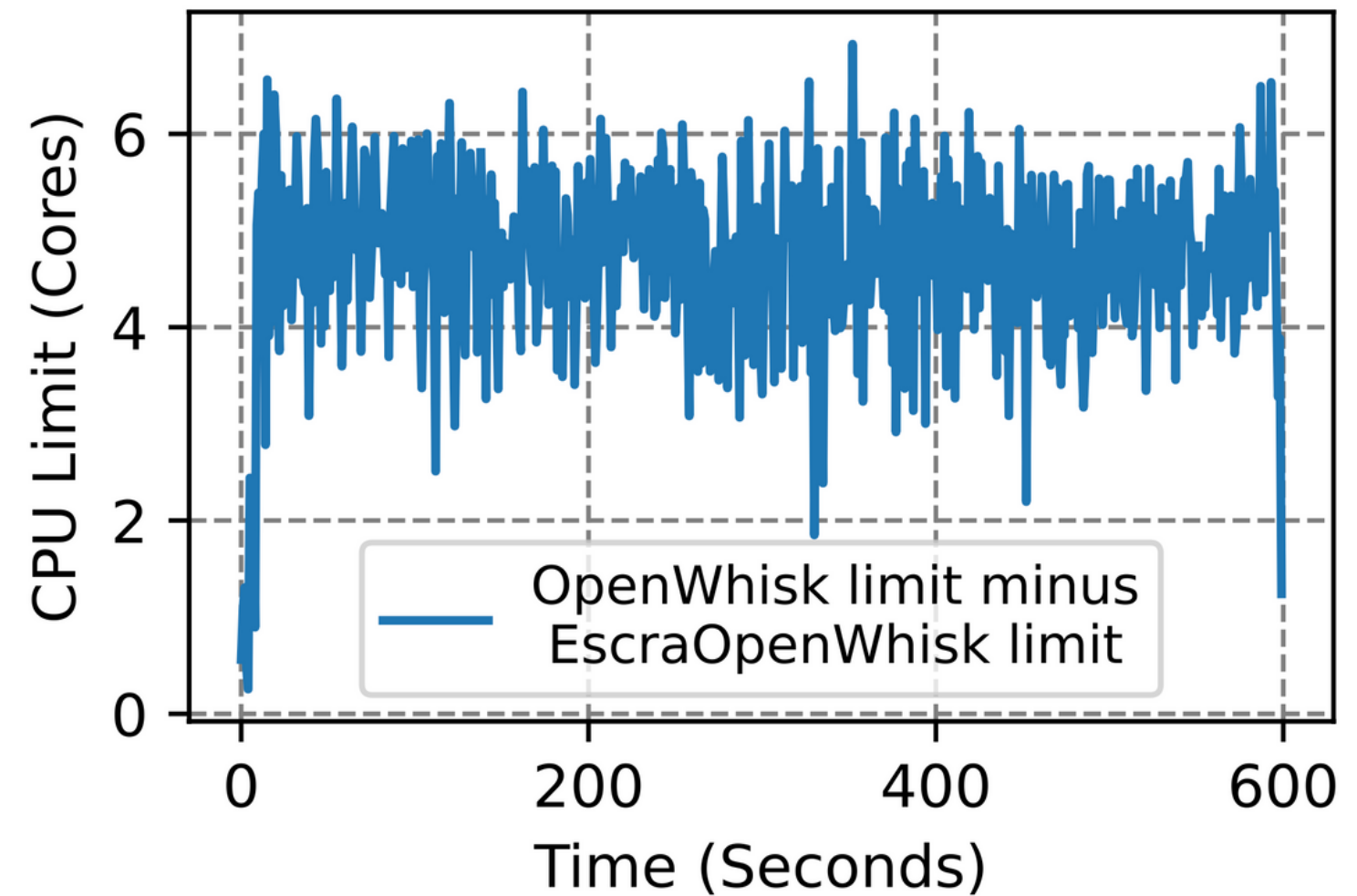
ImageProcess

Serverless

Aggregate CPU Slack



Aggregate CPU Savings

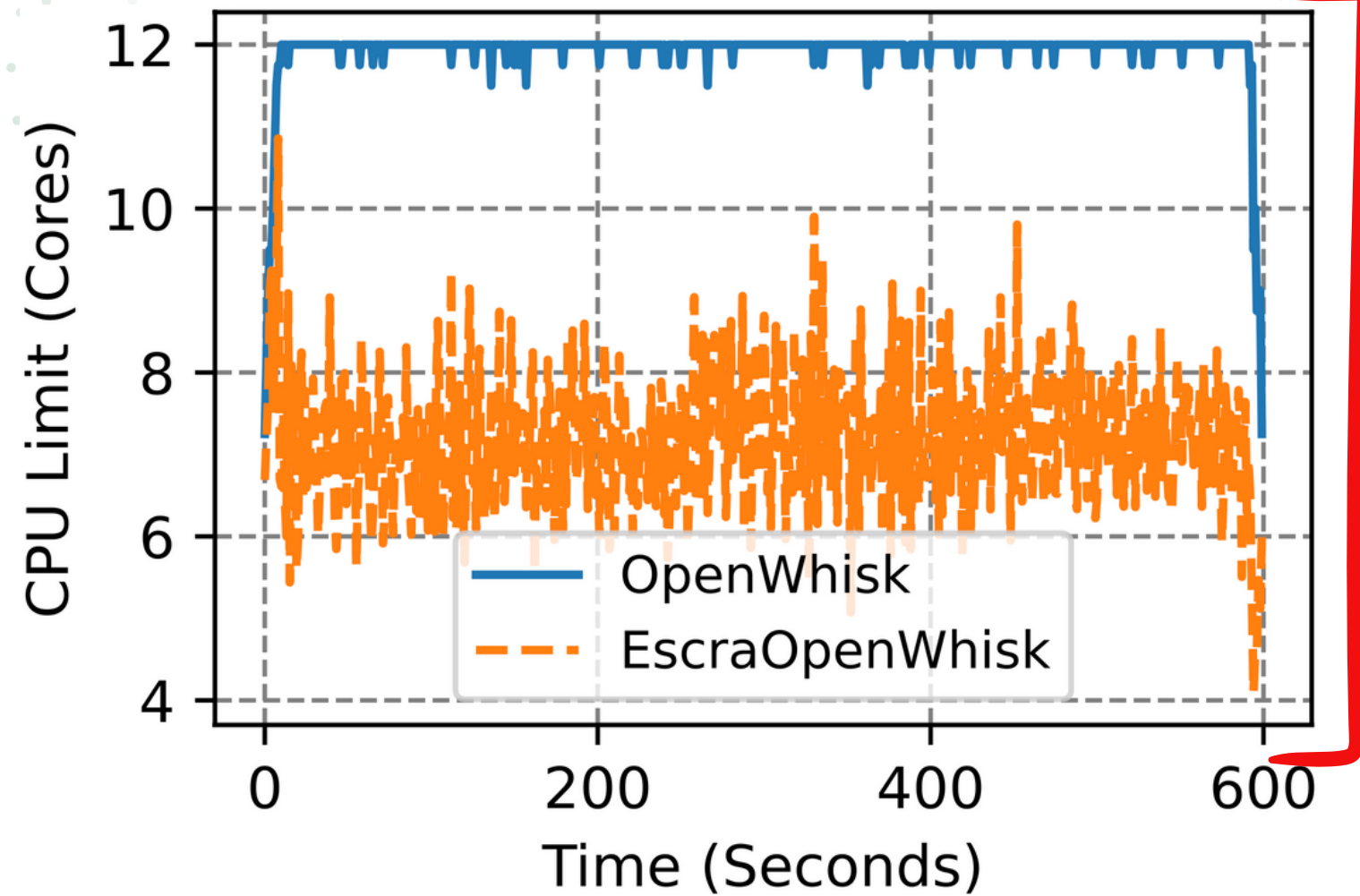


ImageProcess

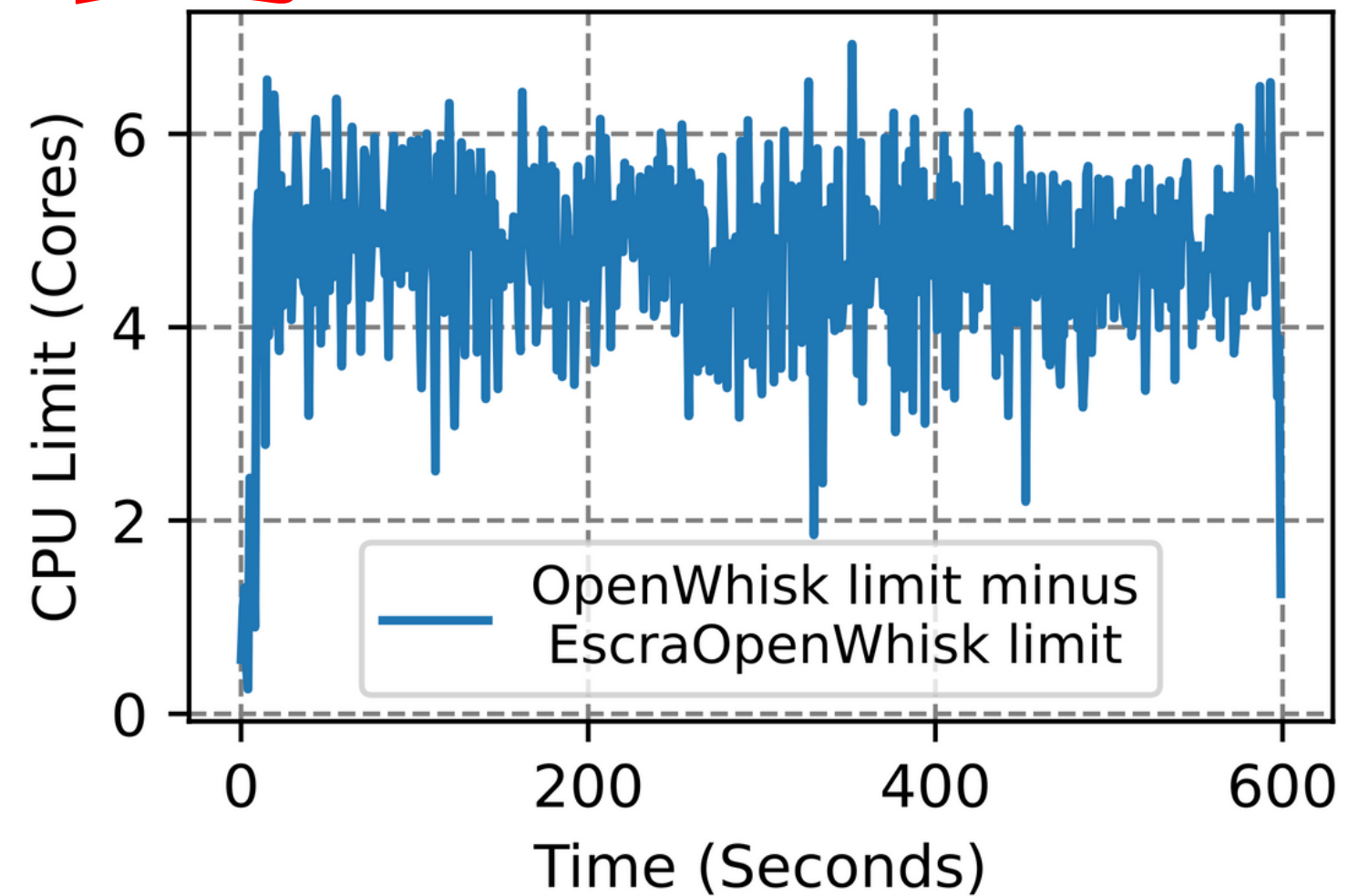
Serverless

Difference!

Aggregate CPU Slack



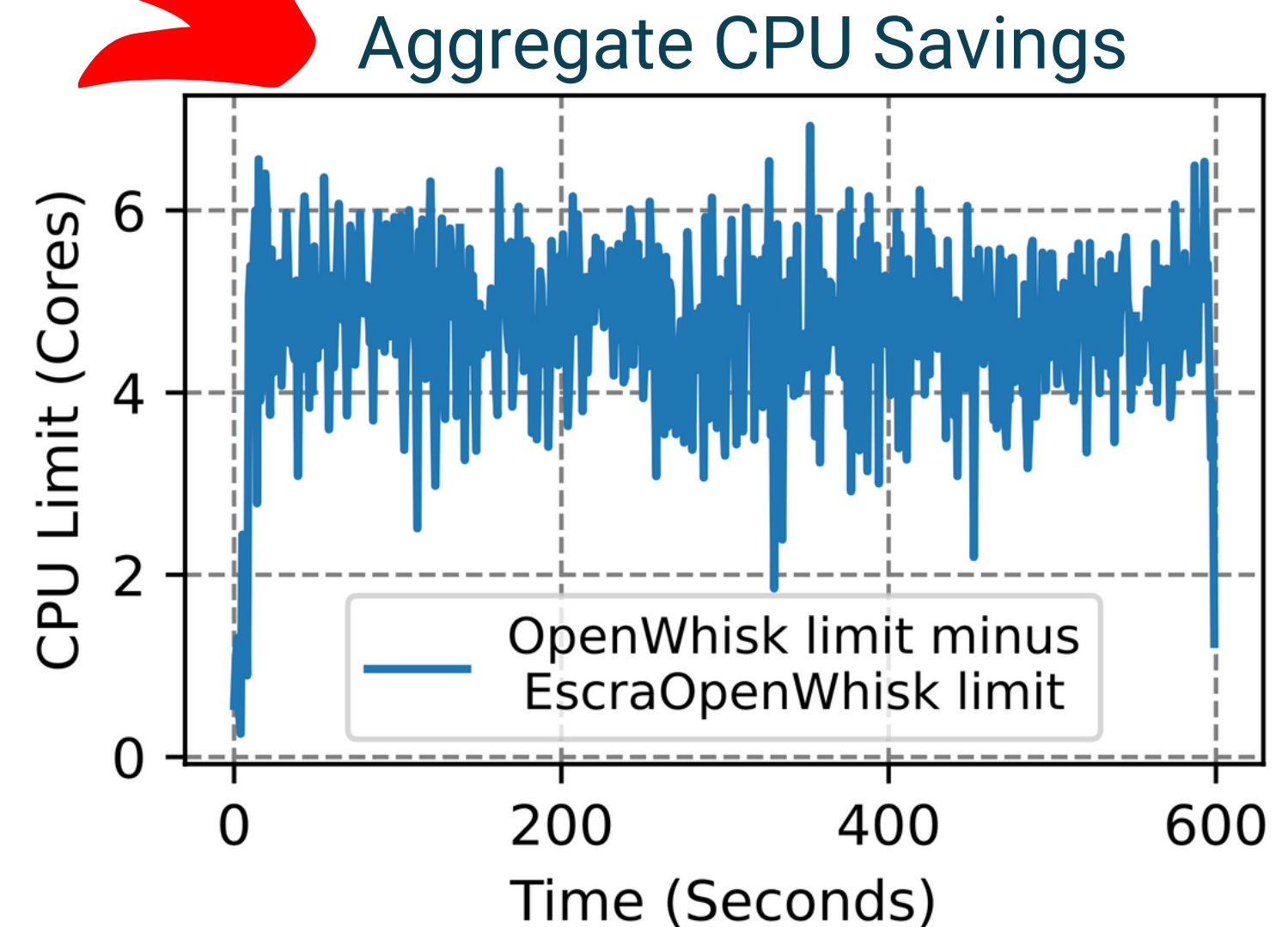
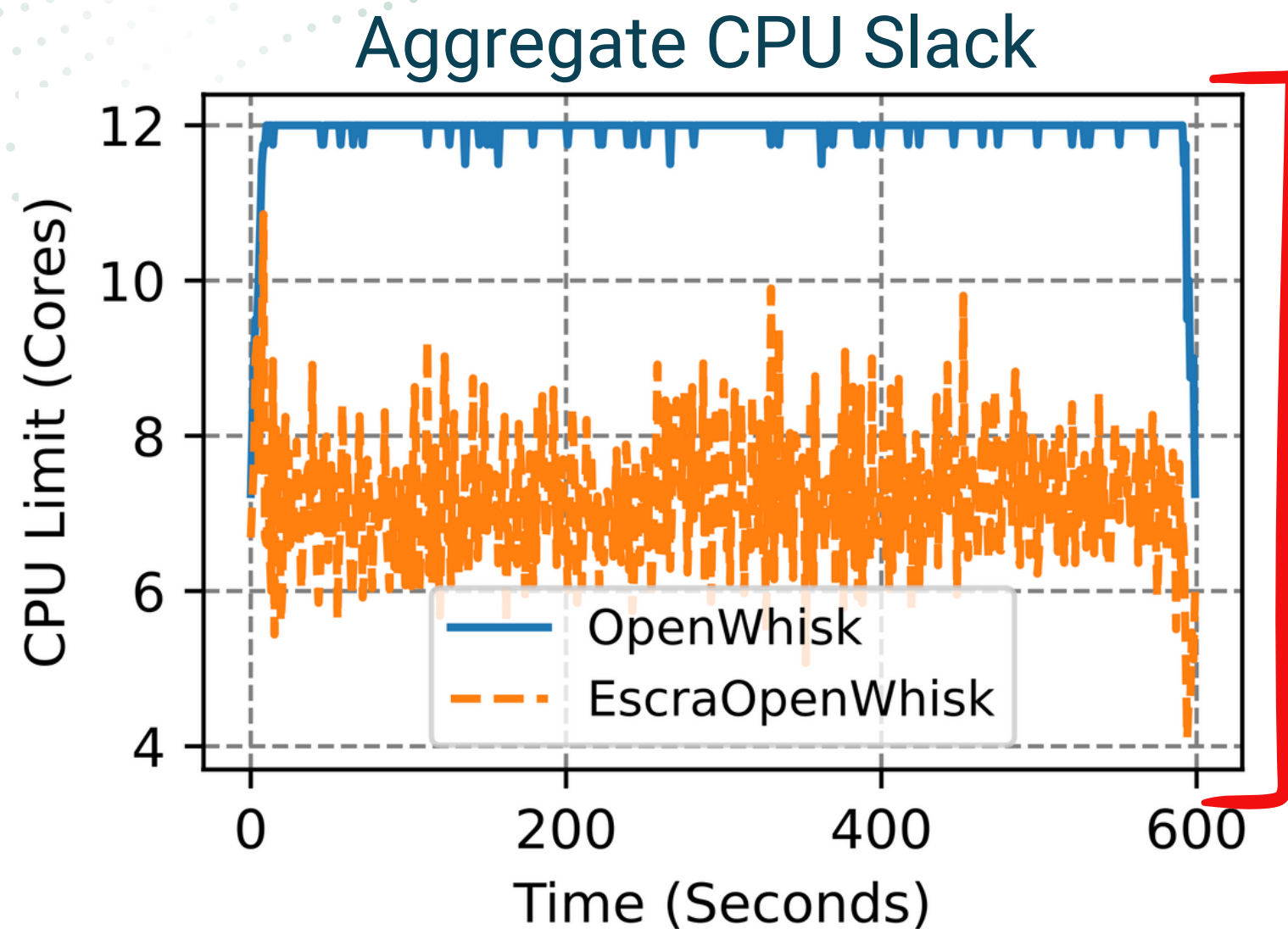
Aggregate CPU Savings



ImageProcess

Serverless

Difference!

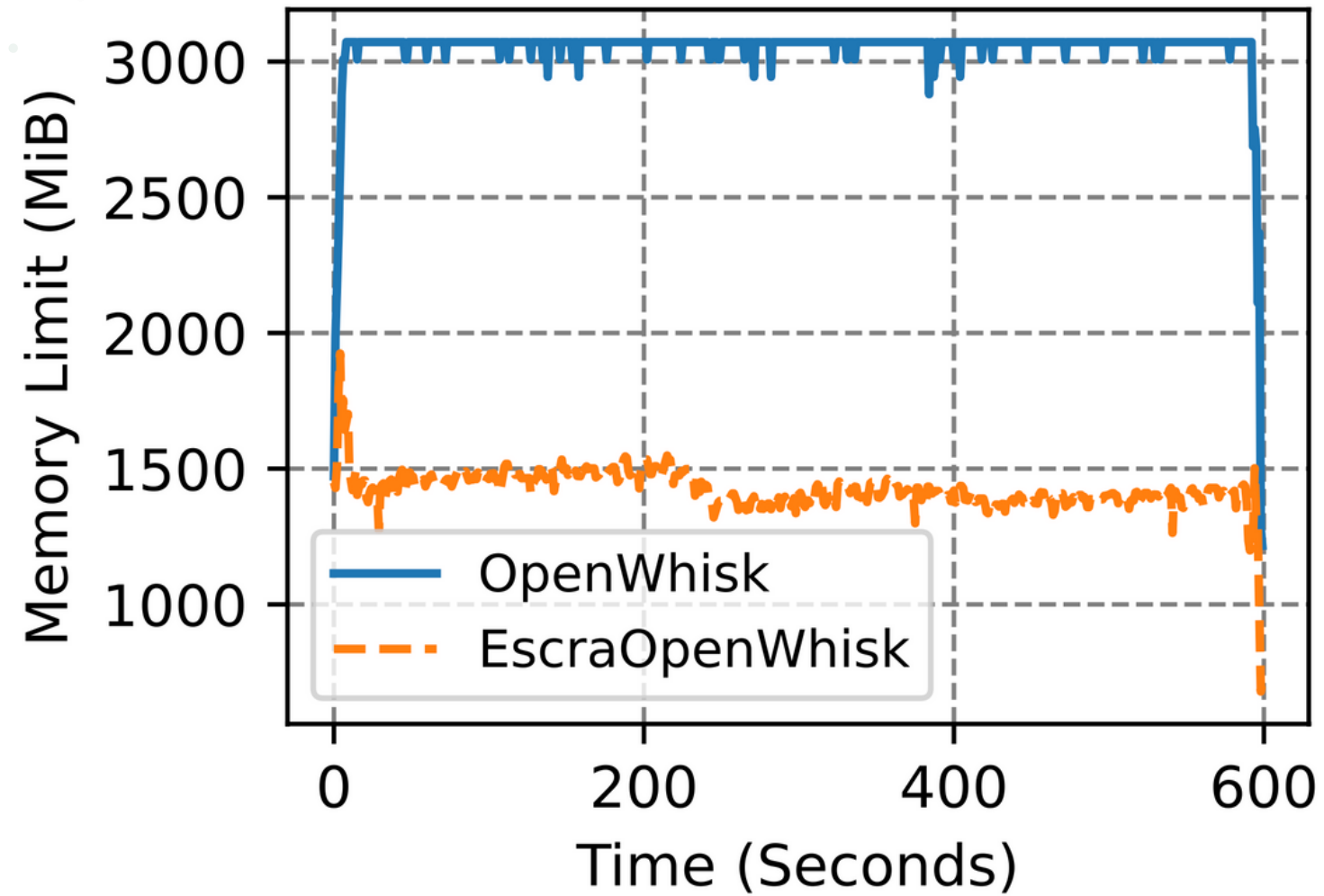


- **Escra + Openwhisk avg vCPU limit: 7 vCPU**
- **Openwhisk avg vCPU limit: 12v CPU**

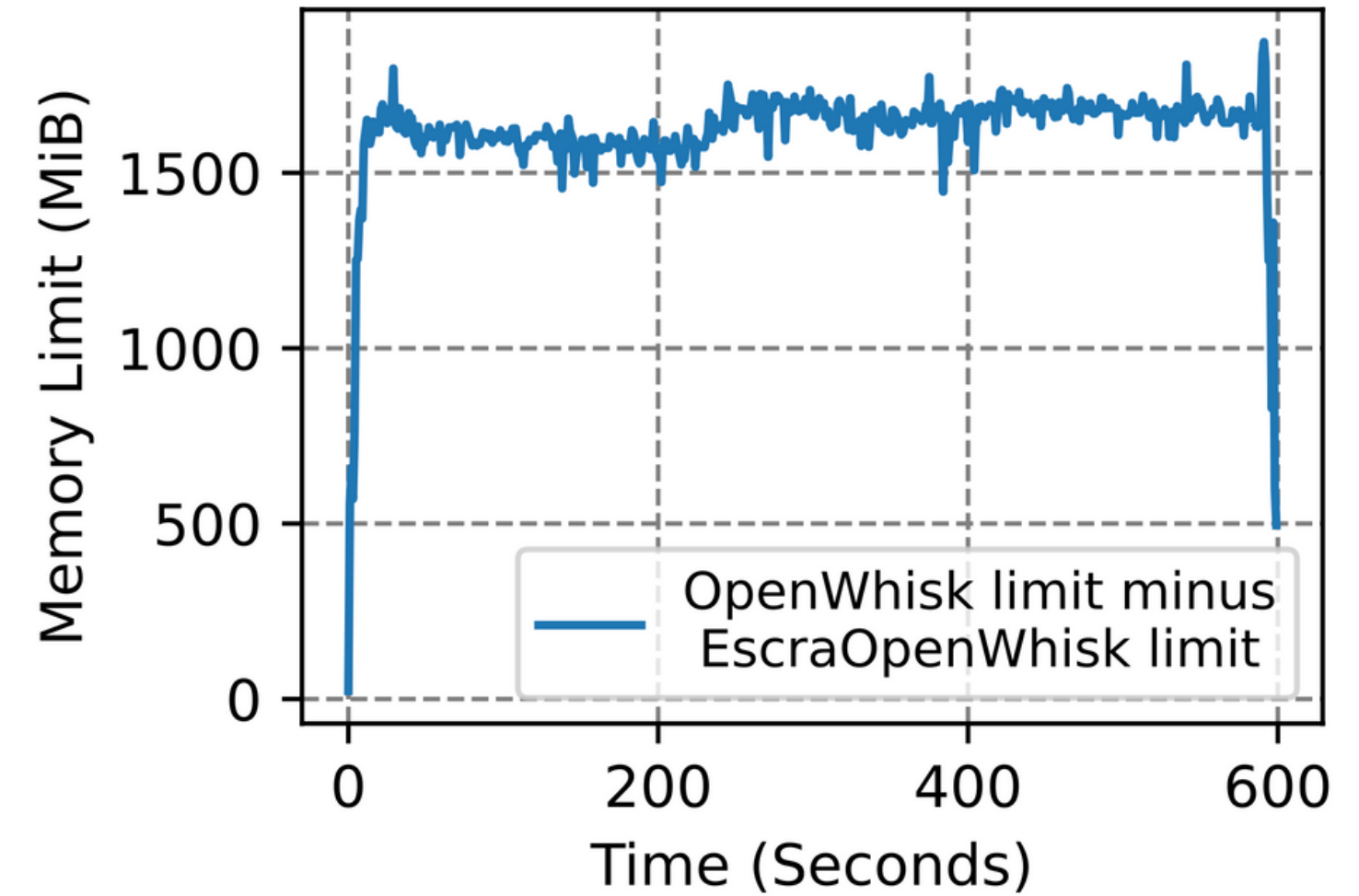
ImageProcess

Serverless

Aggregate Memory Slack



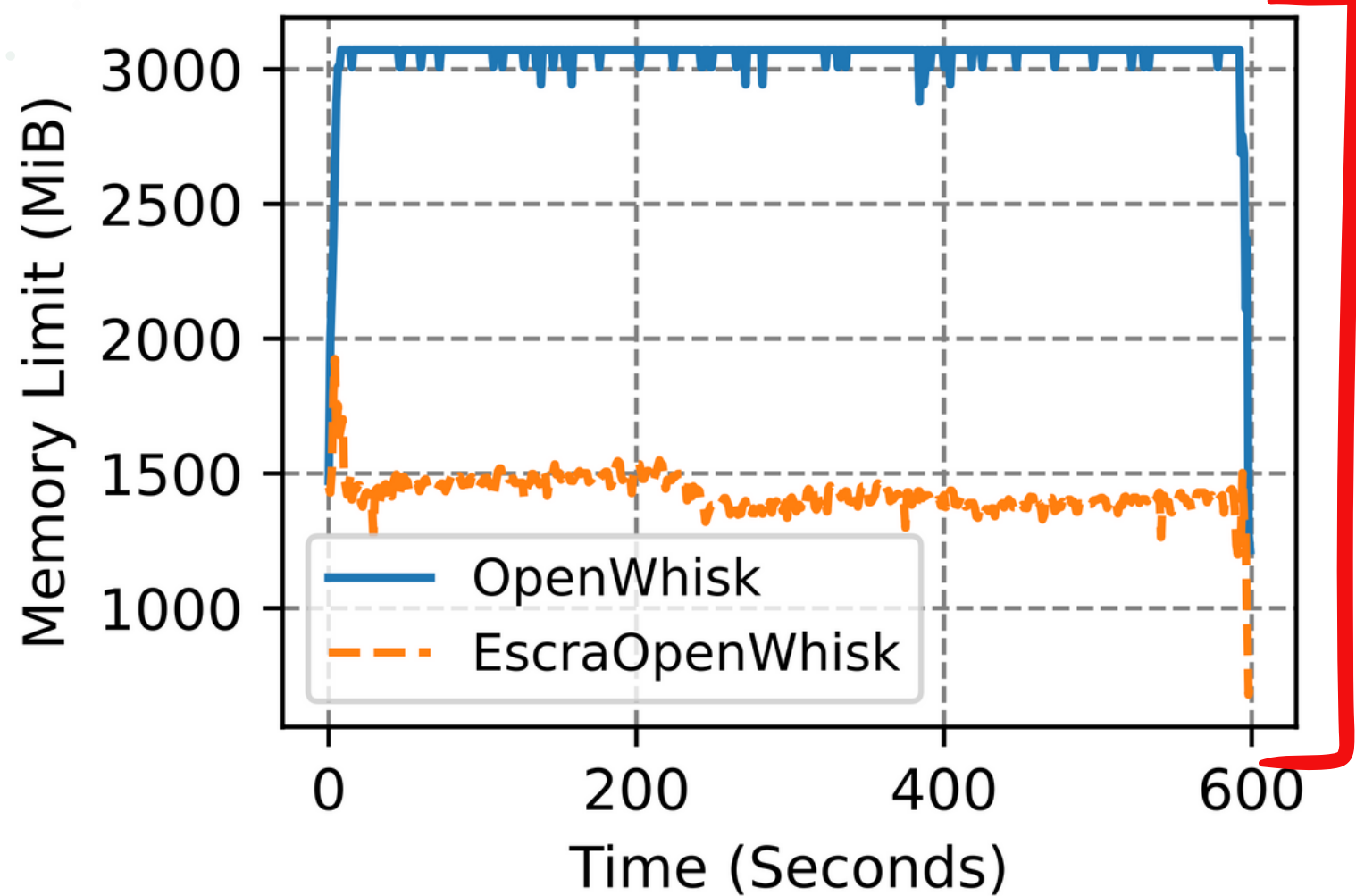
Aggregate Memory Savings



ImageProcess

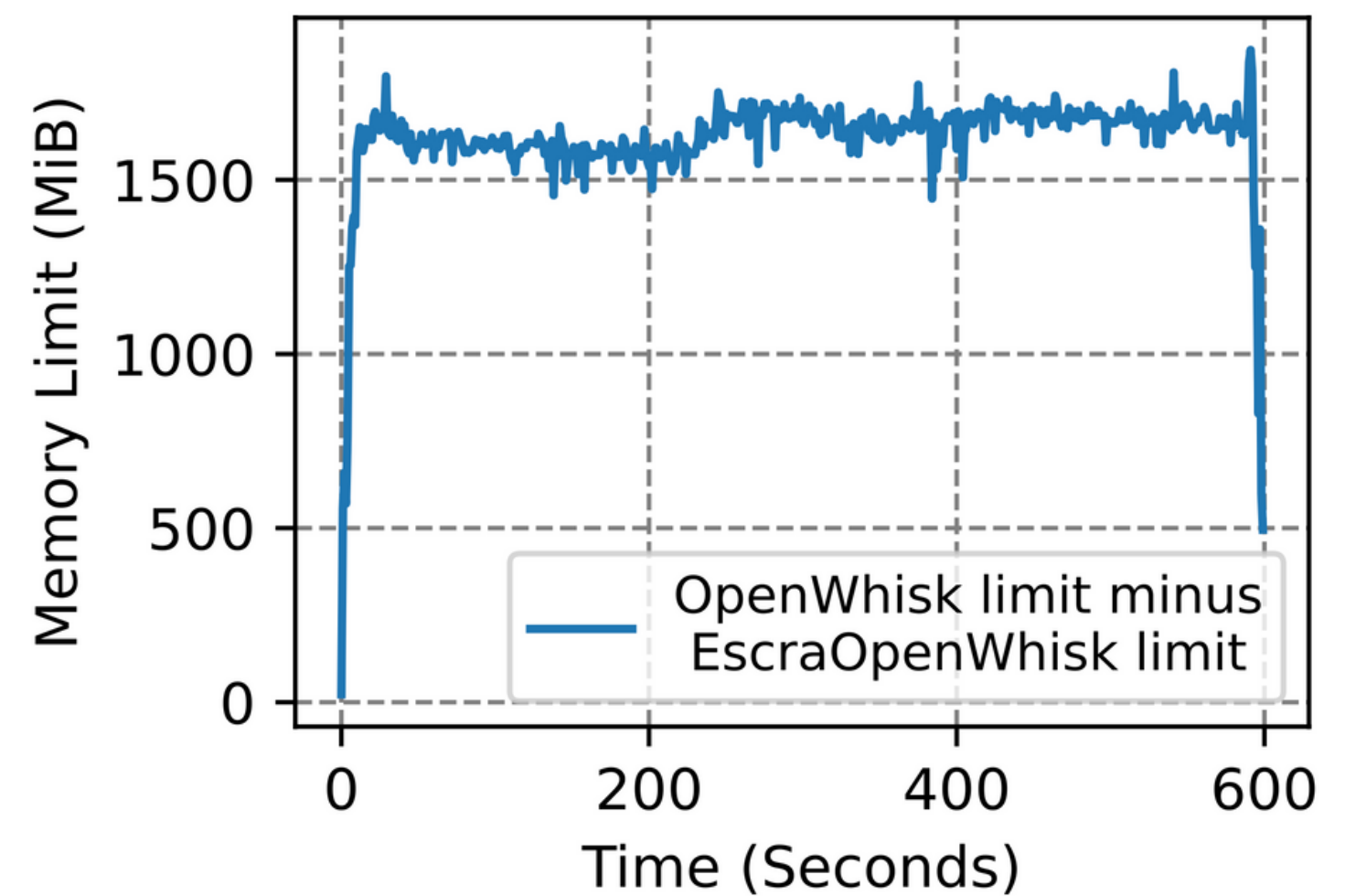
Serverless

Aggregate Memory Slack



Difference!

Aggregate Memory Savings

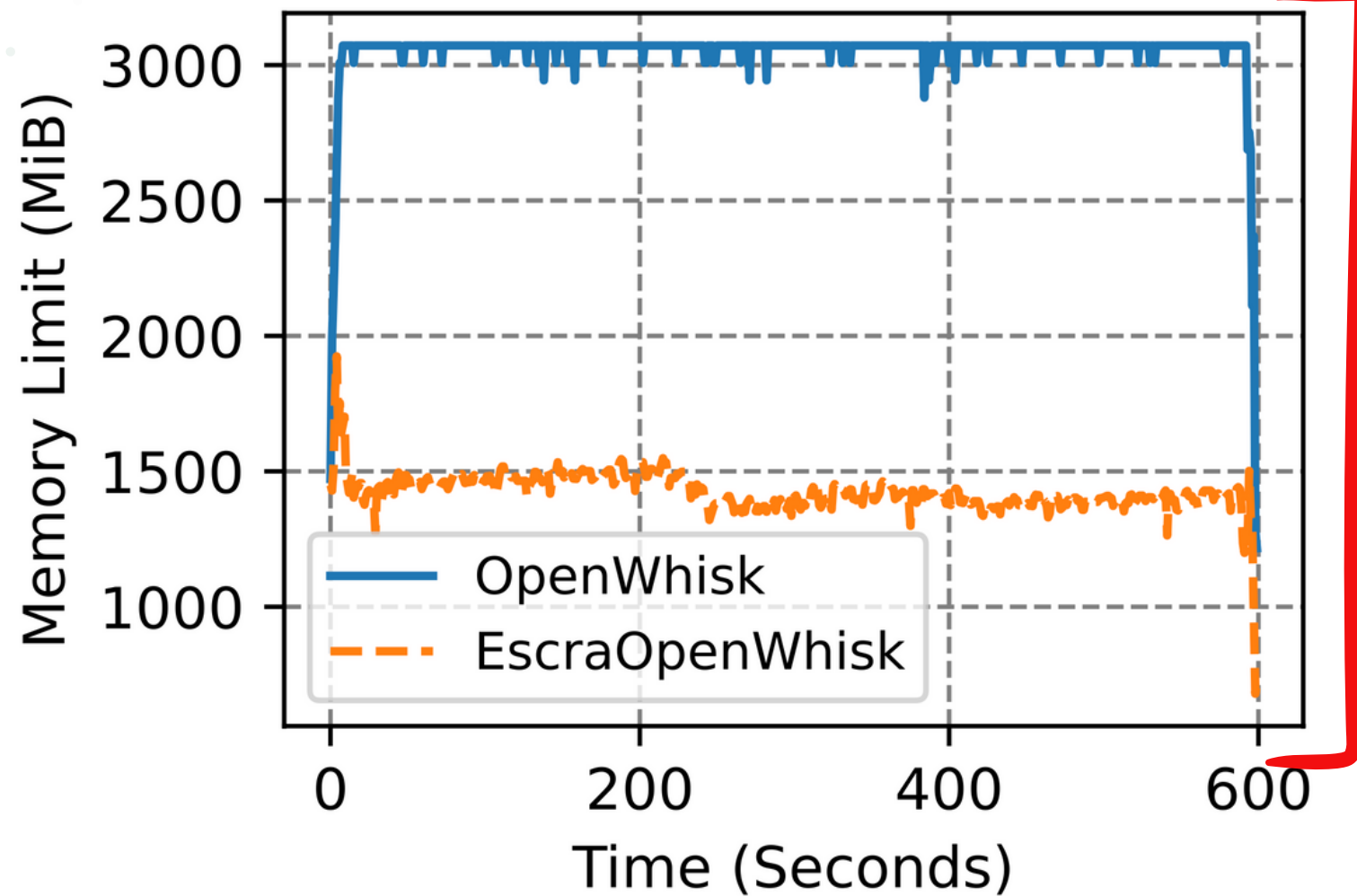


ImageProcess

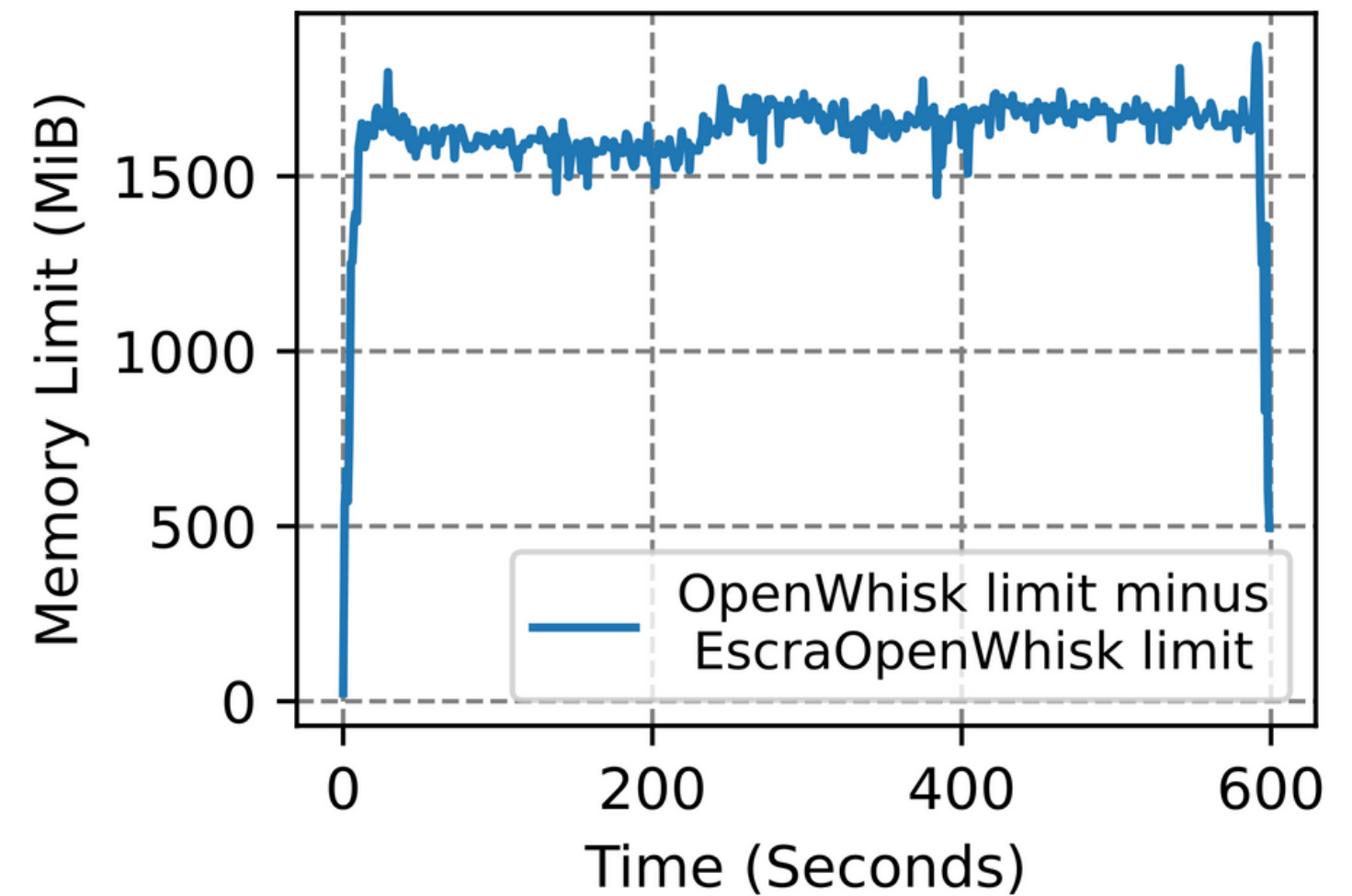
Serverless

Difference!

Aggregate Memory Slack



Aggregate Memory Savings

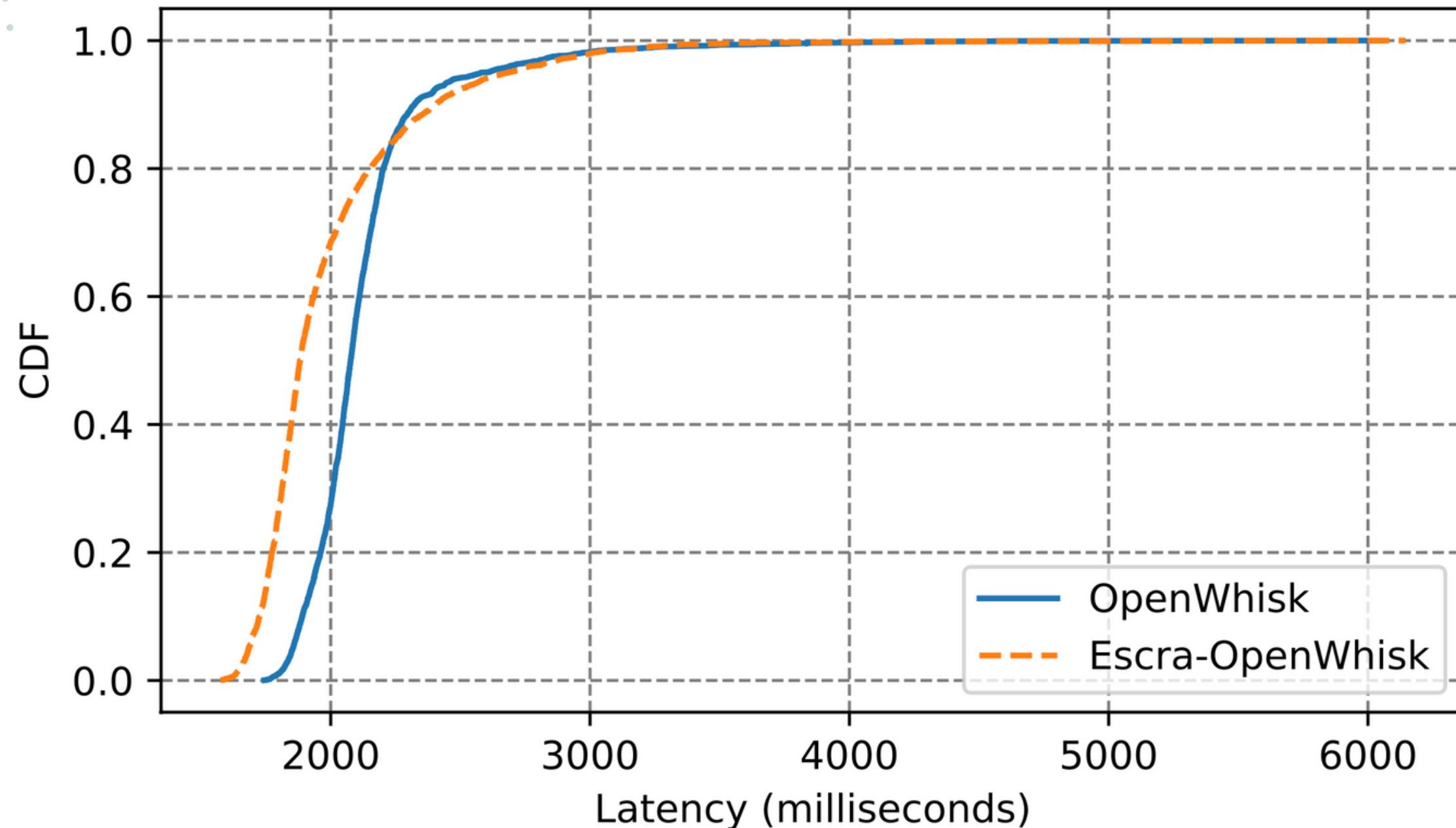


Avg. memory savings: 1550 MiB

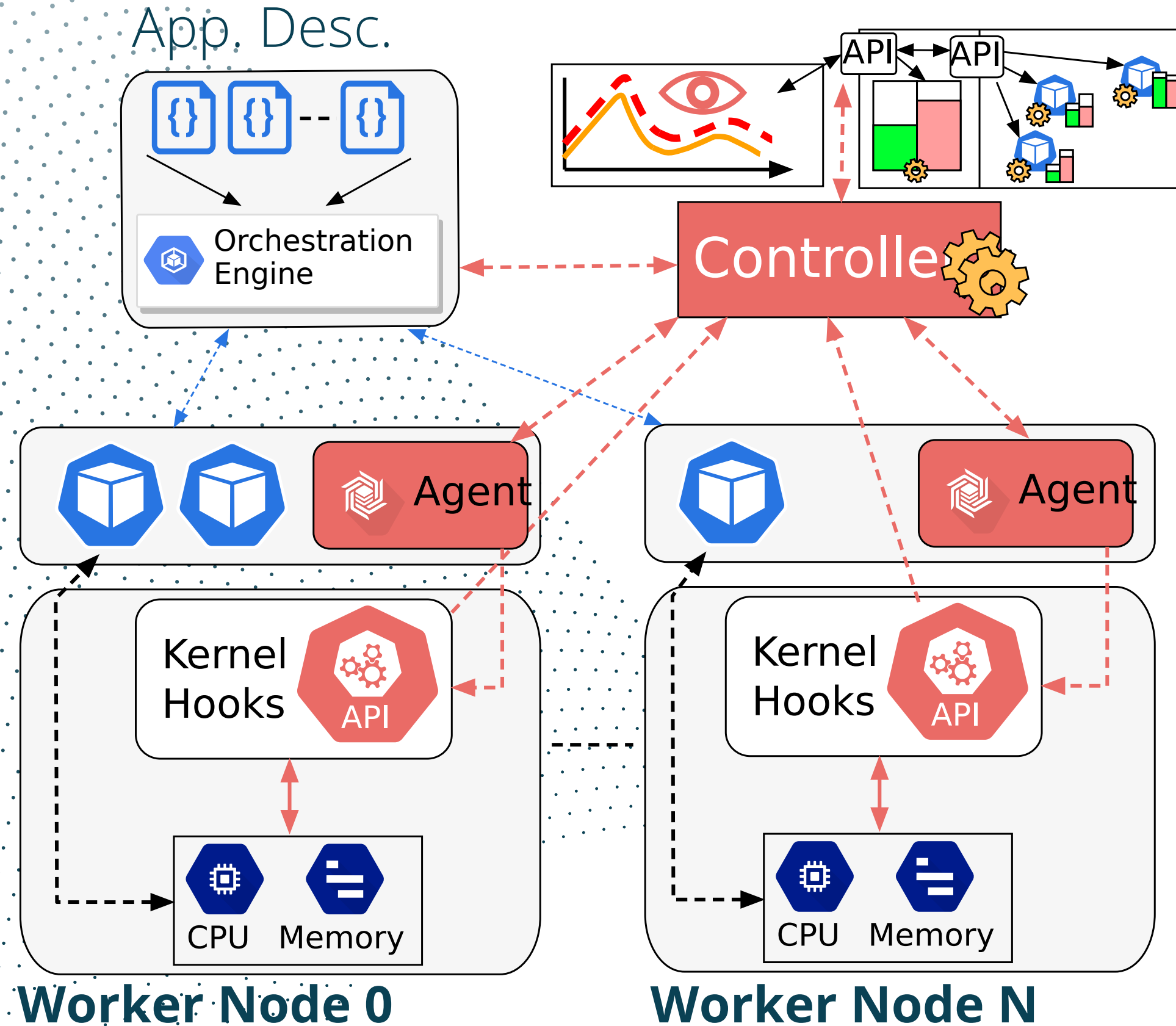
ImageProcess

Serverless

Application Request Latency



- 80% of the time, Escra + Openwhisk sees a strong performance gain
- 99th%-ile latency remains similar for both



In Conclusion...

- Fine-grained allocations allows us to react to, not predict, workload changes
- High efficiency and high performance!
- Outperforms state of the art!
- Significant resource saving in serverless deployments

Questions?



github.com/gregcusack/Esca



[@GregoryCusack](https://twitter.com/GregoryCusack)



gregory.cusack@colorado.edu



linkedin.com/in/gregorycusack

Huge Thank You to:
Maziyar Nazari,
Sepideh Goodarzy,
Erika Hunhoff,
Prerit Oberai,
Eric Keller,
Eric Rozner,
Richard Han