



# Targeted Resource Management in Multi-tenant Distributed Systems

**Jonathan Mace**

*Brown University*



**Peter Bodik**

*MSR Redmond*



**Rodrigo Fonseca**

*Brown University*



**Madanlal Musuvathi**

*MSR Redmond*



BROWN UNIVERSITY

Microsoft®  
**Research**

# *Resource Management in Multi-Tenant Systems*

# *Resource Management in Multi-Tenant Systems*



Failure in availability zone **cascades to shared control plane**, causes thread pool starvation for all zones

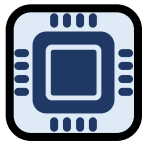
- April 2011 – Amazon EBS Failure

# *Resource Management in Multi-Tenant Systems*



Failure in availability zone **cascades to shared control plane**, causes thread pool starvation for all zones

- April 2011 – Amazon EBS Failure



Aggressive **background task** responds to increased hardware capacity with deluge of warnings and logging

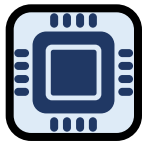
- Aug. 2012 – Azure Storage Outage

# *Resource Management in Multi-Tenant Systems*



Failure in availability zone **cascades to shared control plane**, causes thread pool starvation for all zones

- April 2011 – Amazon EBS Failure



Aggressive **background task** responds to increased hardware capacity with deluge of warnings and logging

- Aug. 2012 – Azure Storage Outage



Code change increases database usage, shifts bottleneck to unmanaged **application-level lock**

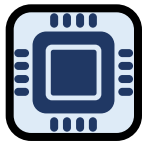
- Nov. 2014 – Visual Studio Online outage

# Resource Management in Multi-Tenant Systems



Failure in availability zone **cascades to shared control plane**, causes thread pool starvation for all zones

- April 2011 – Amazon EBS Failure



Aggressive **background task** responds to increased hardware capacity with deluge of warnings and logging

- Aug. 2012 – Azure Storage Outage



Code change increases database usage, shifts bottleneck to unmanaged **application-level lock**

- Nov. 2014 – Visual Studio Online outage



**Shared storage layer** bottlenecks circumvent resource management layer

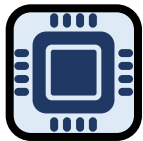
- 2014 – Communication with Cloudera

# *Resource Management in Multi-Tenant Systems*



Failure in availability zone **cascades to shared control plane**, causes thread pool starvation for all zones

- April 2011 – Amazon EBS Failure



Aggressive **background task** responds to increased hardware capacity with deluge of warnings and logging

- Aug. 2012 – Azure Storage Outage



Code change increases database usage, shifts bottleneck to unmanaged **application-level lock**

- Nov. 2014 – Visual Studio Online outage

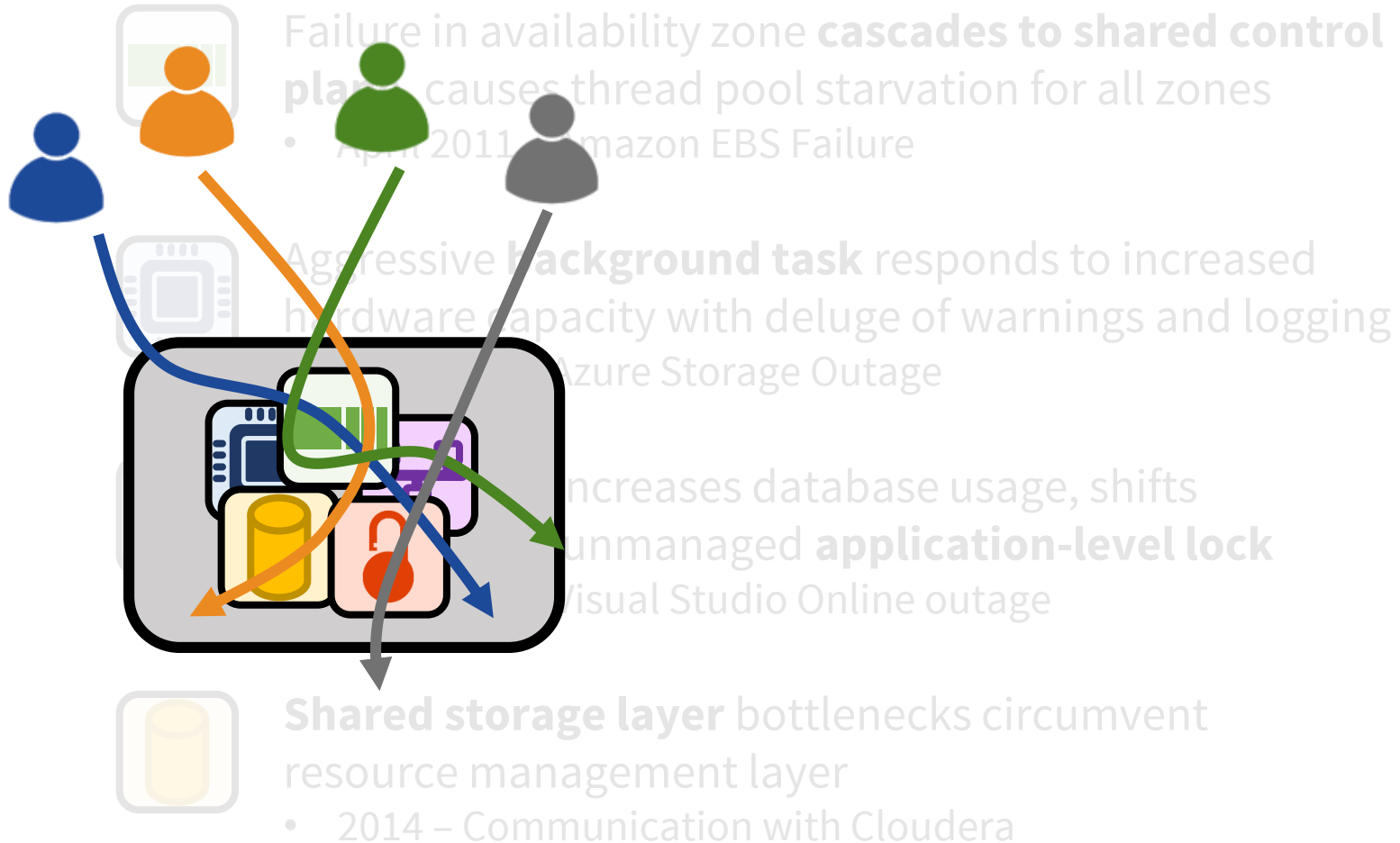


**Shared storage layer** bottlenecks circumvent resource management layer

- 2014 – Communication with Cloudera

Degraded performance, Violated SLOs, system outages

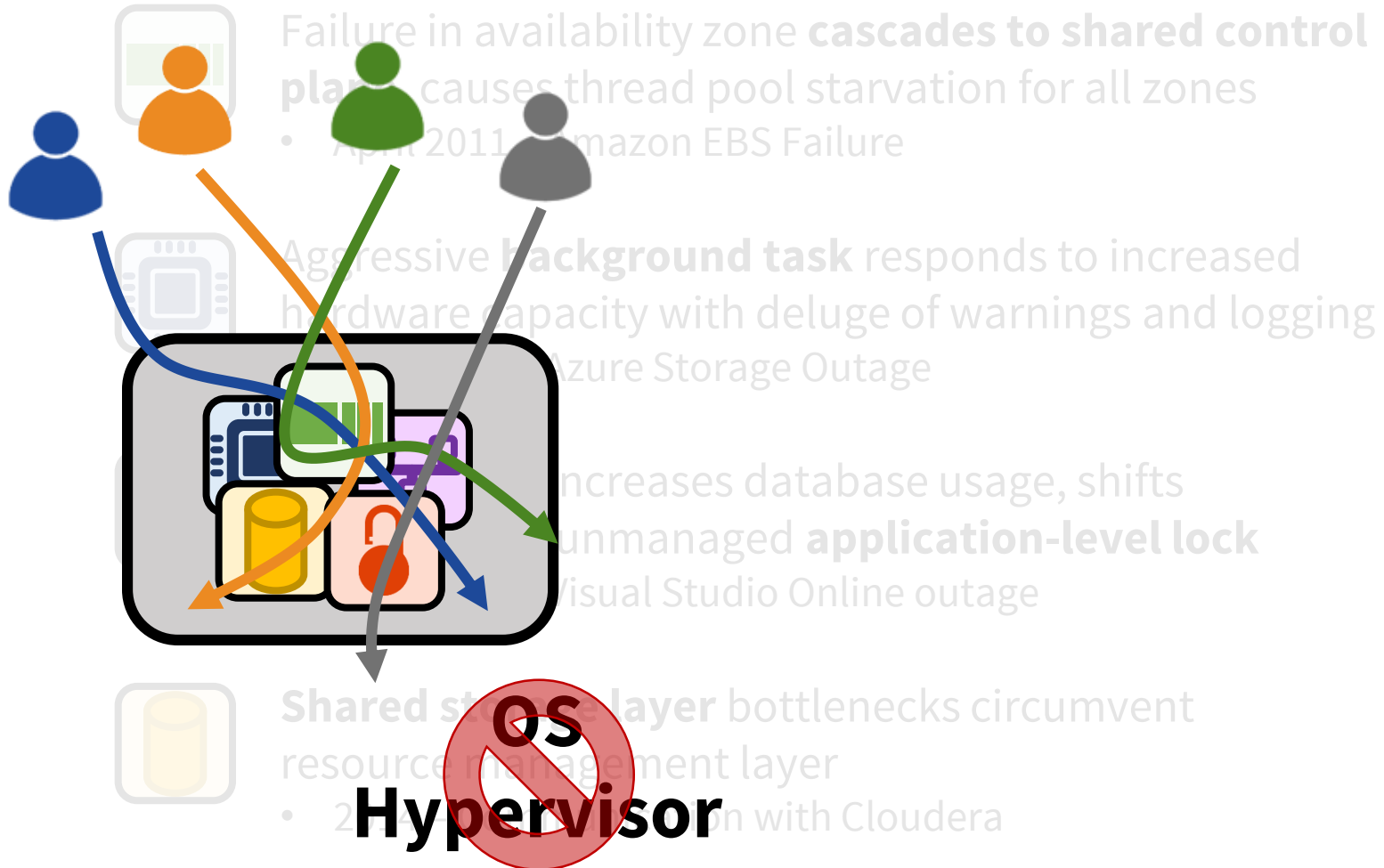
# Resource Management in Multi-Tenant Systems



Degraded performance, Violated SLOs, system outages



# Resource Management in Multi-Tenant Systems



Degraded performance, Violated SLOs, system outages

Failure in availability zone cascades to shared control plane causes thread pool starvation for all zones

- April 2011 Amazon EBS Failure

Aggressive background task increases database hardware capacity with deluged Azure Storage Outage

Increases database unmanaged application Visual Studio Online outage

Shared storage layer bottlenecks circuit resource management layer

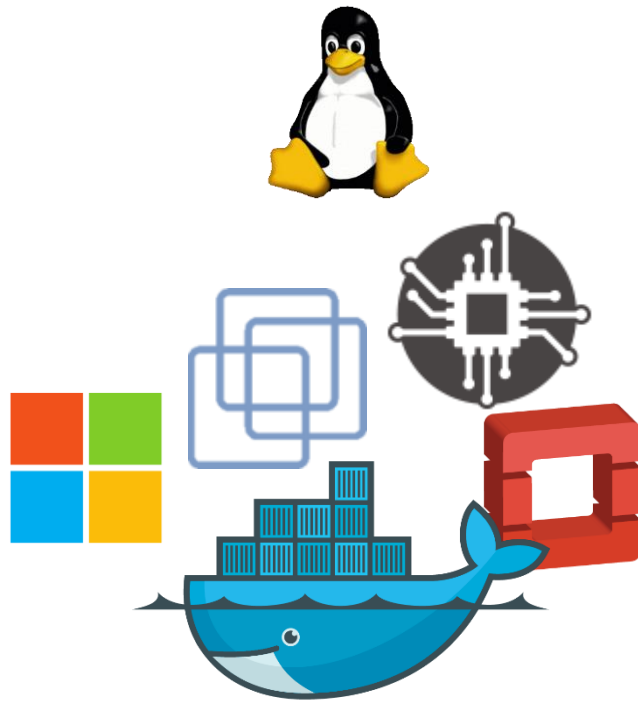
- 2014 outage with Cloudera

**OS layer**

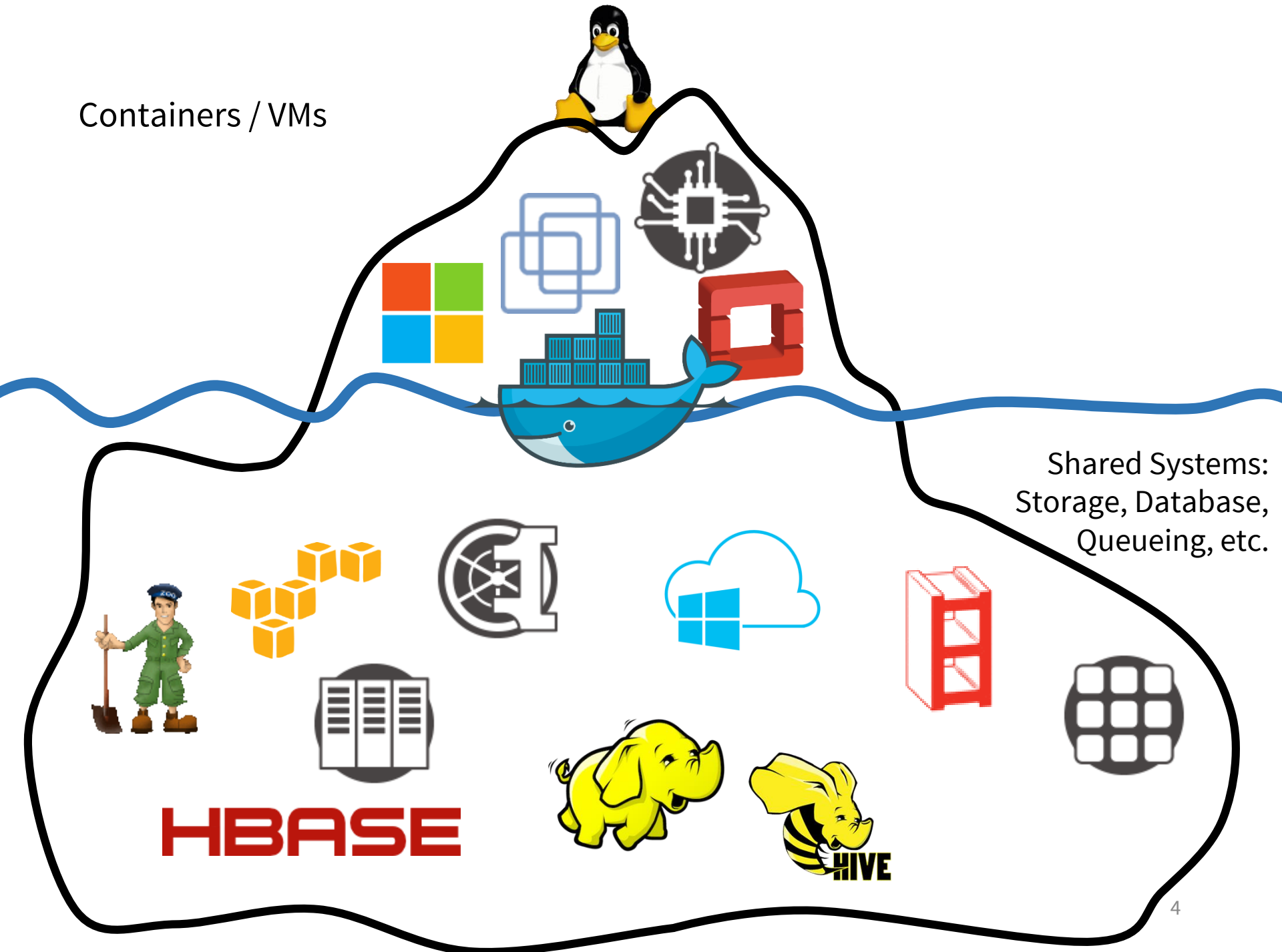
**Hypervisor**

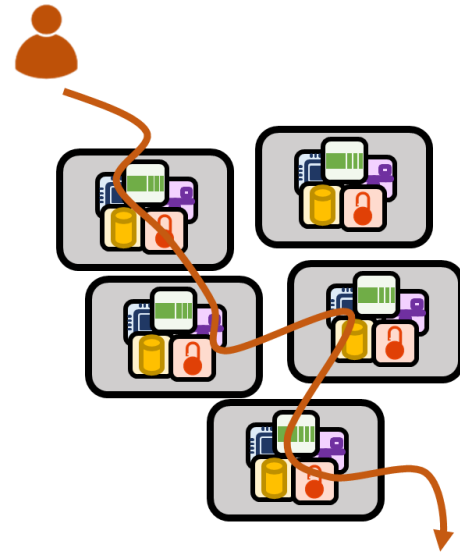
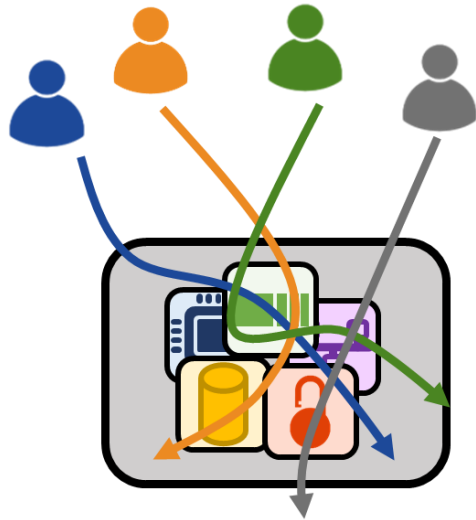
3

## Containers / VMs



Containers / VMs









Monitors resource usage of each tenant in near real-time  
Actively schedules tenants and activities



Monitors resource usage of each tenant in near real-time  
Actively schedules tenants and activities

High-level, centralized policies:  
Encapsulates resource management logic





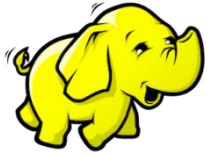
Monitors resource usage of each tenant in near real-time  
Actively schedules tenants and activities

High-level, centralized policies:

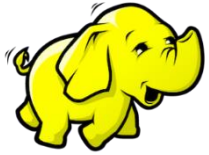
- Encapsulates resource management logic

- Abstractions – not specific to resource type, system

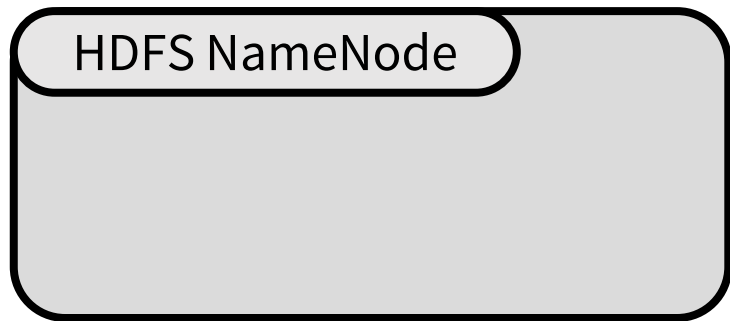
- Achieve different goals: guarantee average latencies, fair share a resource, etc.



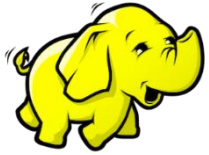
# Hadoop Distributed File System (HDFS)



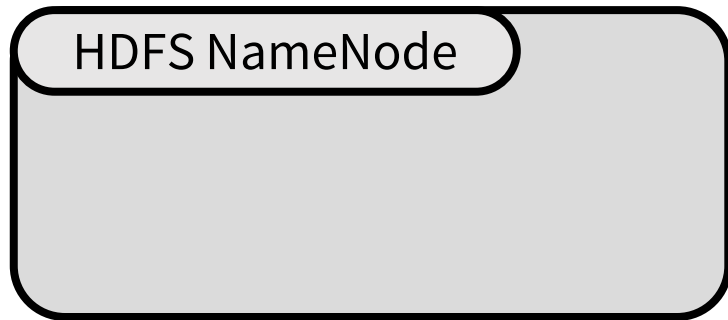
# Hadoop Distributed File System (HDFS)



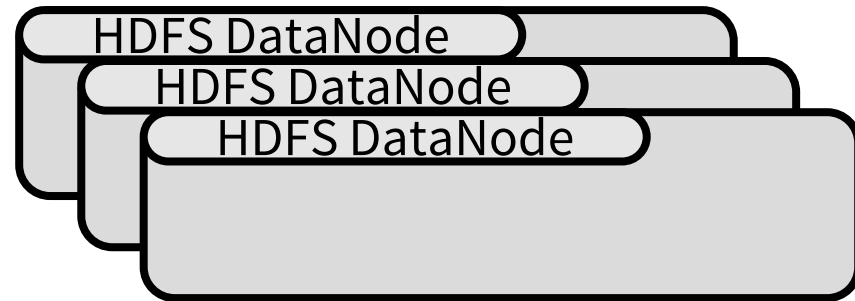
Filesystem metadata



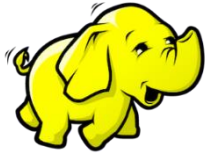
# Hadoop Distributed File System (HDFS)



Filesystem metadata

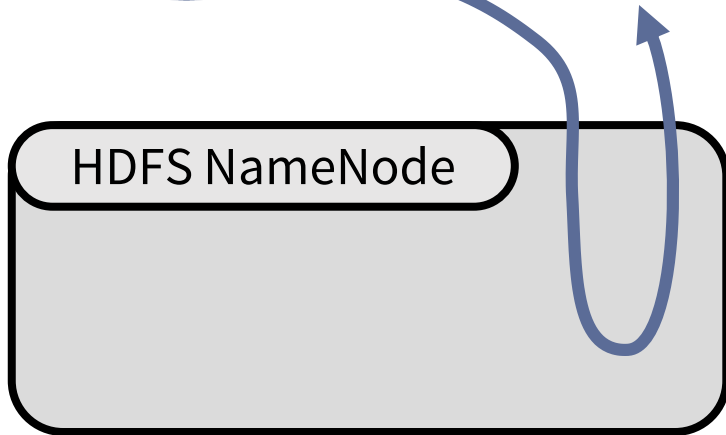


Replicated block storage



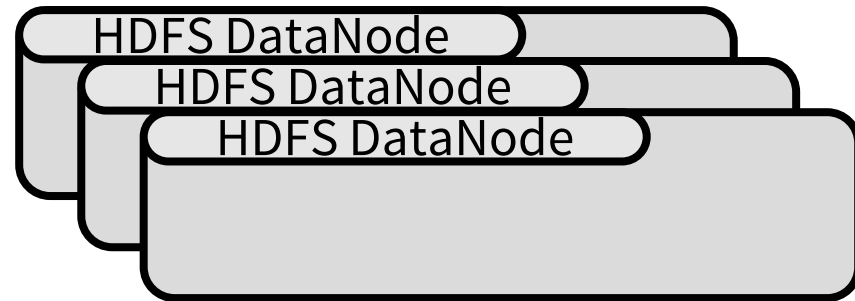
# Hadoop Distributed File System (HDFS)

Rename

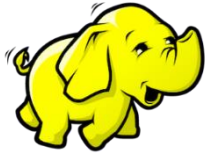


HDFS NameNode

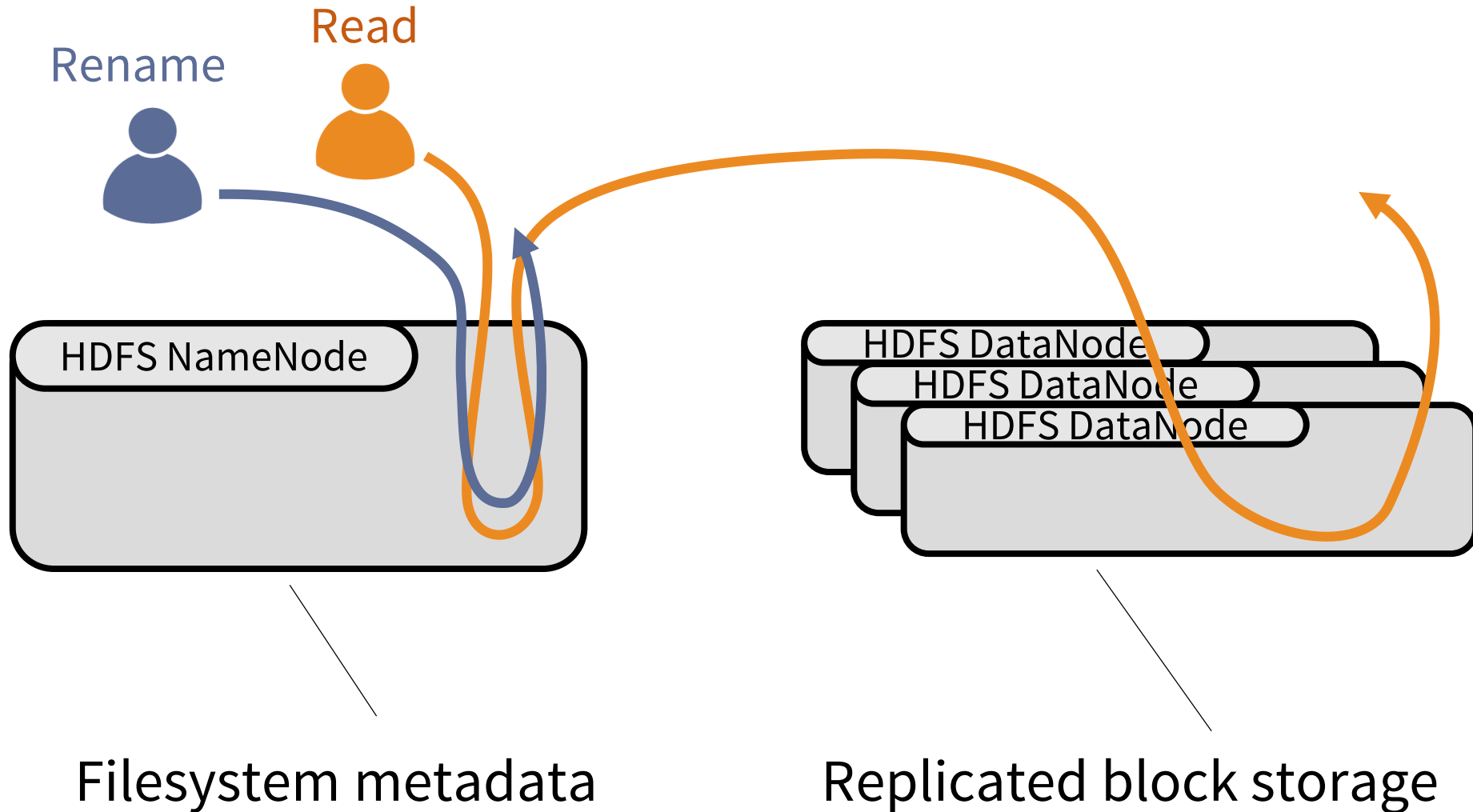
Filesystem metadata

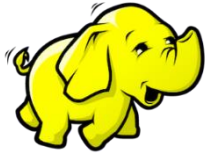


Replicated block storage

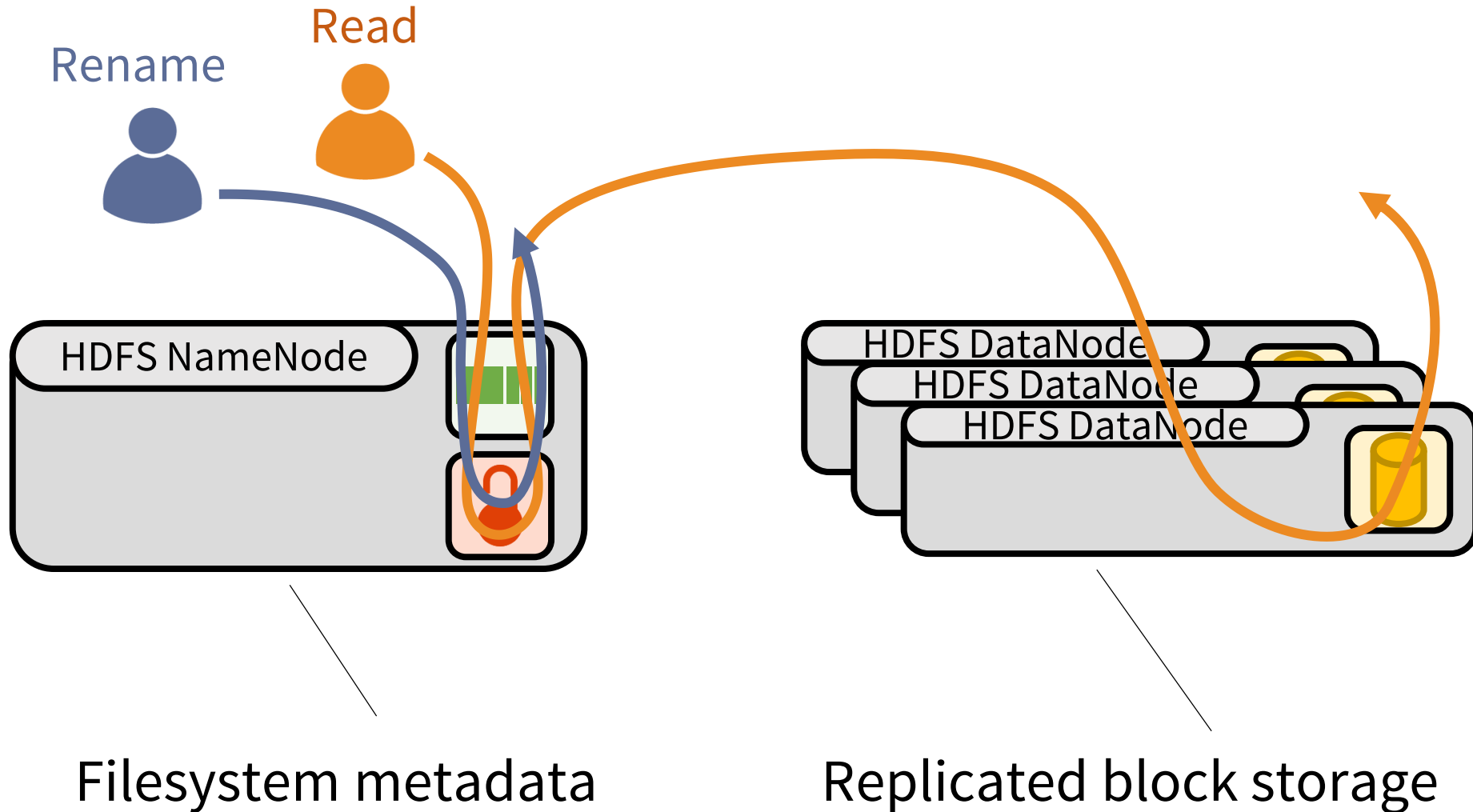


# Hadoop Distributed File System (HDFS)

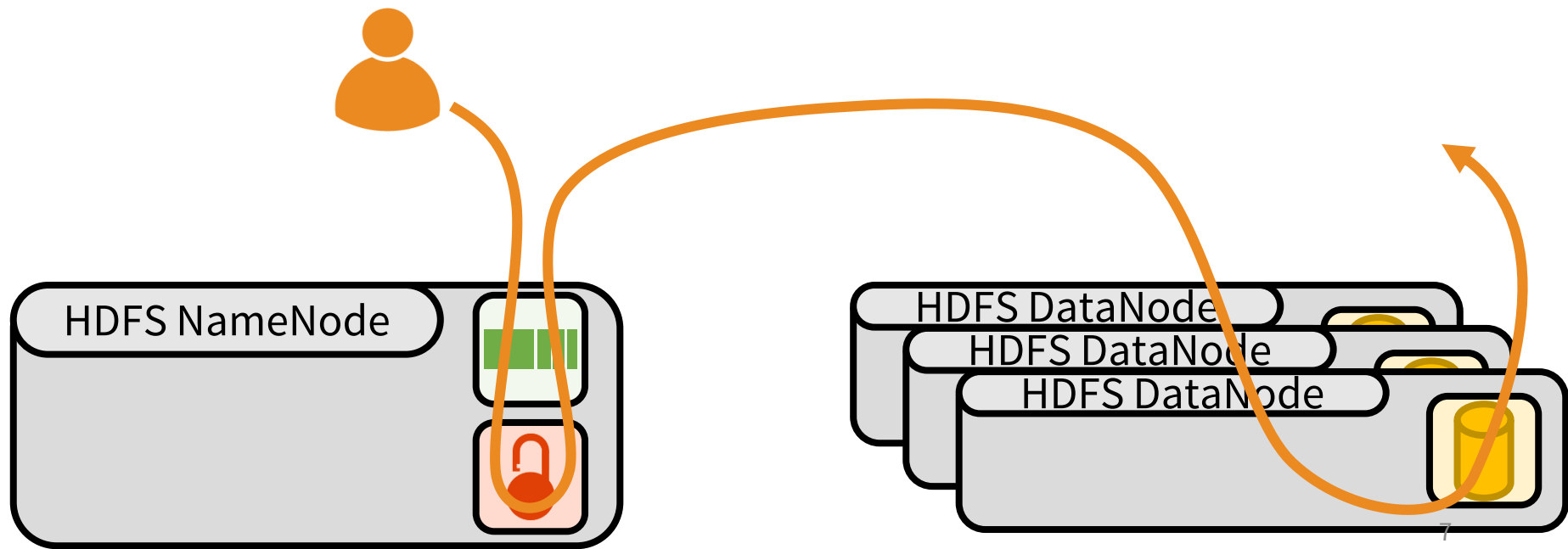




# Hadoop Distributed File System (HDFS)



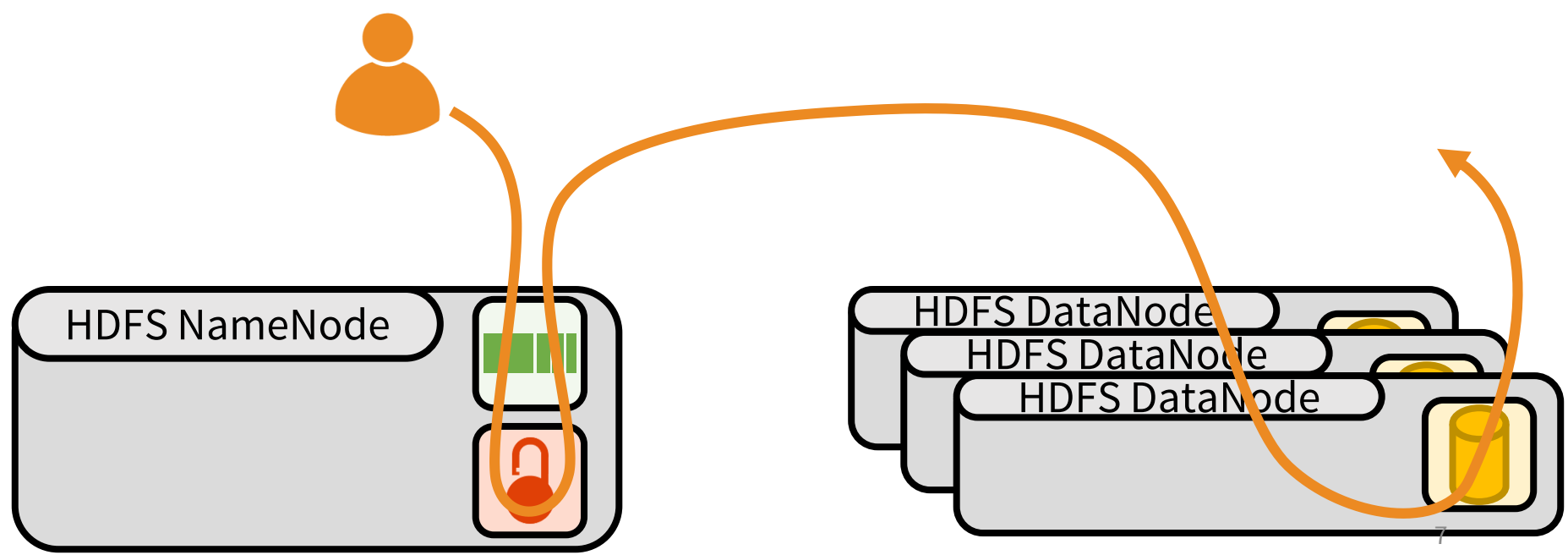
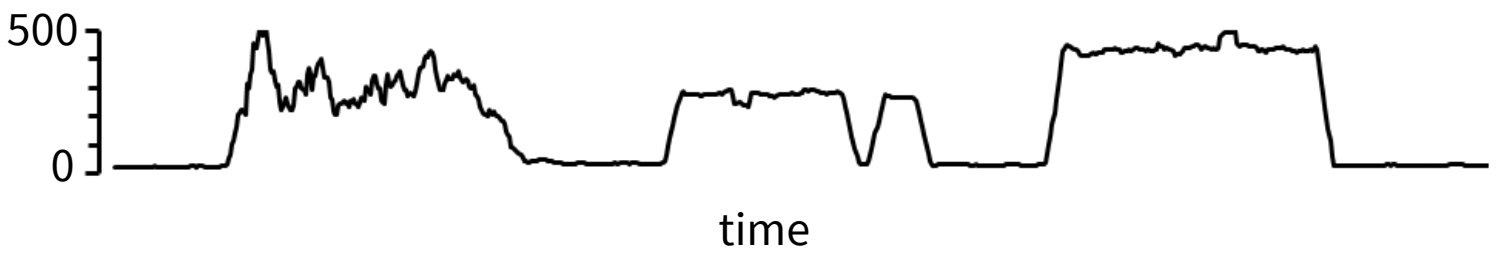
random 8kb reads





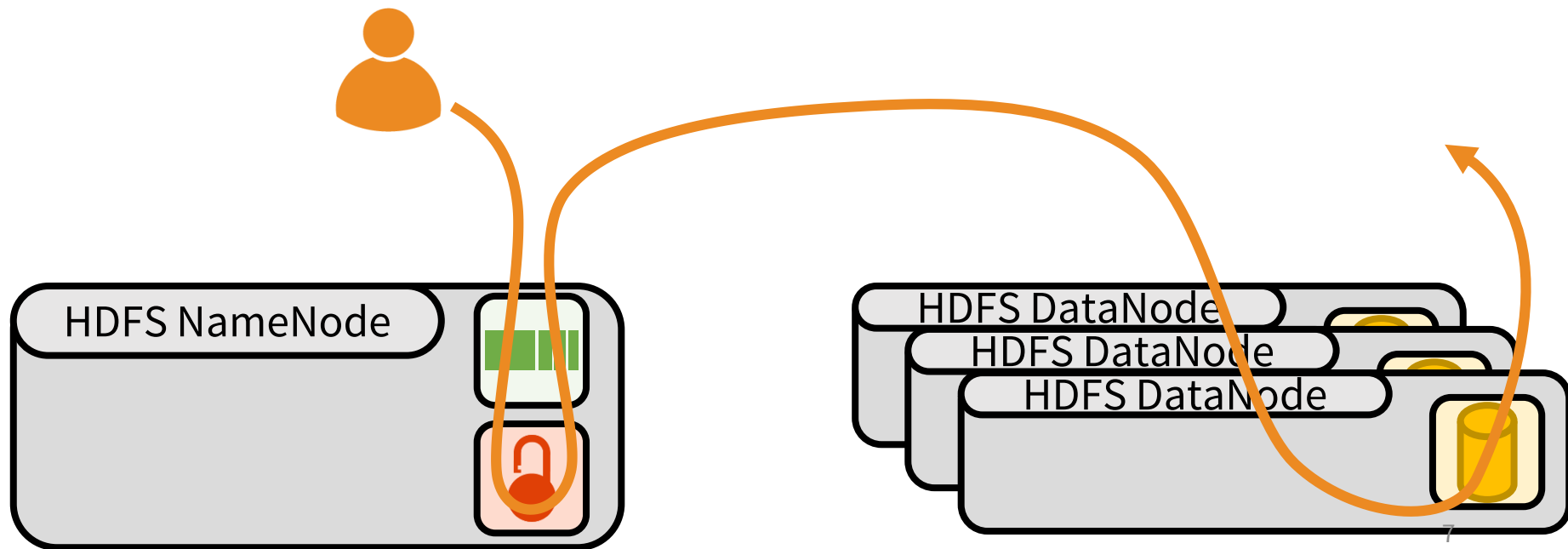
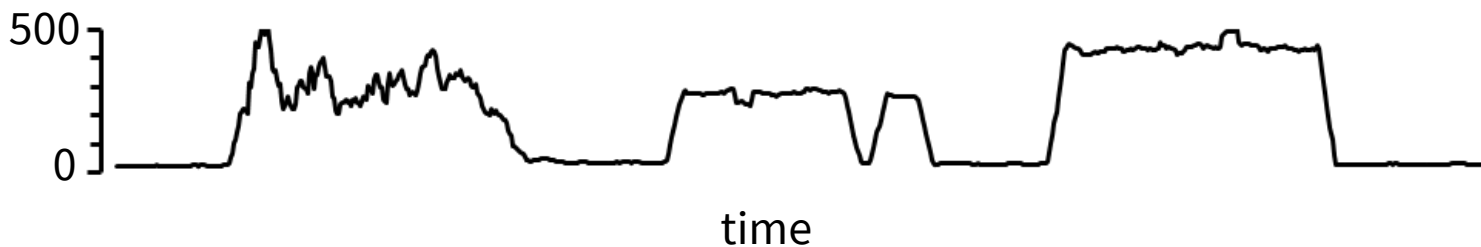
random 8kb reads

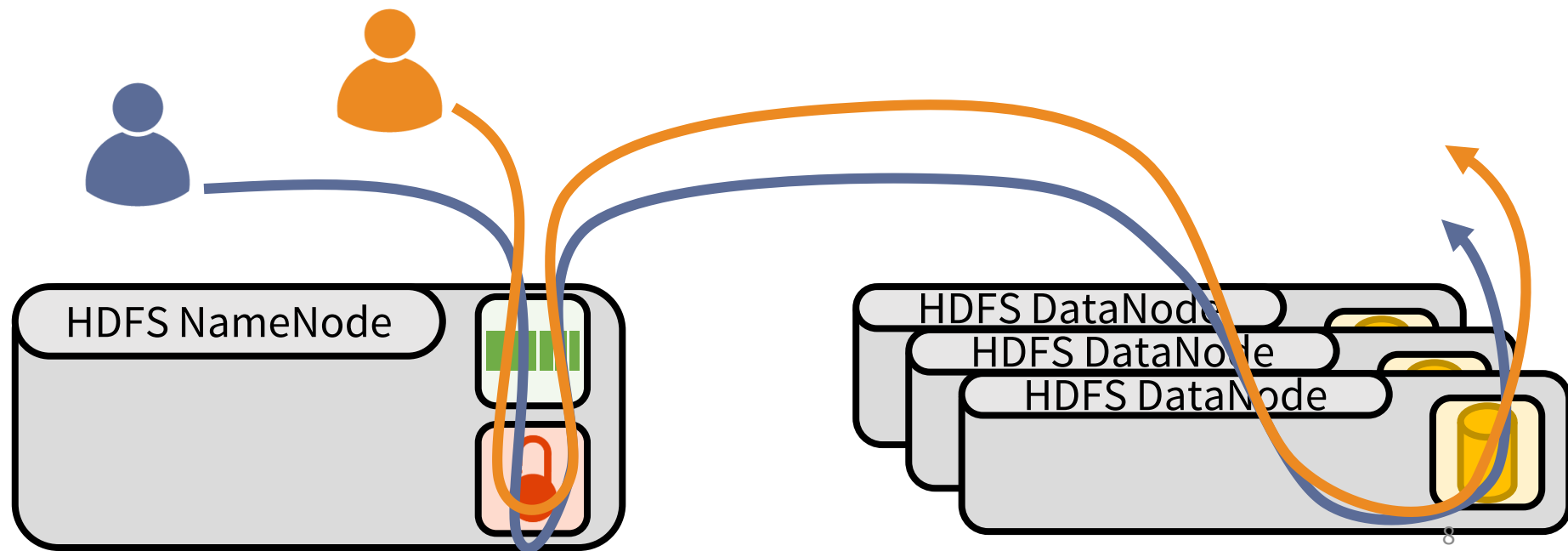
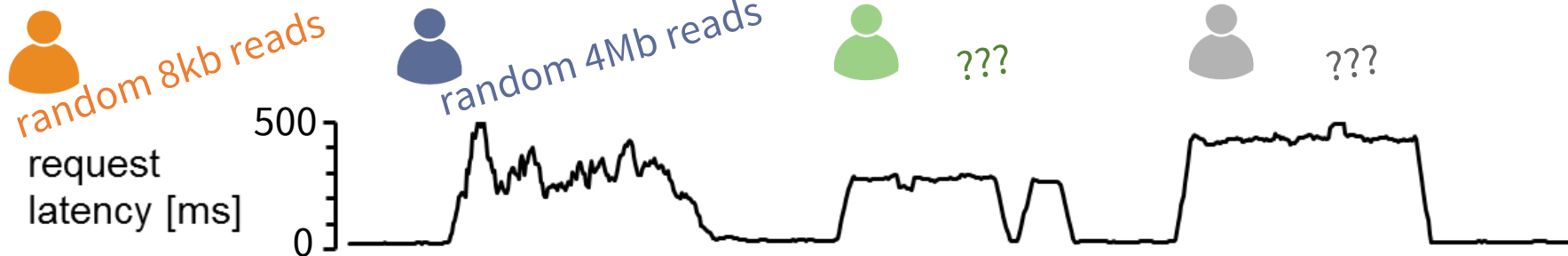
request  
latency [ms]

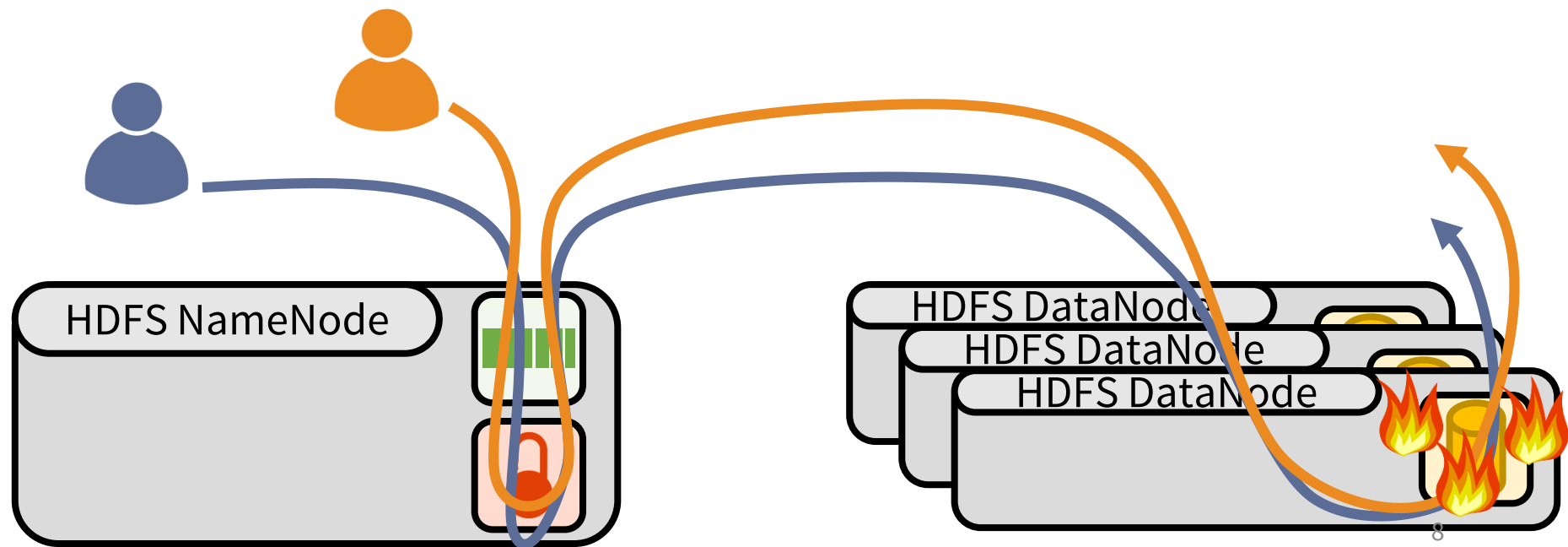


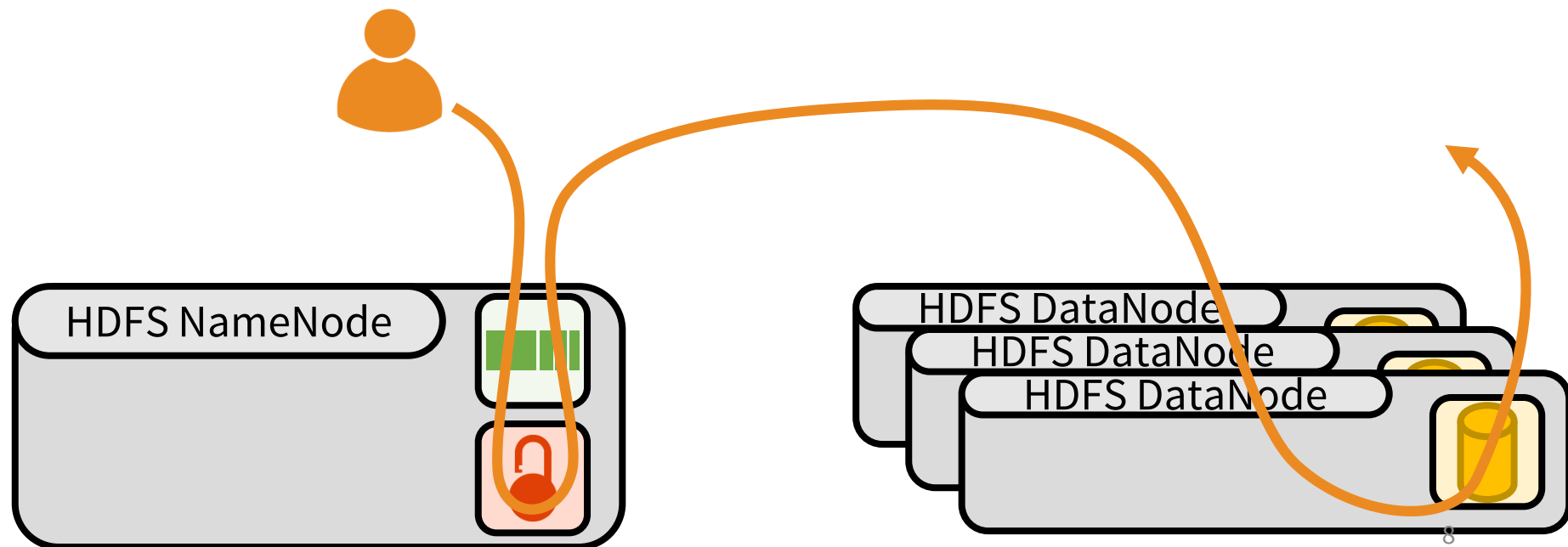
random 8kb reads

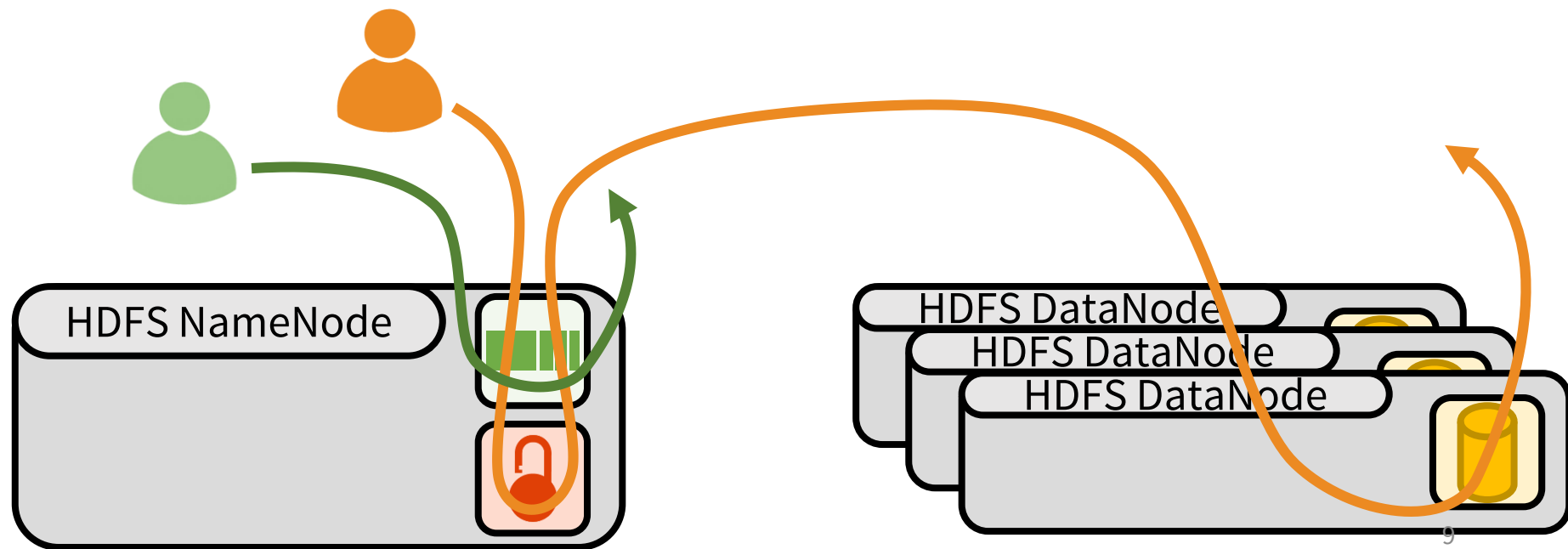
request  
latency [ms]

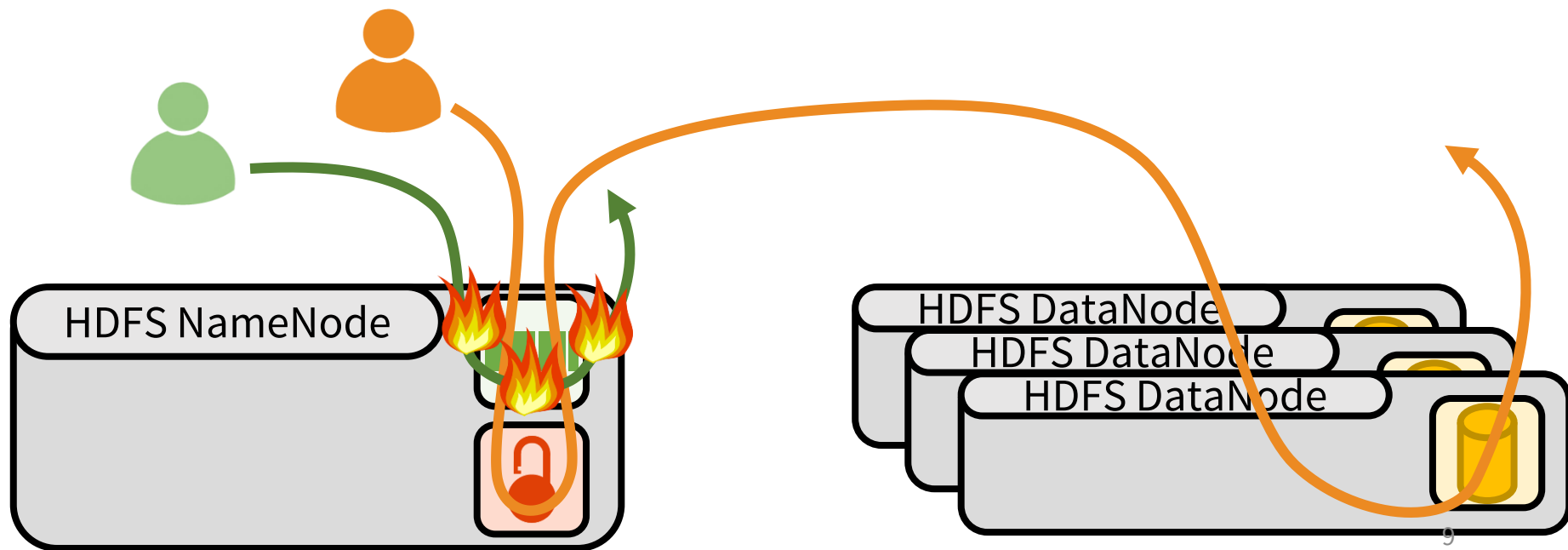


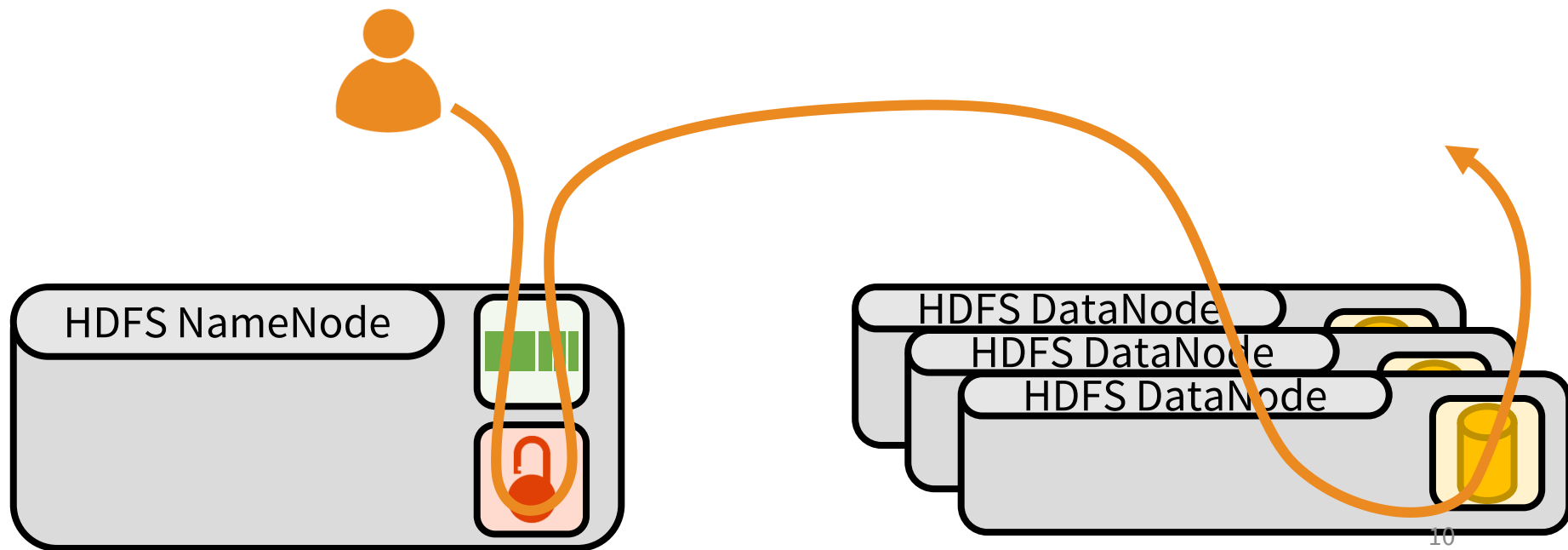
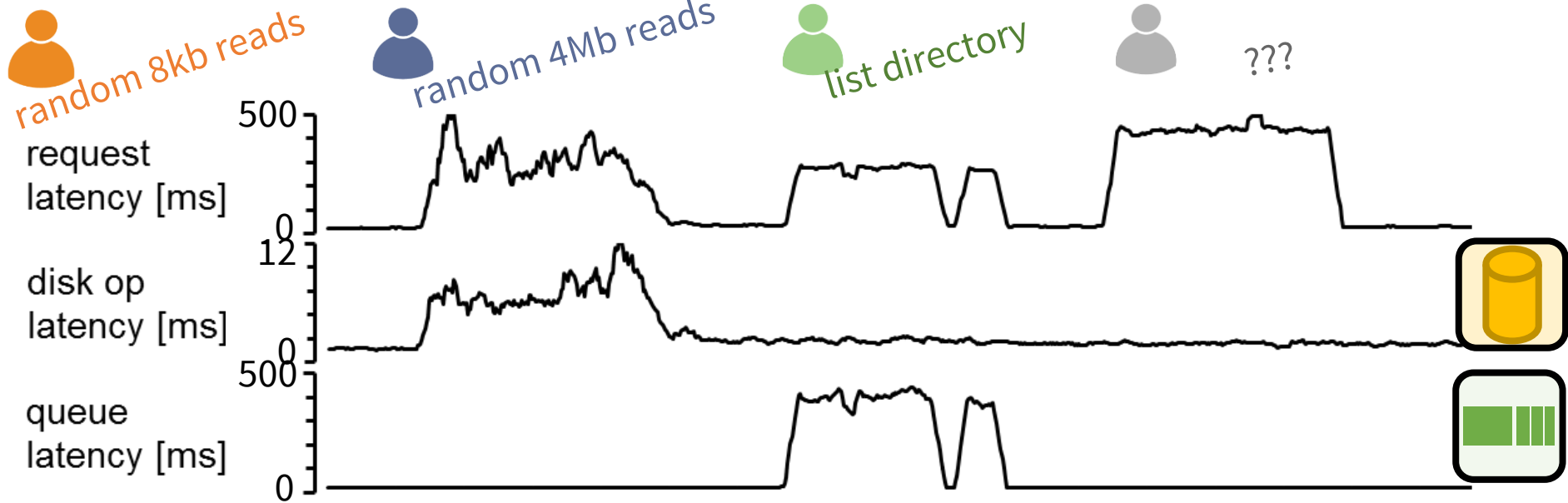




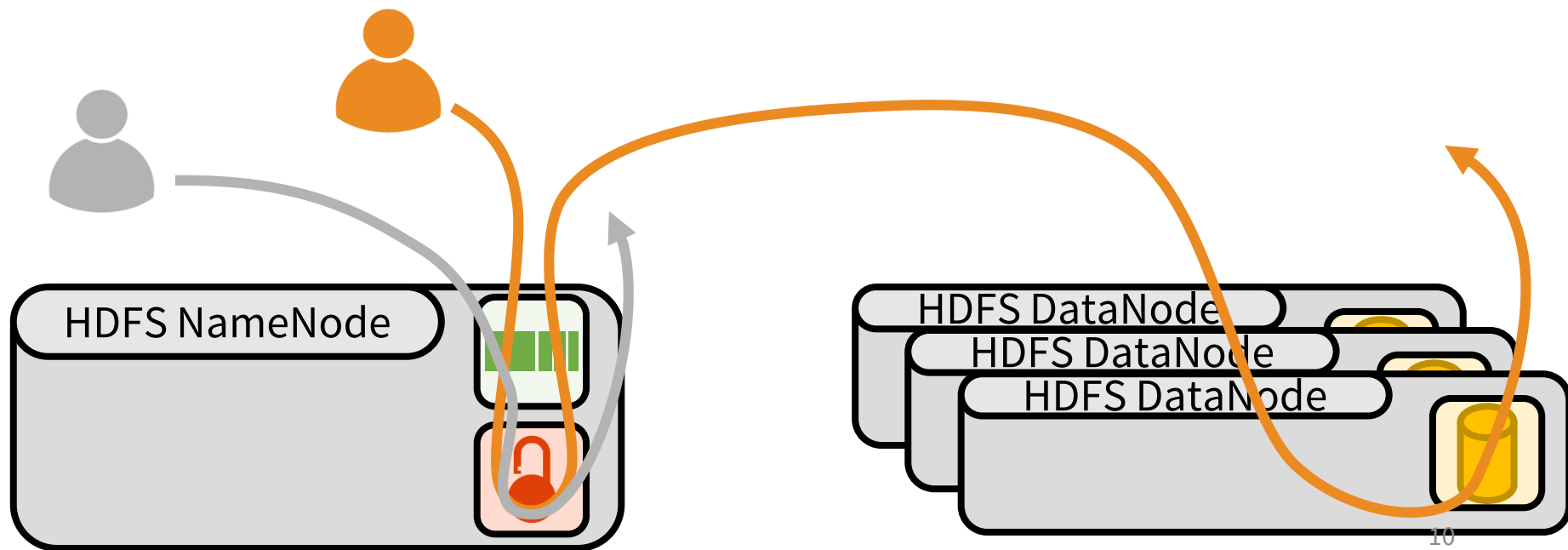
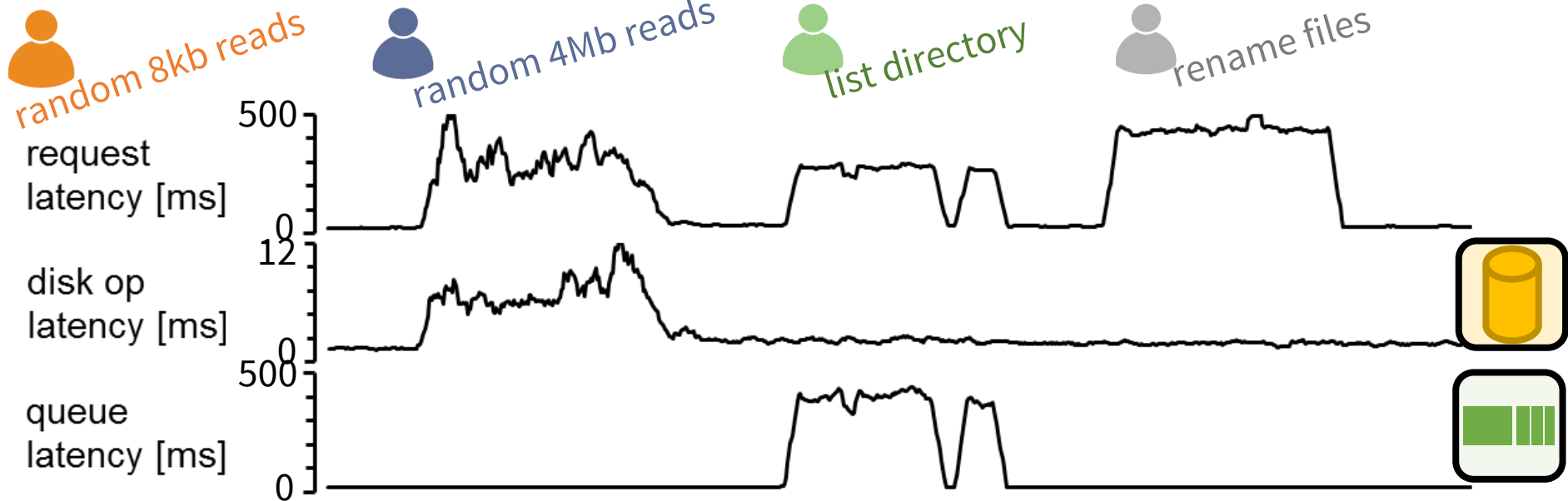


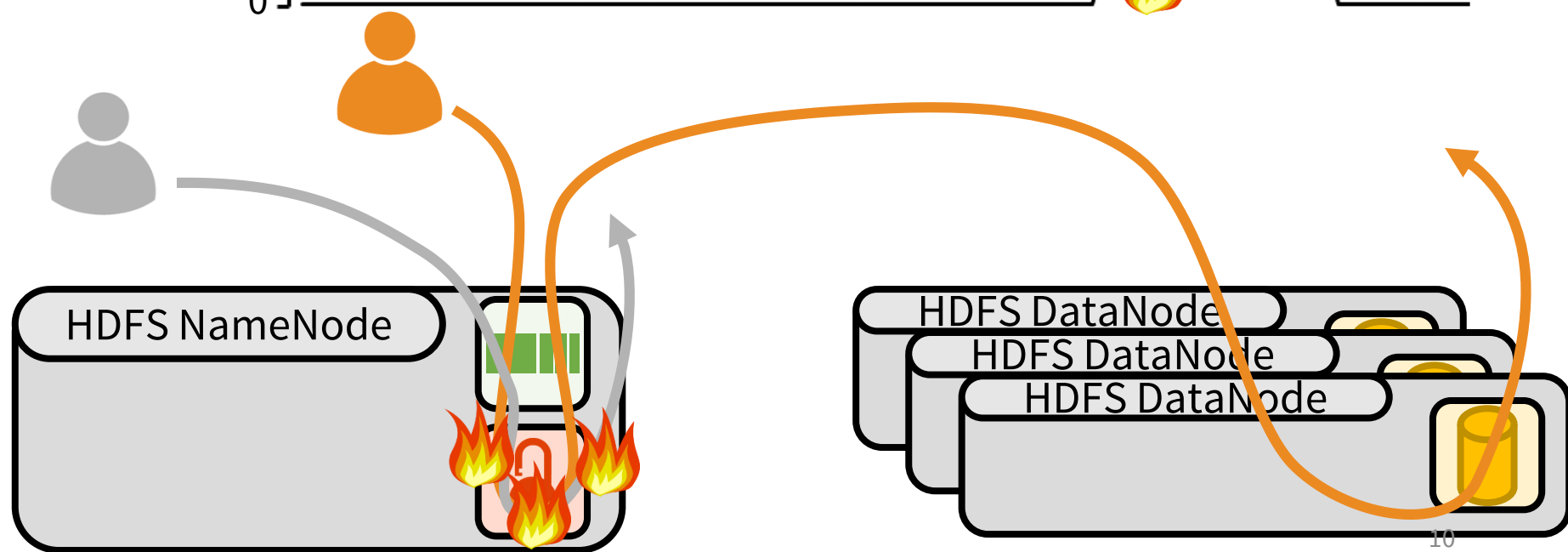
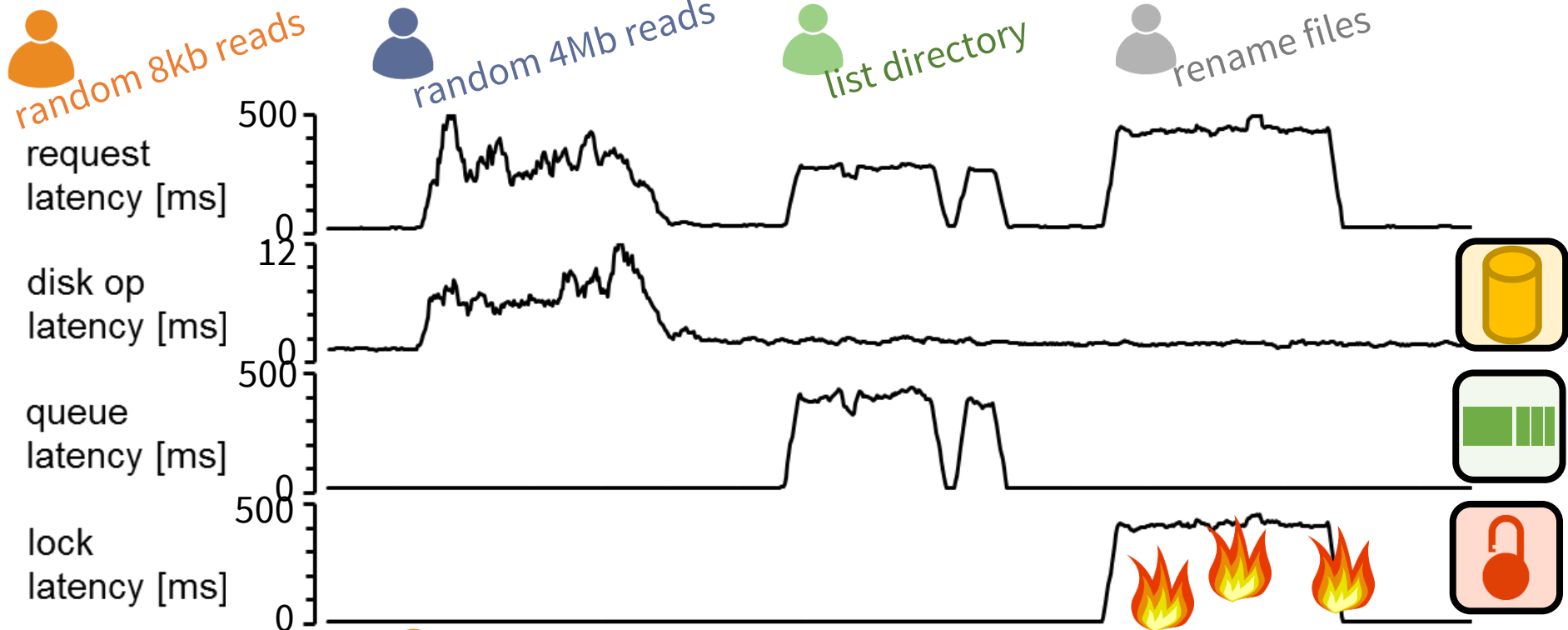


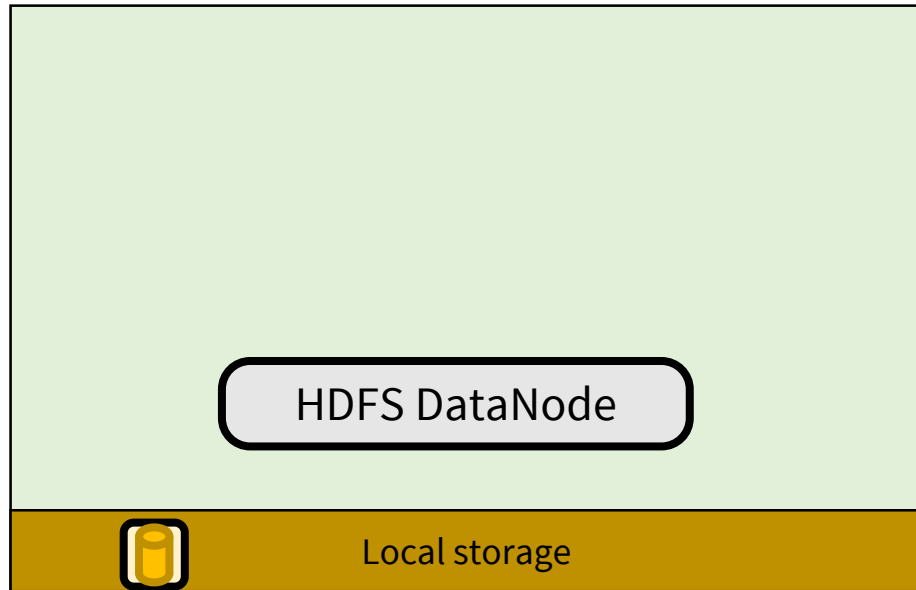


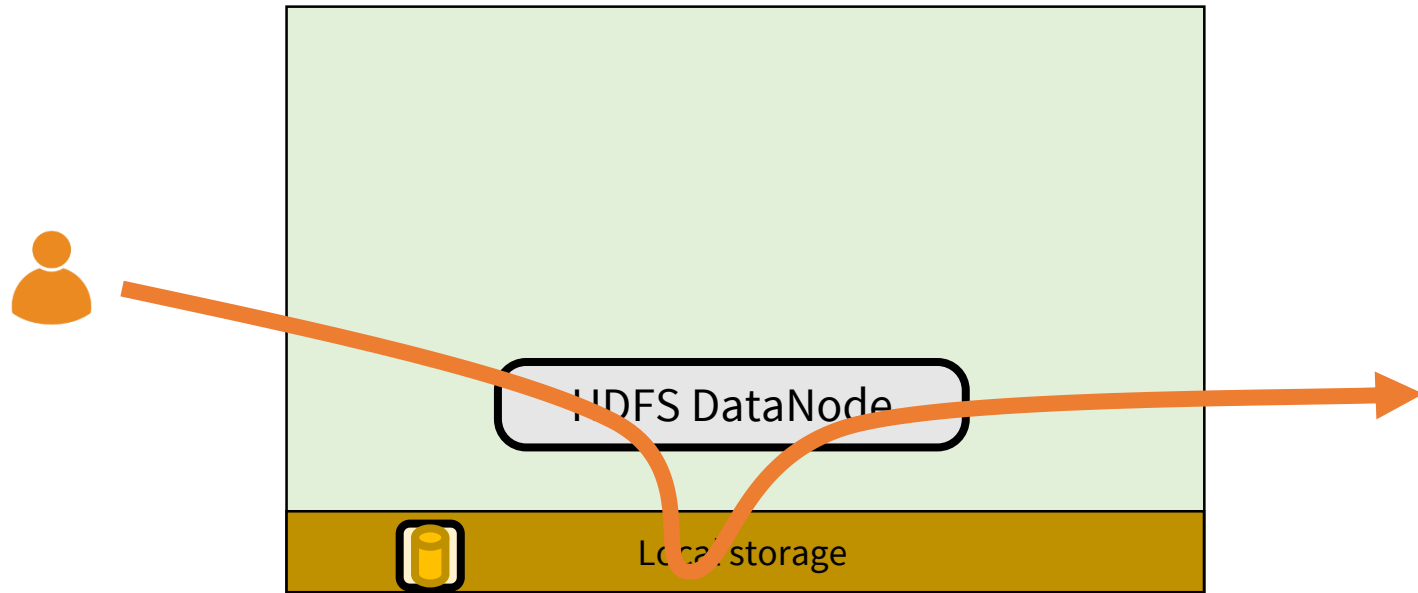


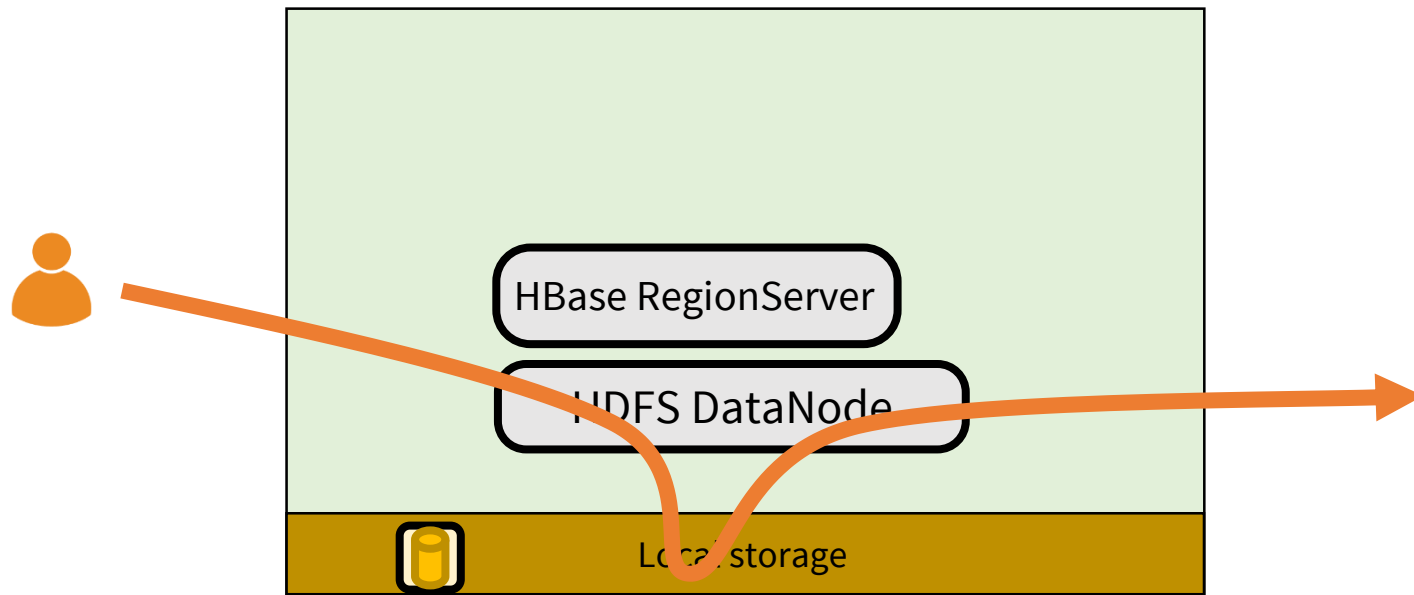


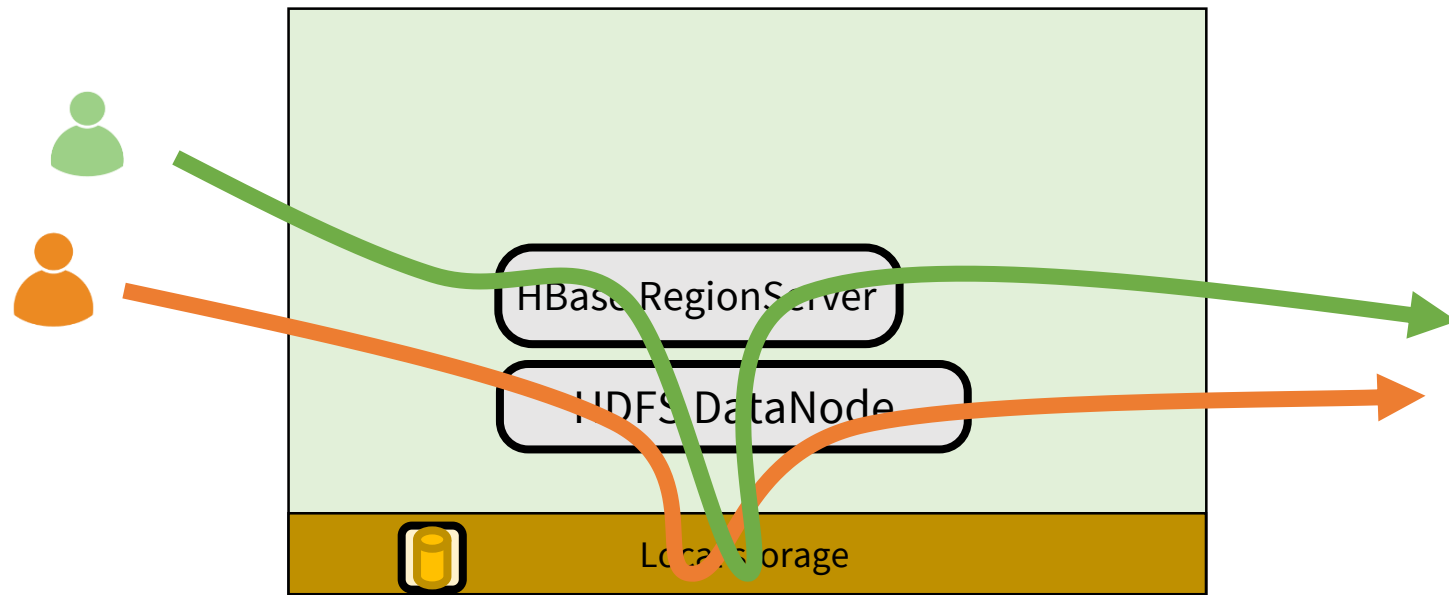


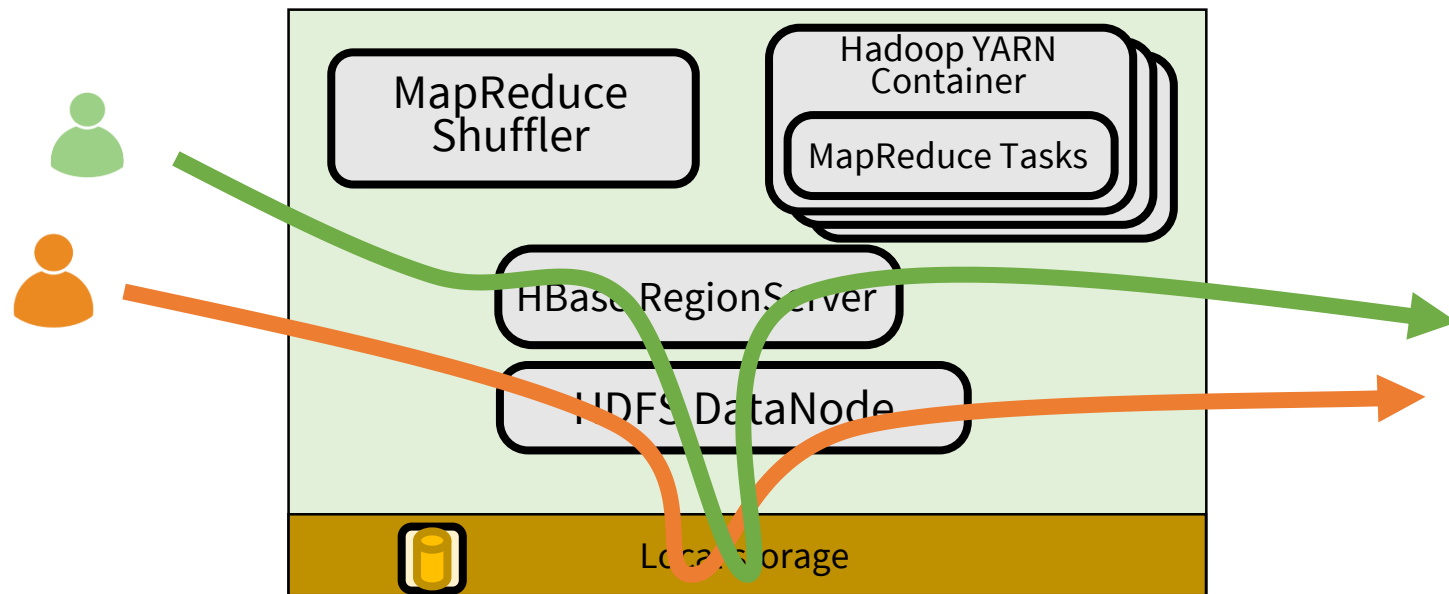


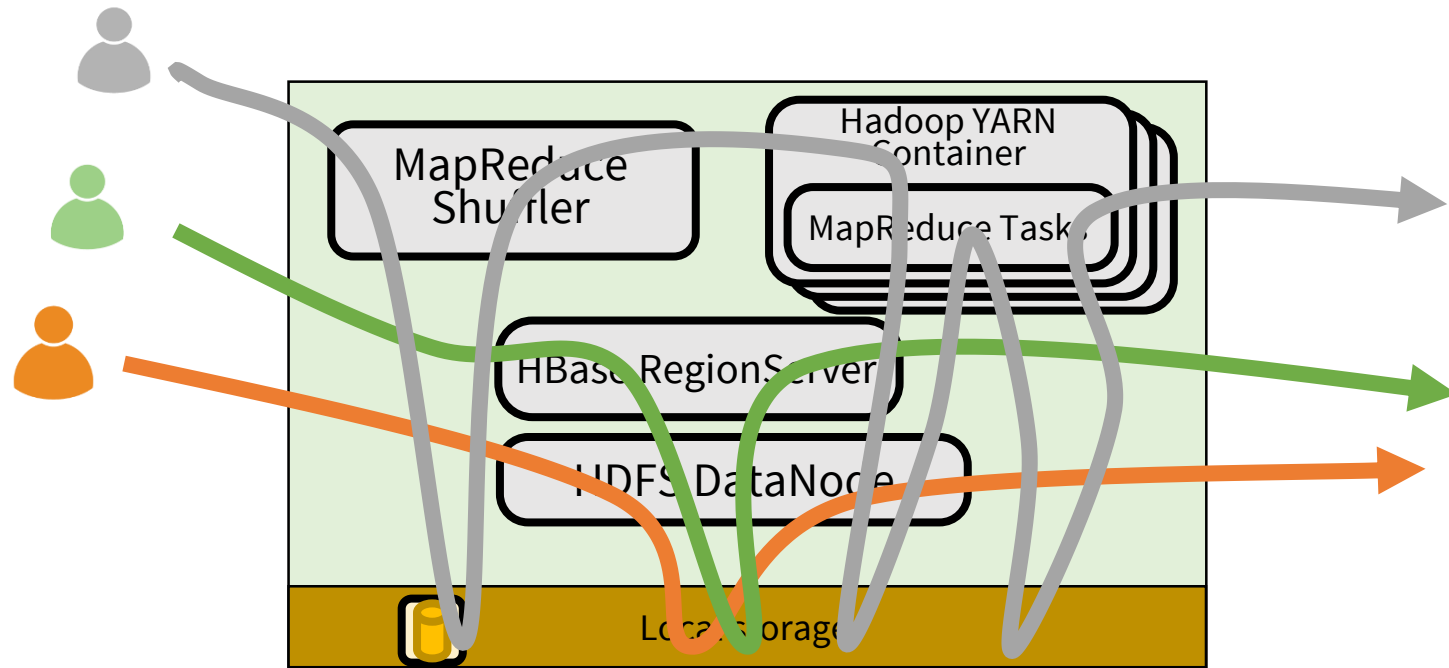




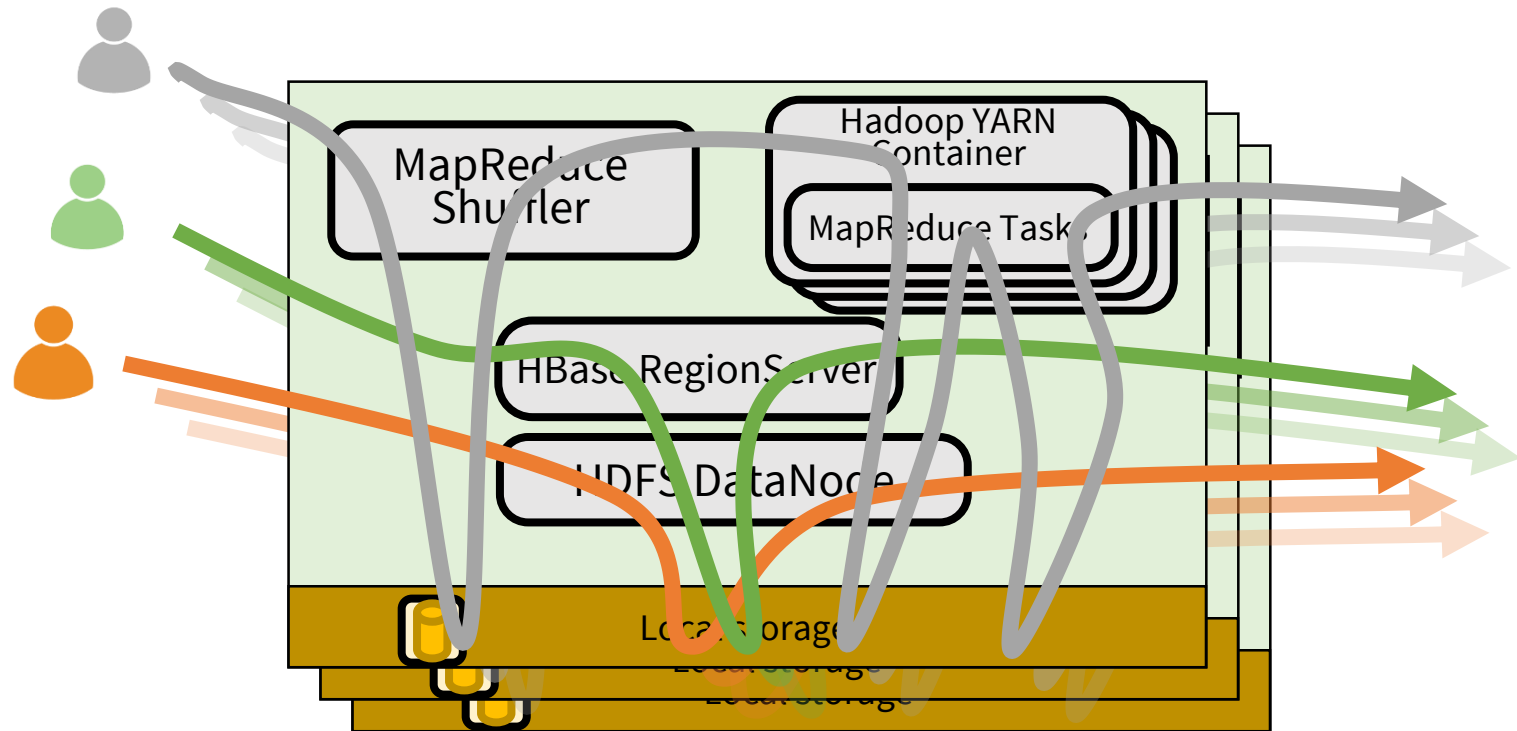












# *Retro* **Goals**

# *Retro* **Goals**

Co-ordinated control across processes, machines, and services

# *Retro* **Goals**

Co-ordinated control across processes, machines, and services

Handle system and application level resources

# *Retro* **Goals**

Co-ordinated control across processes, machines, and services

Handle system and application level resources

Principals: tenants, background tasks

# *Retro* **Goals**

Co-ordinated control across processes, machines, and services

Handle system and application level resources

Principals: tenants, background tasks

Real-time and reactive

# *Retro* **Goals**

Co-ordinated control across processes, machines, and services

Handle system and application level resources

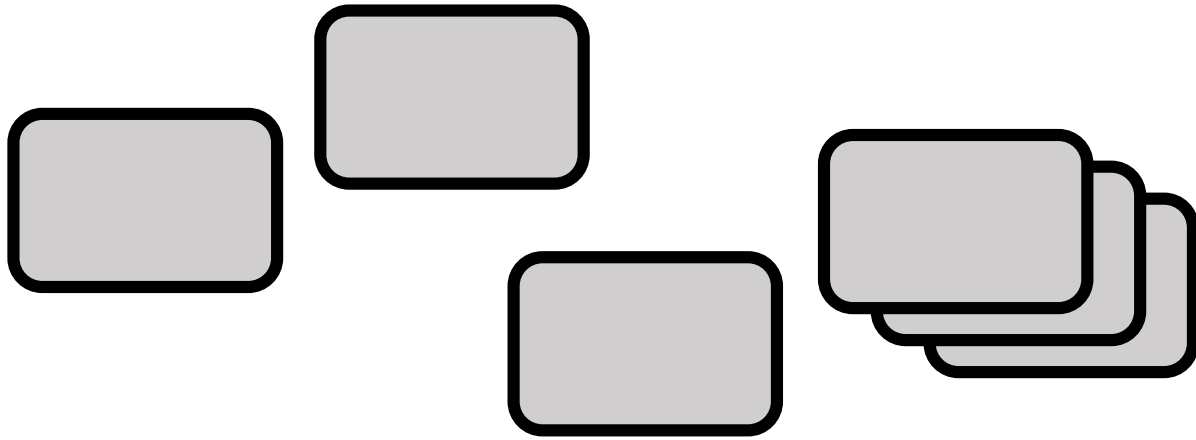
Principals: tenants, background tasks

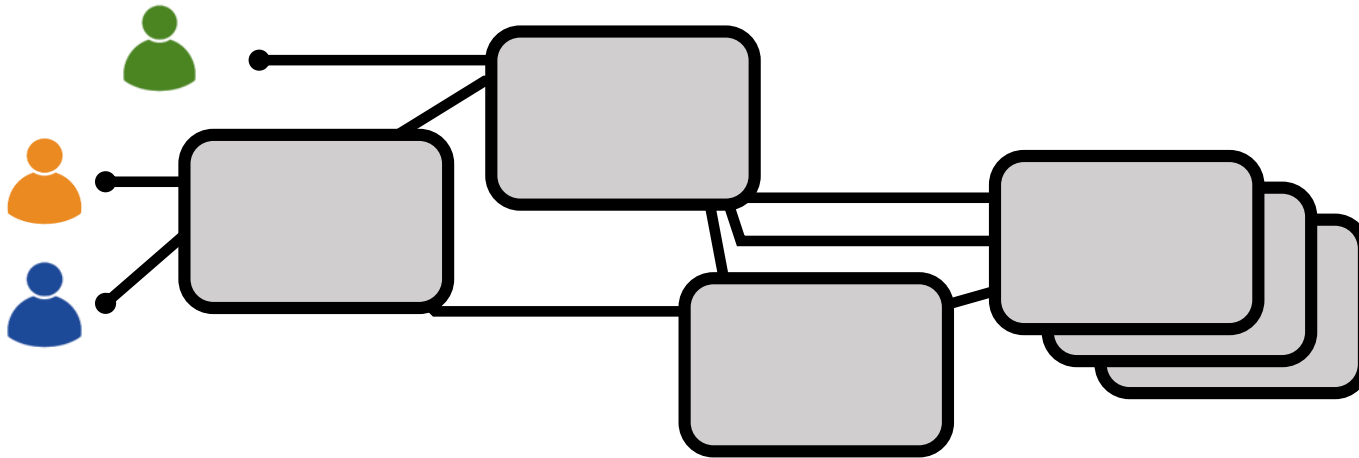
Real-time and reactive

Efficient: Only control what is needed

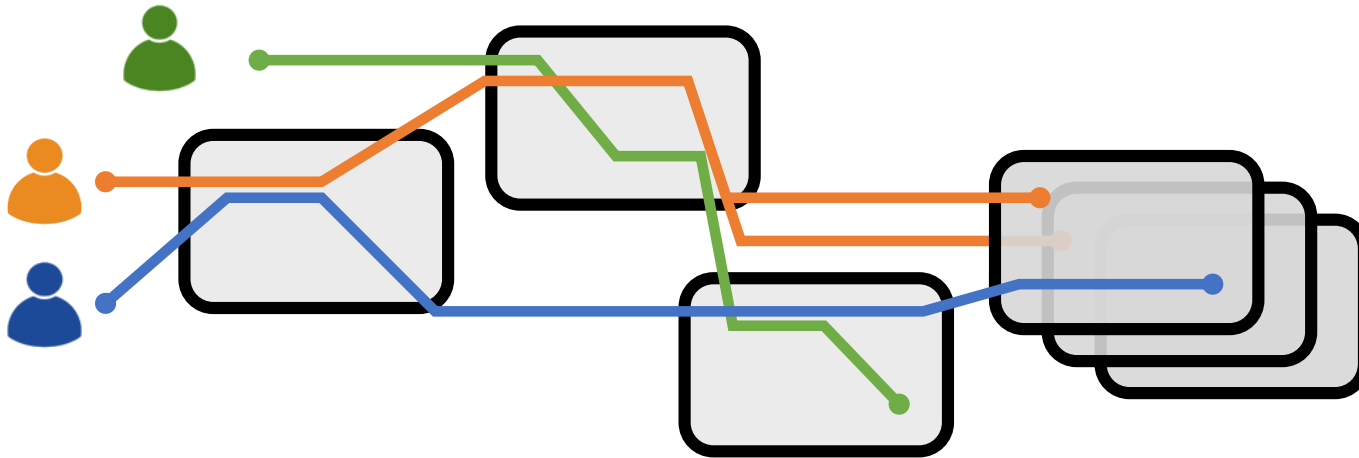
# *Retro* **Architecture**





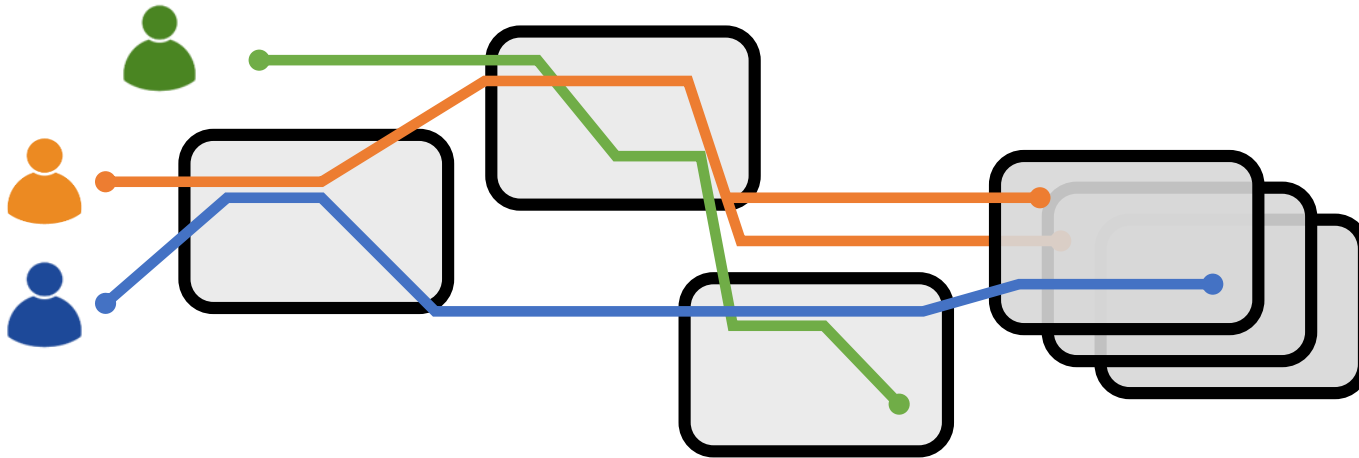


 Tenant Requests



 Workflows

Purpose: identify requests from different users, background activities  
eg, all requests from a tenant over time  
eg, data balancing in HDFS



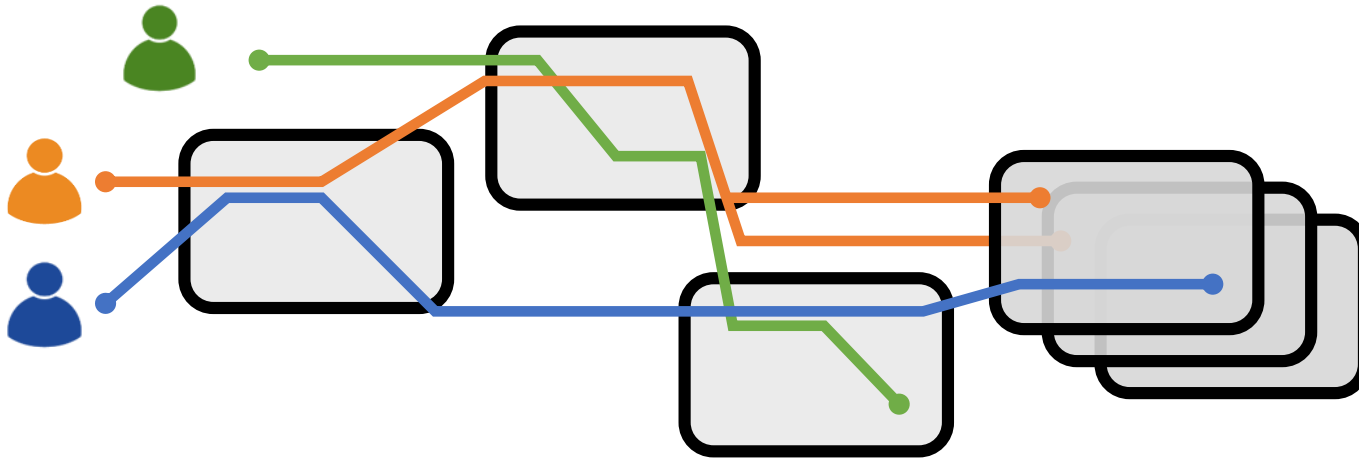
#### Workflows

Purpose: identify requests from different users, background activities

eg, all requests from a tenant over time

eg, data balancing in HDFS

Unit of resource measurement, attribution, and enforcement



### Workflows

Purpose: identify requests from different users, background activities

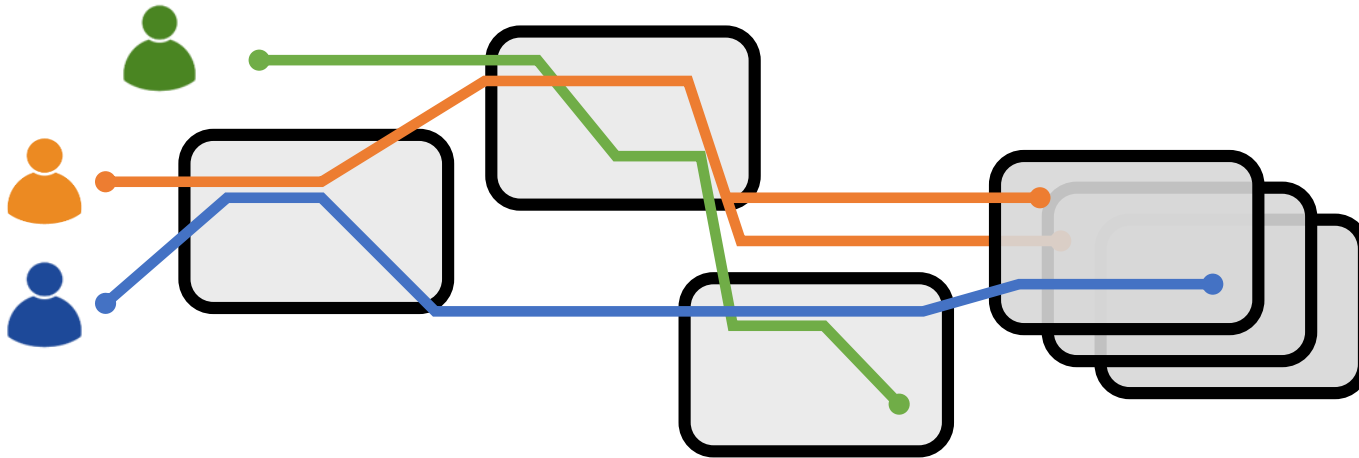
eg, all requests from a tenant over time

eg, data balancing in HDFS

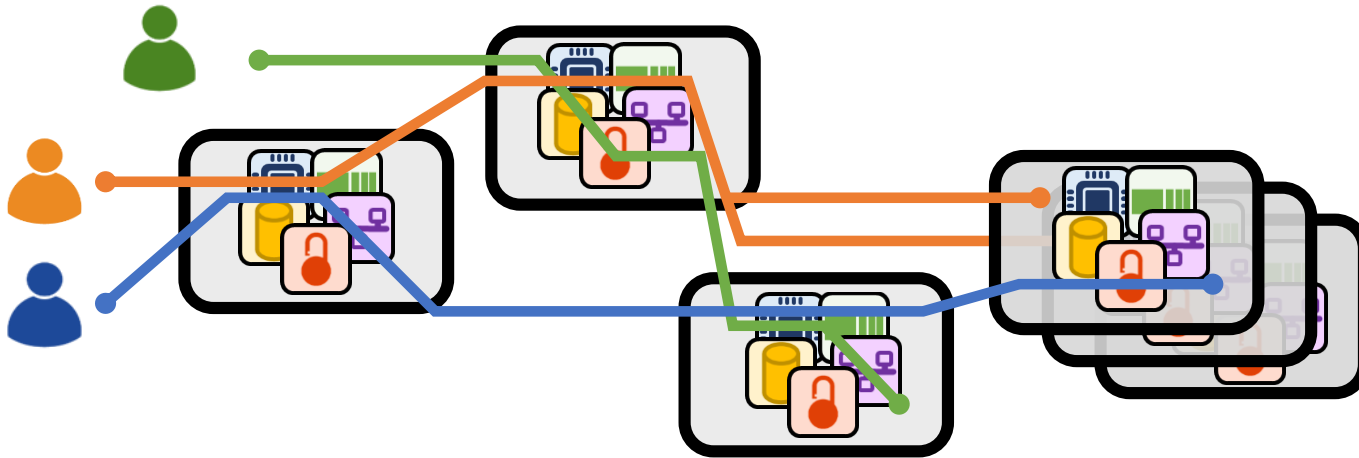
Unit of resource measurement, attribution, and enforcement

Tracks a request across varying levels of granularity

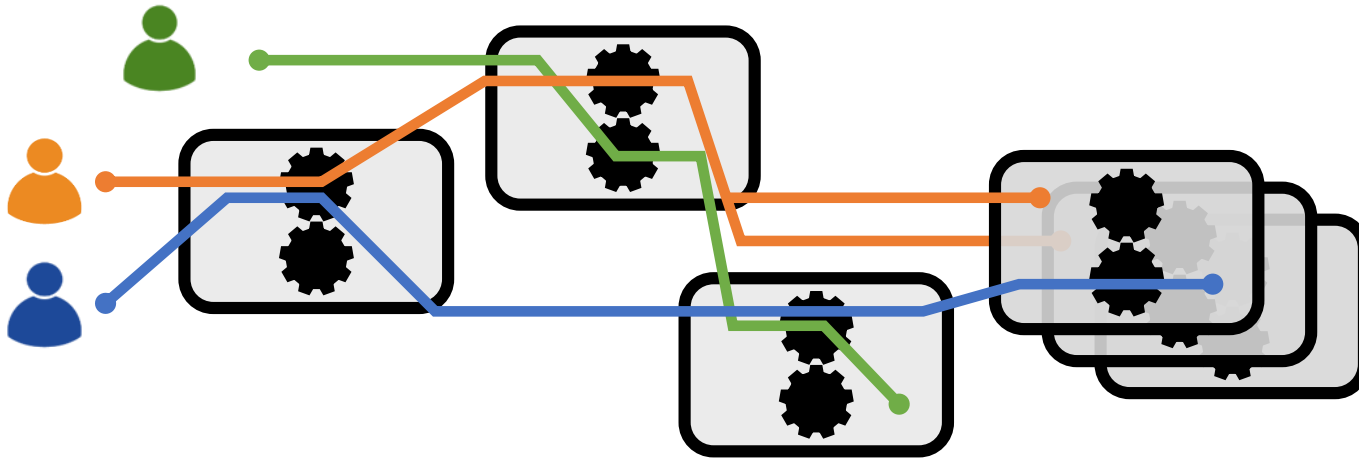
Orthogonal to threads, processes, network flows, etc.



 Workflows



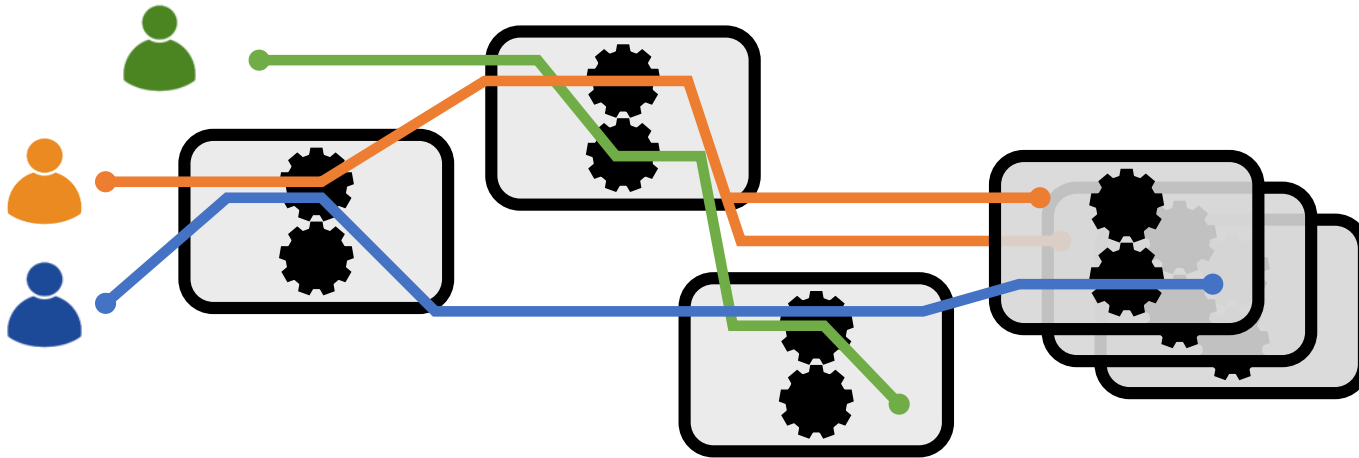
 Workflows





   Workflows       Resources

Purpose: cope with diversity of resources



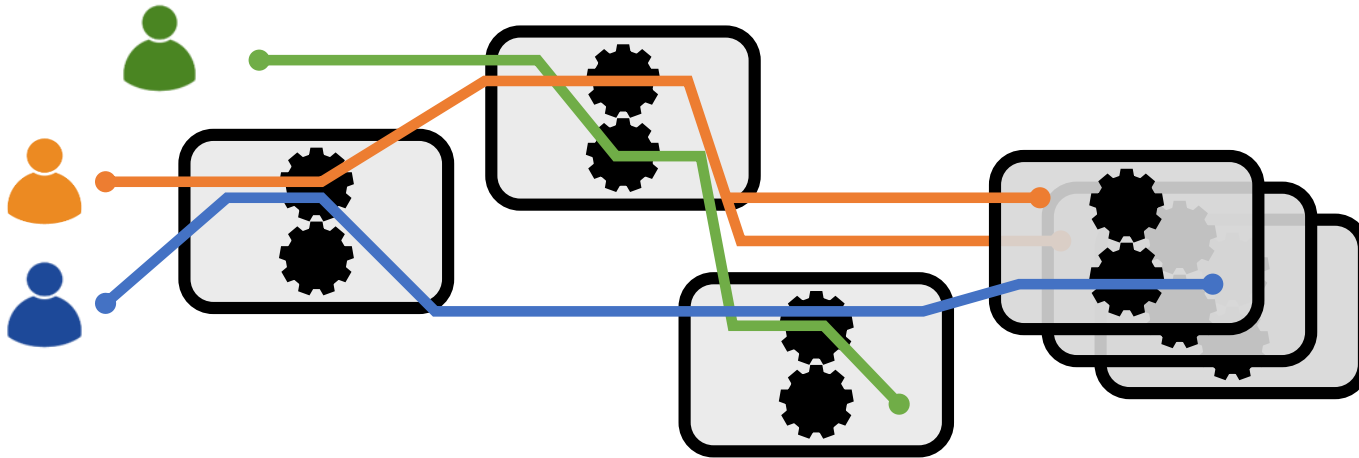



 Workflows       Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources

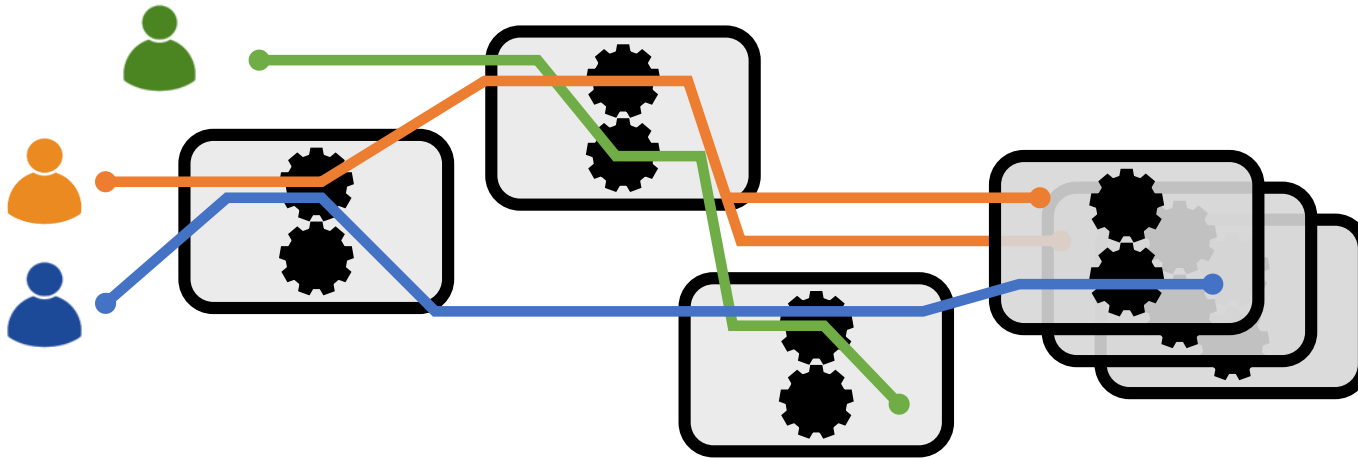




 Workflows       Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources
2. Identify culprit workflows



 Workflows       Resources

Purpose: cope with diversity of resources

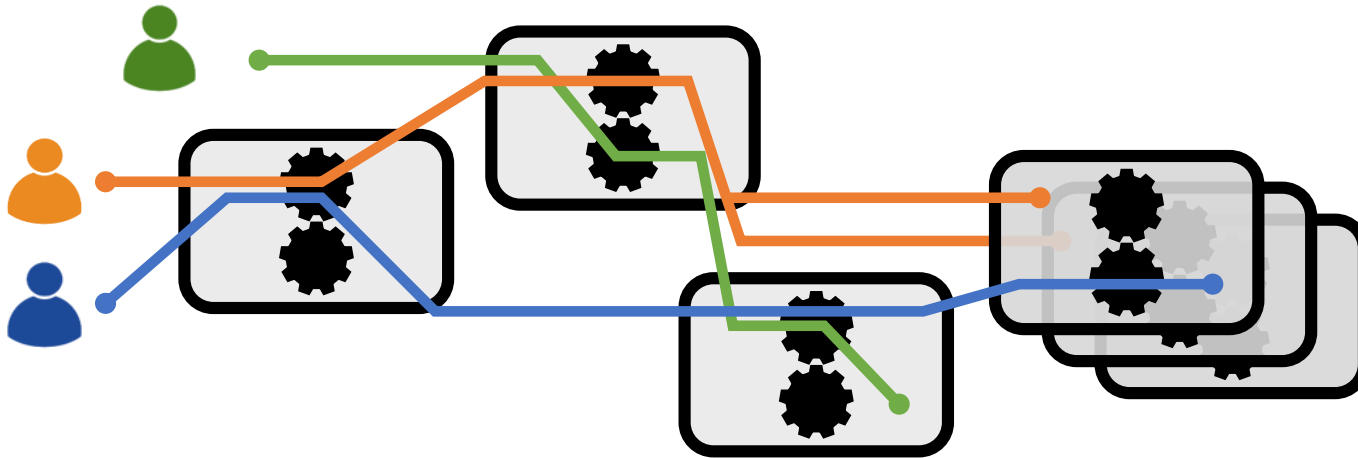
What we need:

1. Identify overloaded resources

**Slowdown**

Ratio of how slow the resource is now compared to its baseline performance with no contention.

2. Identify culprit workflows



 Workflows       Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources

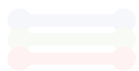
**Slowdown**

Ratio of how slow the resource is now compared to its baseline performance with no contention.

2. Identify culprit workflows

**Load**

Fraction of current utilization that we can attribute to each workflow



Workflows



Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources

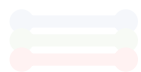
**Slowdown**

Ratio of how slow the resource is now compared to its baseline performance with no contention.

2. Identify culprit workflows

**Load**

Fraction of current utilization that we can attribute to each workflow



Workflows



Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources

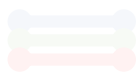
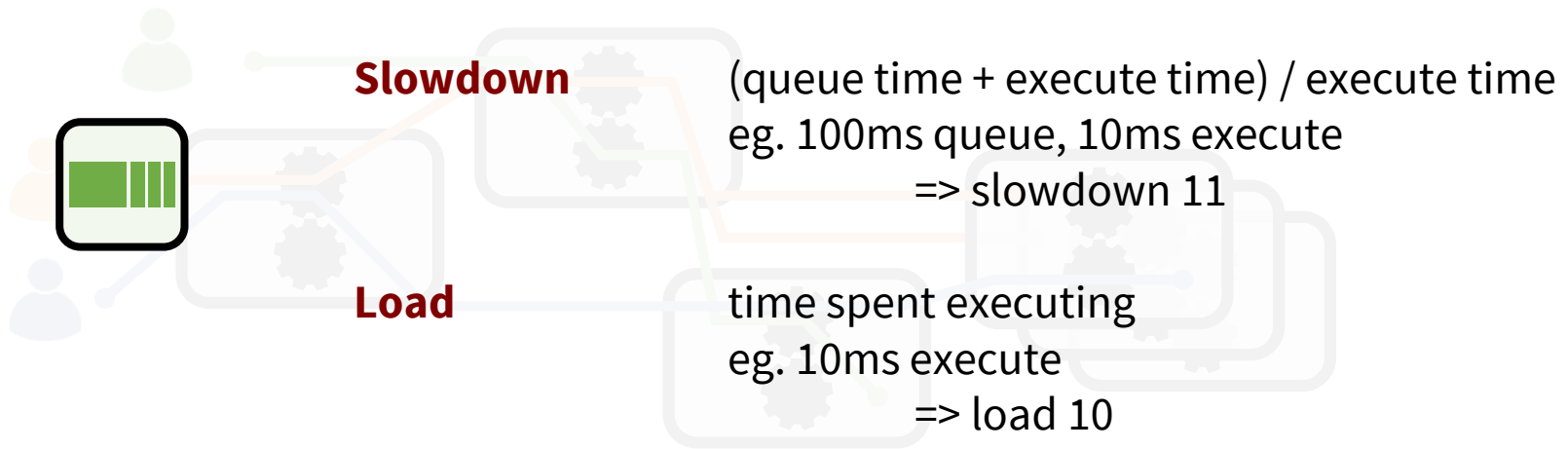
**Slowdown**

Ratio of how slow the resource is now compared to its baseline performance with no contention.

2. Identify culprit workflows

**Load**

Fraction of current utilization that we can attribute to each workflow



Workflows



Resources

Purpose: cope with diversity of resources

What we need:

1. Identify overloaded resources

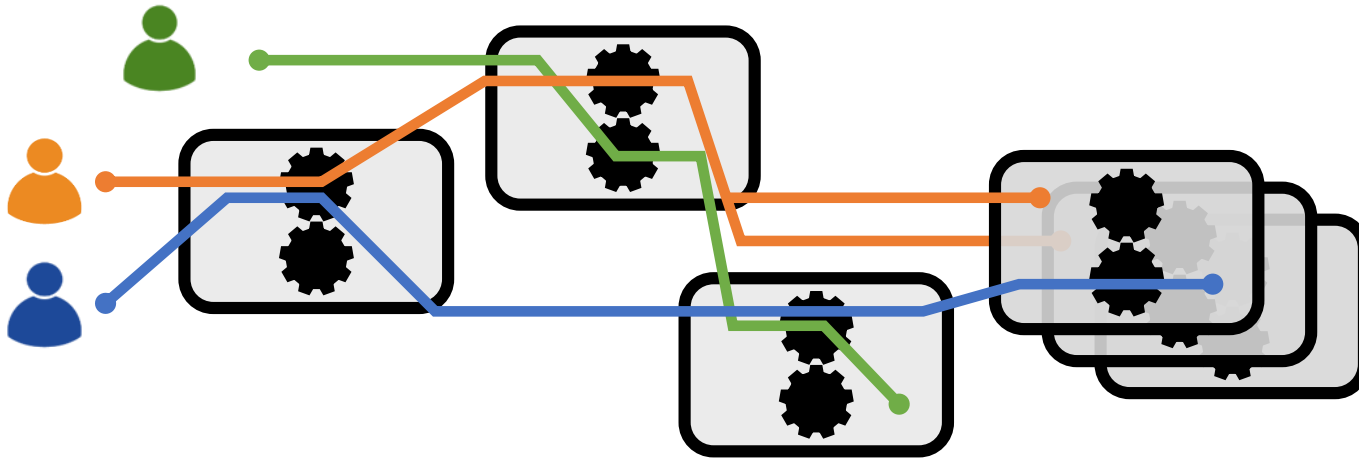
### **Slowdown**

Ratio of how slow the resource is now compared to its baseline performance with no contention.

2. Identify culprit workflows

### **Load**

Fraction of current utilization that we can attribute to each workflow

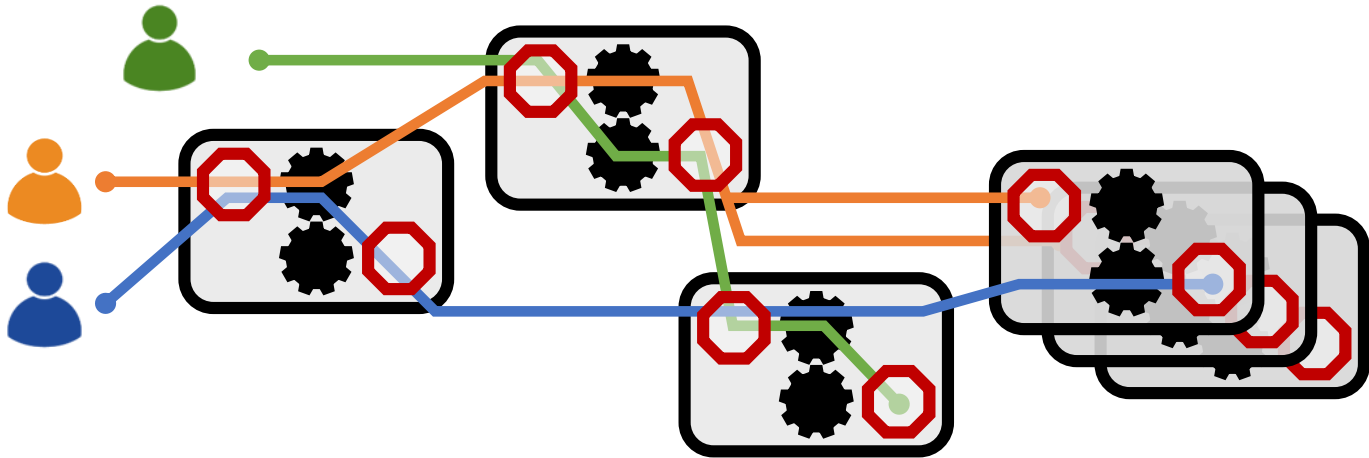


Workflows



Resources





Workflows

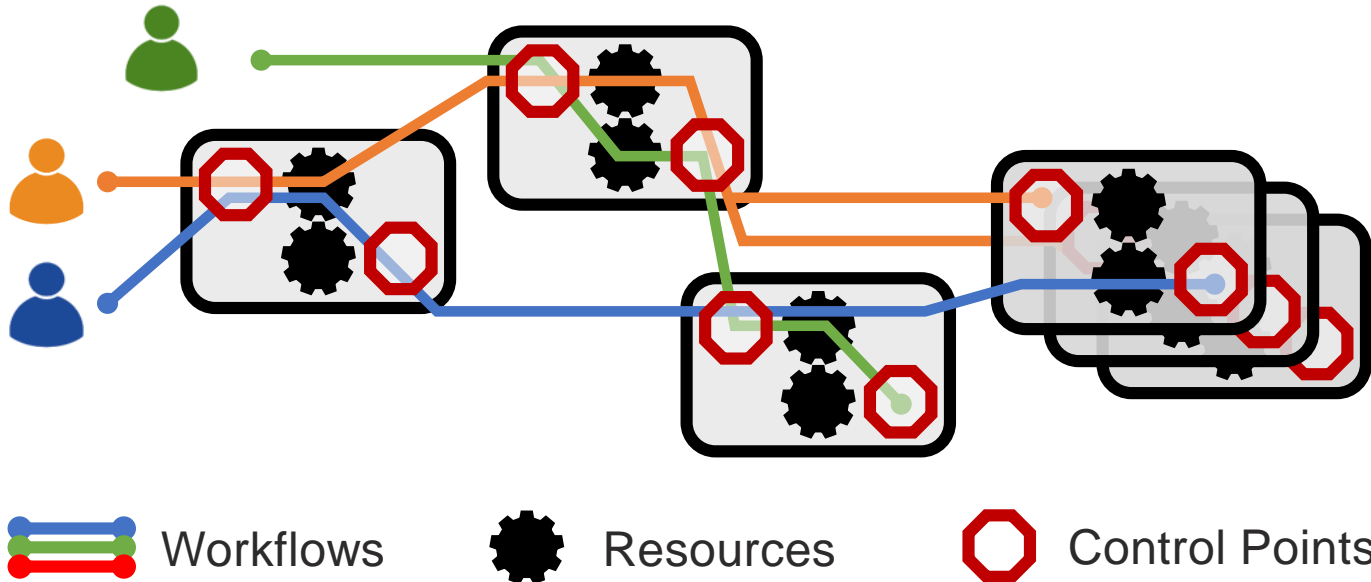


Resources



Control Points

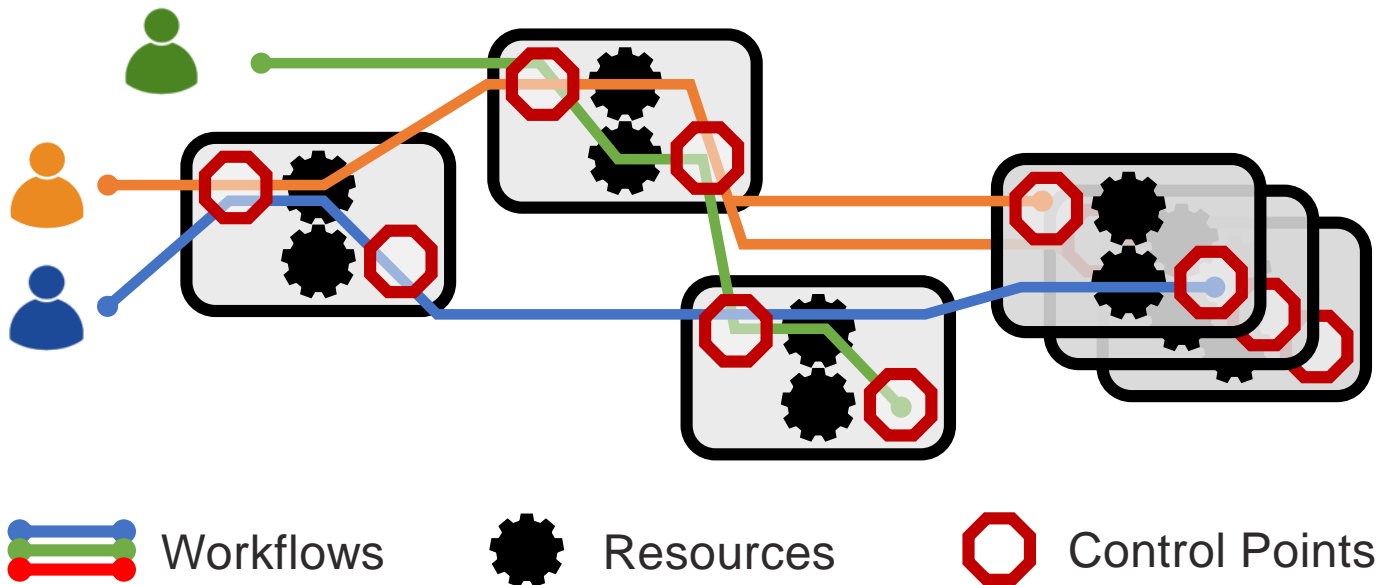
Goal: enforce resource management decisions



Goal: enforce resource management decisions

Decoupled from resources

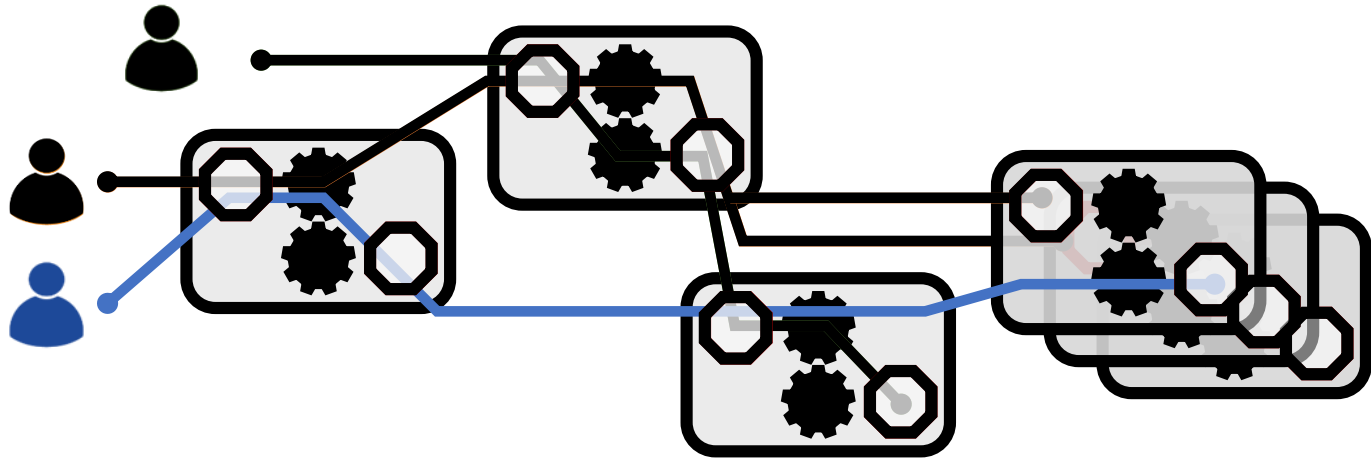
Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue



Goal: enforce resource management decisions

Decoupled from resources

Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue

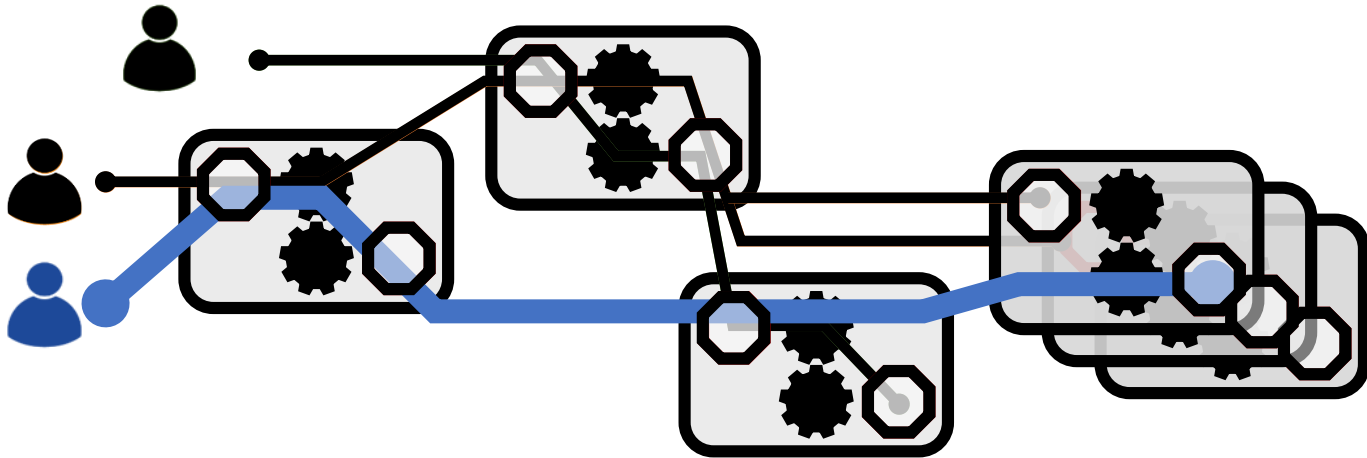


 Workflows       Resources       Control Points

Goal: enforce resource management decisions

Decoupled from resources

Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue

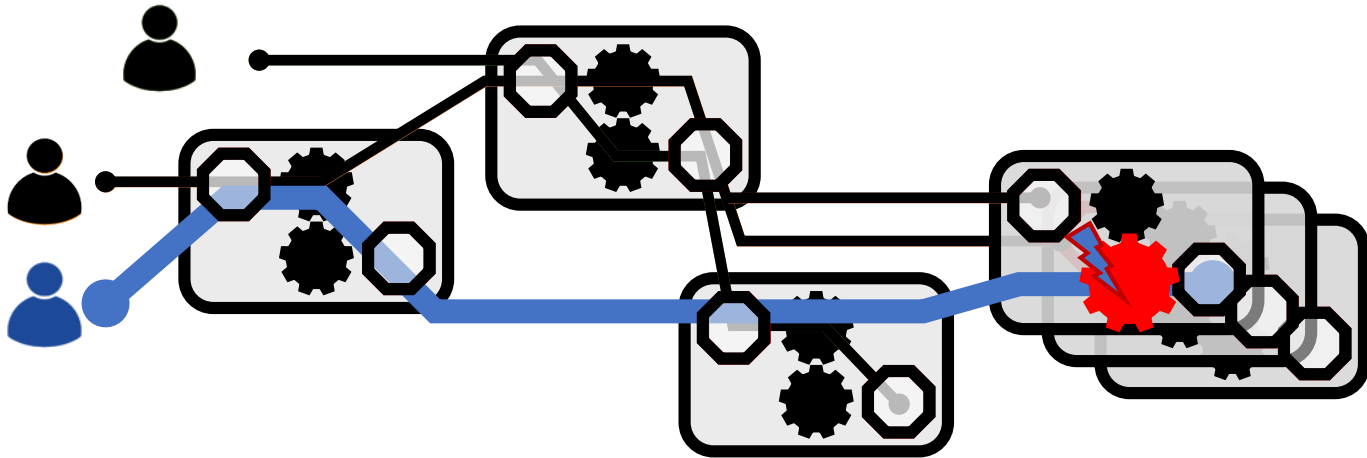


 Workflows       Resources       Control Points

Goal: enforce resource management decisions

Decoupled from resources

Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue



 Workflows       Resources       Control Points

Goal: enforce resource management decisions

Decoupled from resources

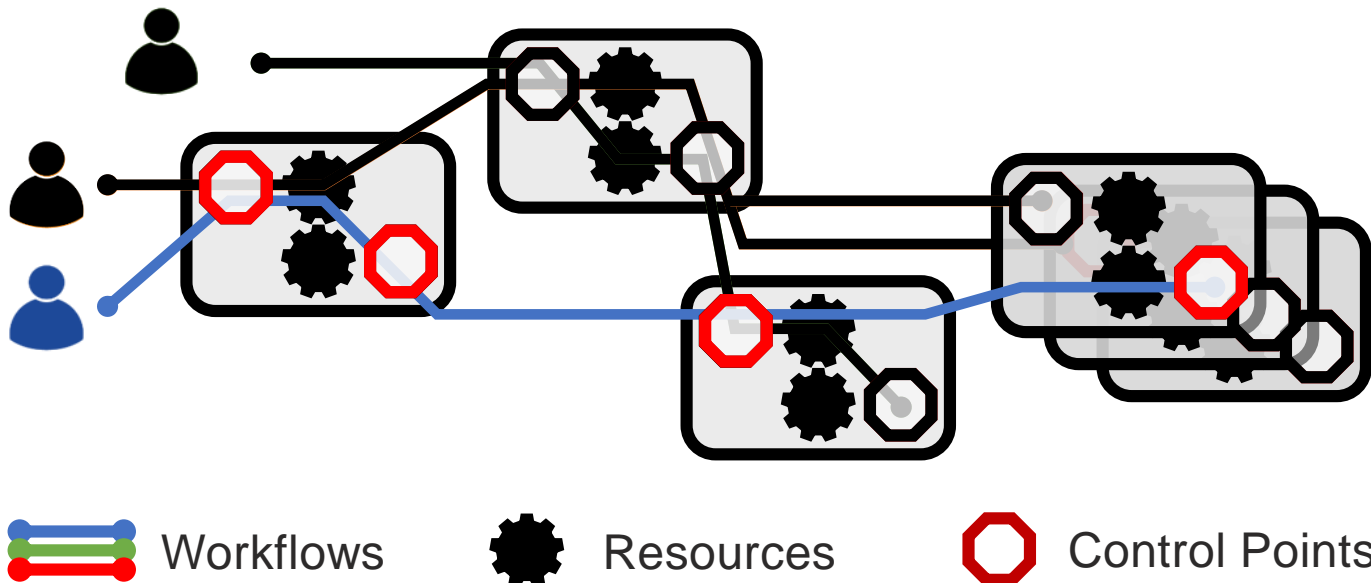
Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue



## Decoupled from resources

# token bucket

## priority queue

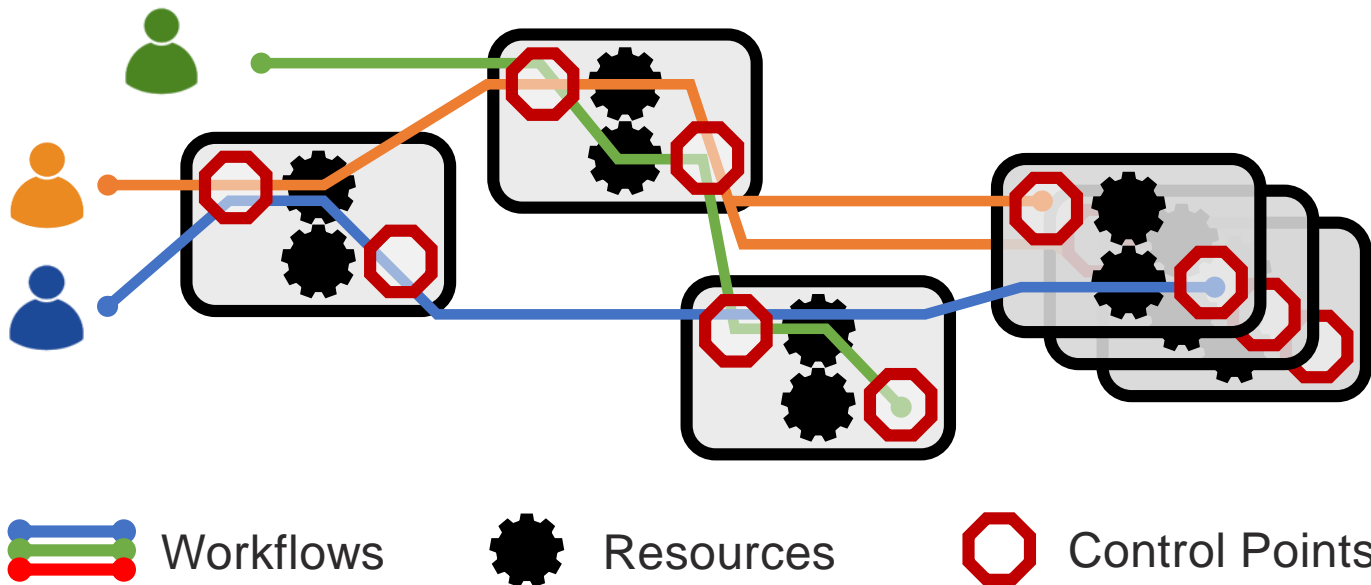


Goal: enforce resource management decisions

Decoupled from resources

Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue

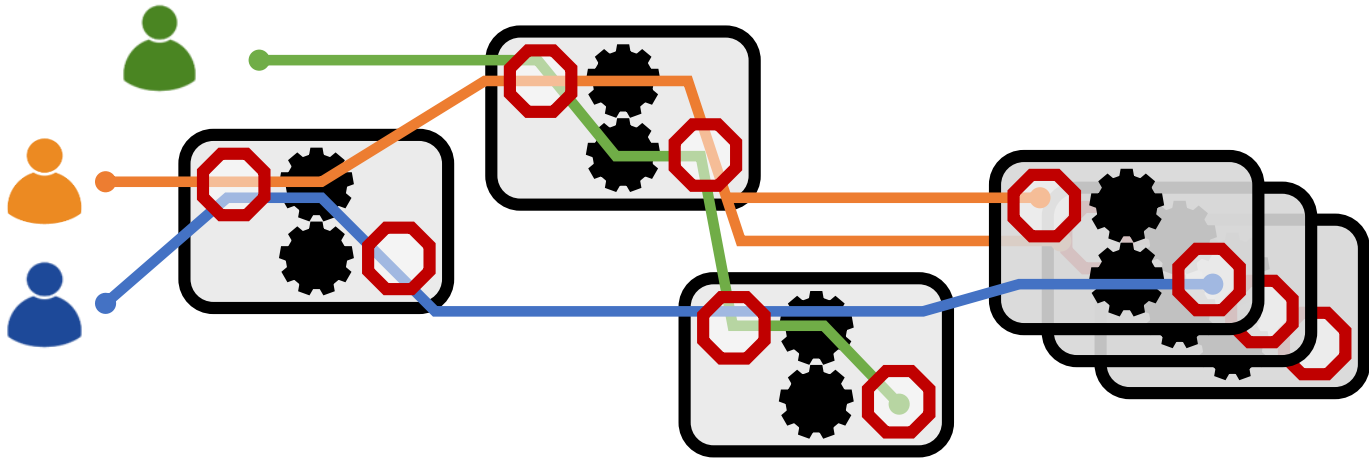




Goal: enforce resource management decisions

Decoupled from resources

Rate-limits workflows, agnostic to underlying implementation e.g.,  
token bucket  
priority queue



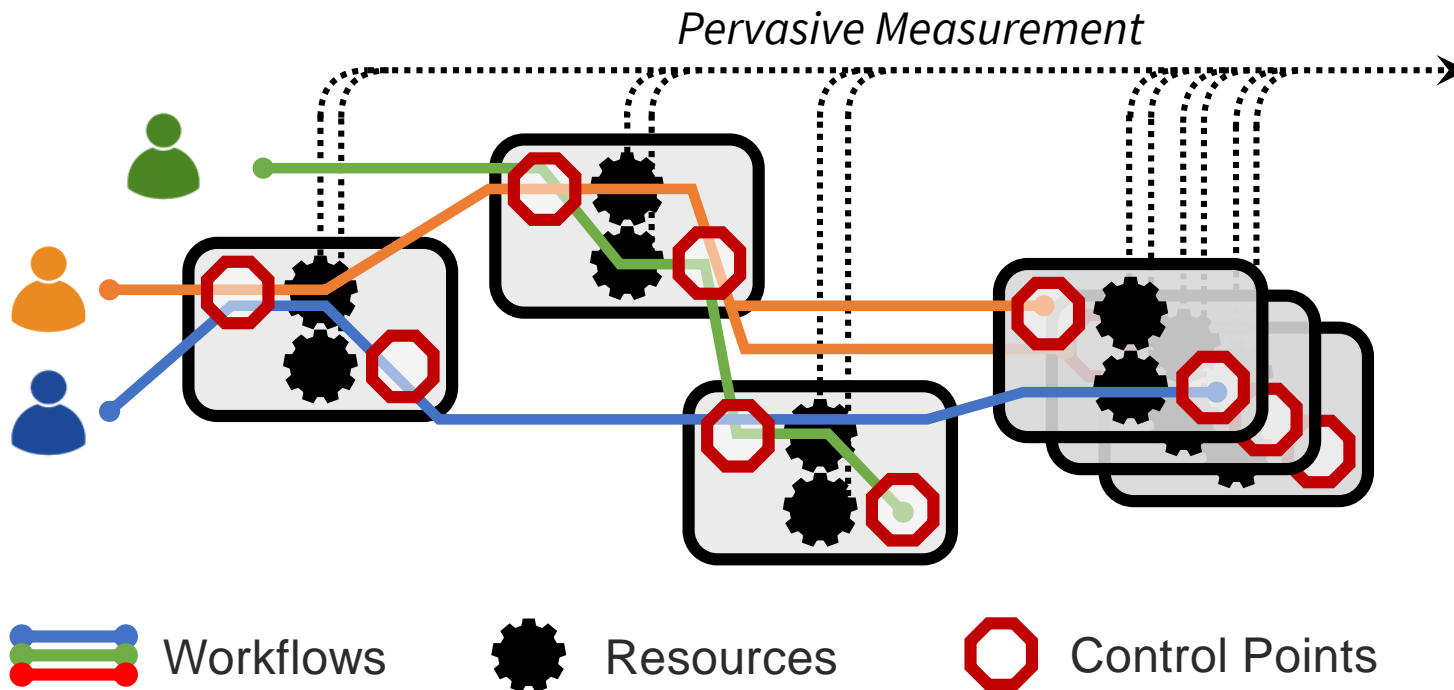
Workflows



Resources

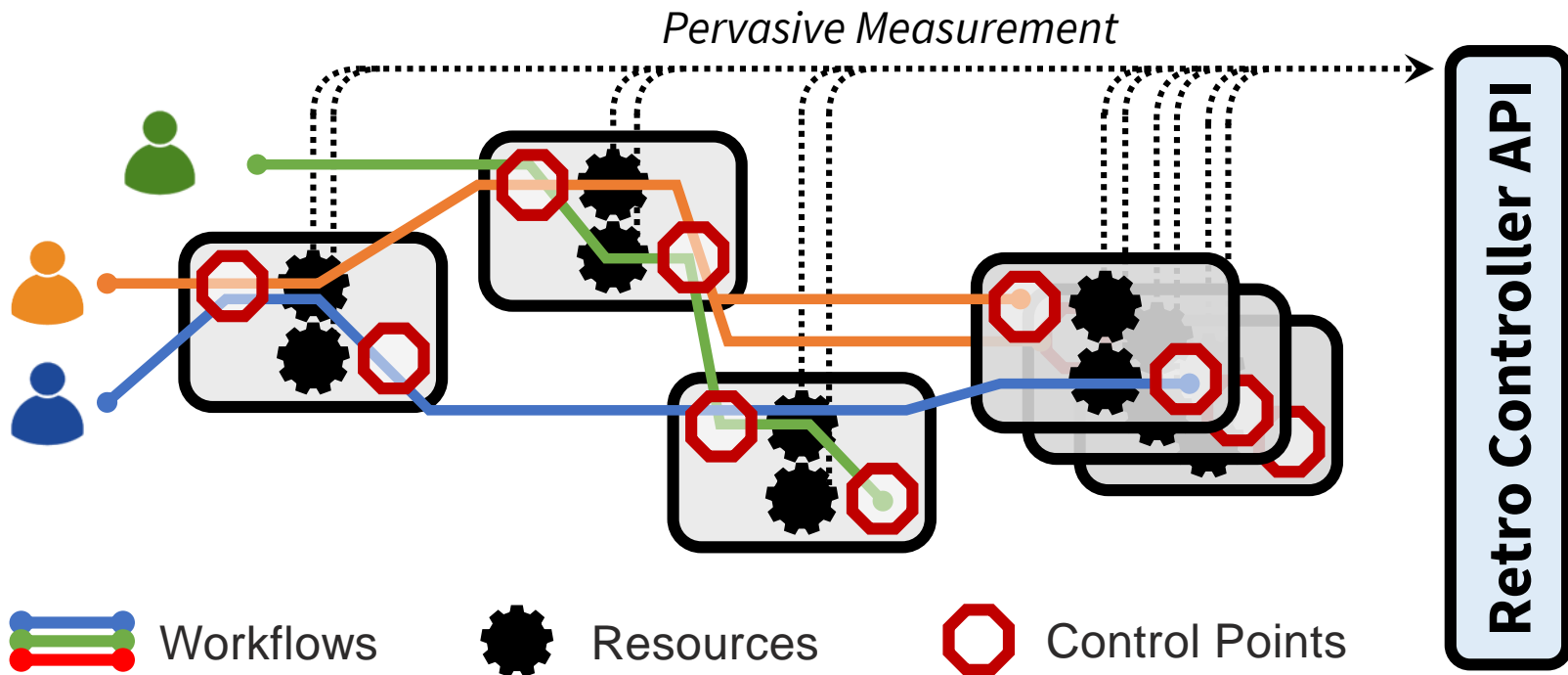


Control Points



## 1. Pervasive Measurement

Aggregated locally then reported centrally once per second

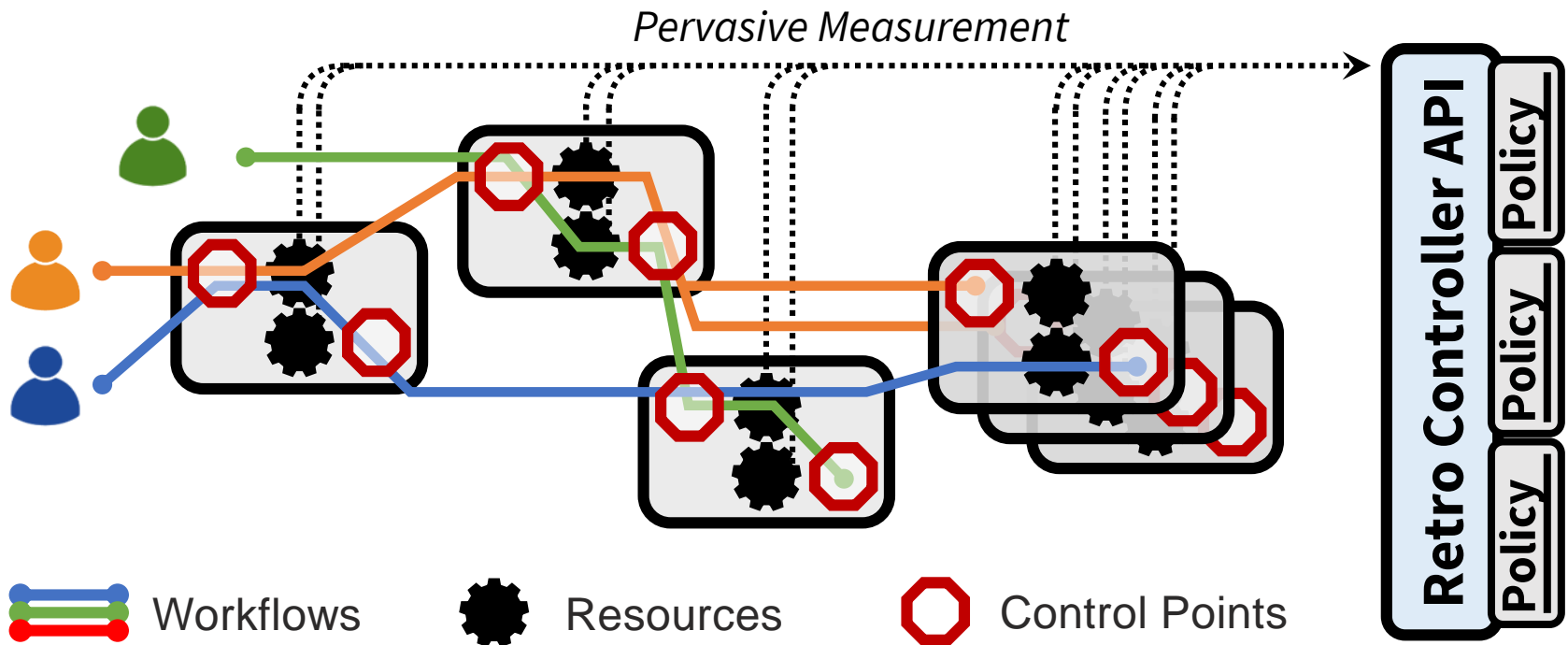


## 1. Pervasive Measurement

Aggregated locally then reported centrally once per second

## 2. Centralized Controller

Global, abstracted view of the system



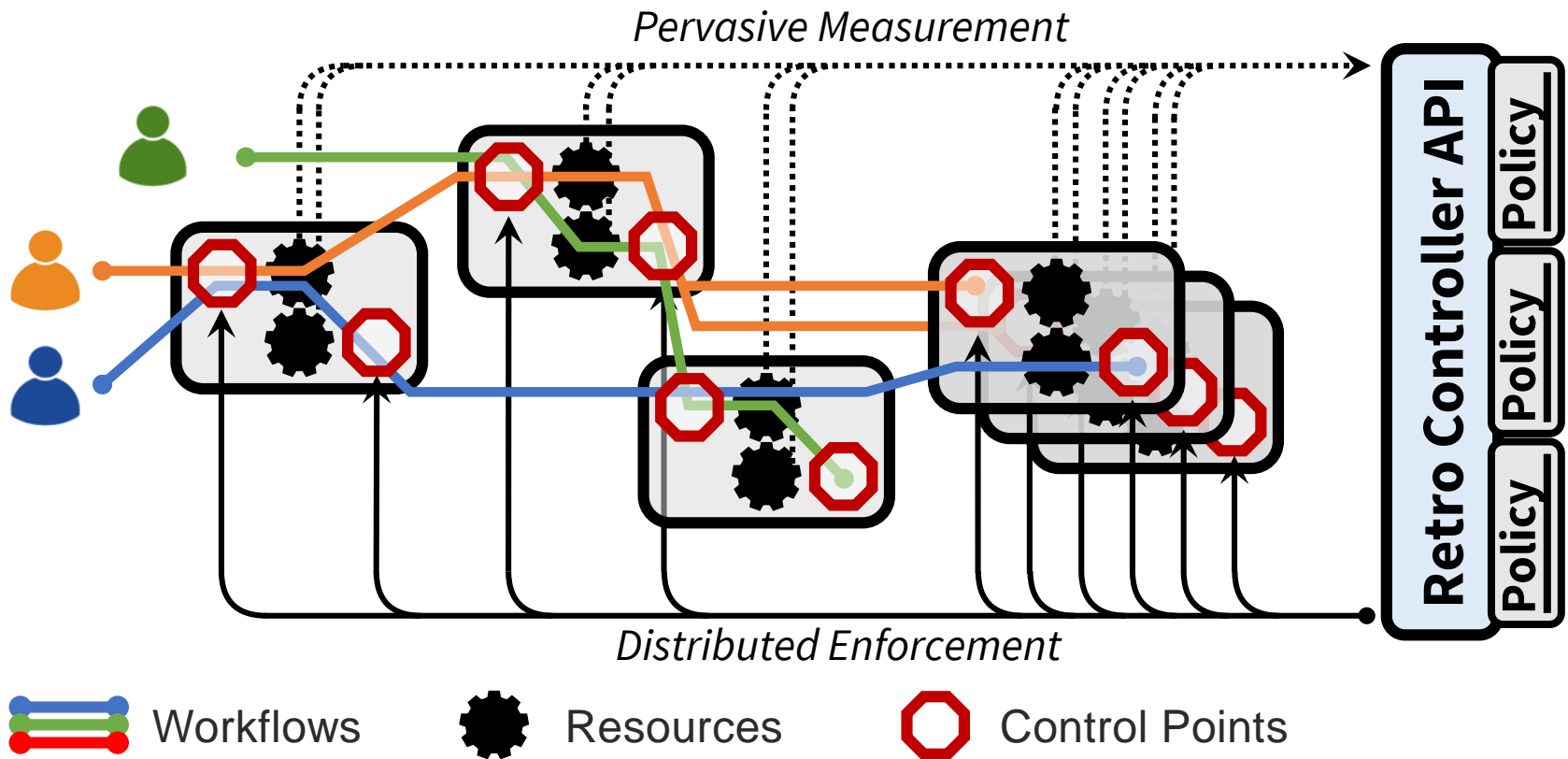
## 1. Pervasive Measurement

Aggregated locally then reported centrally once per second

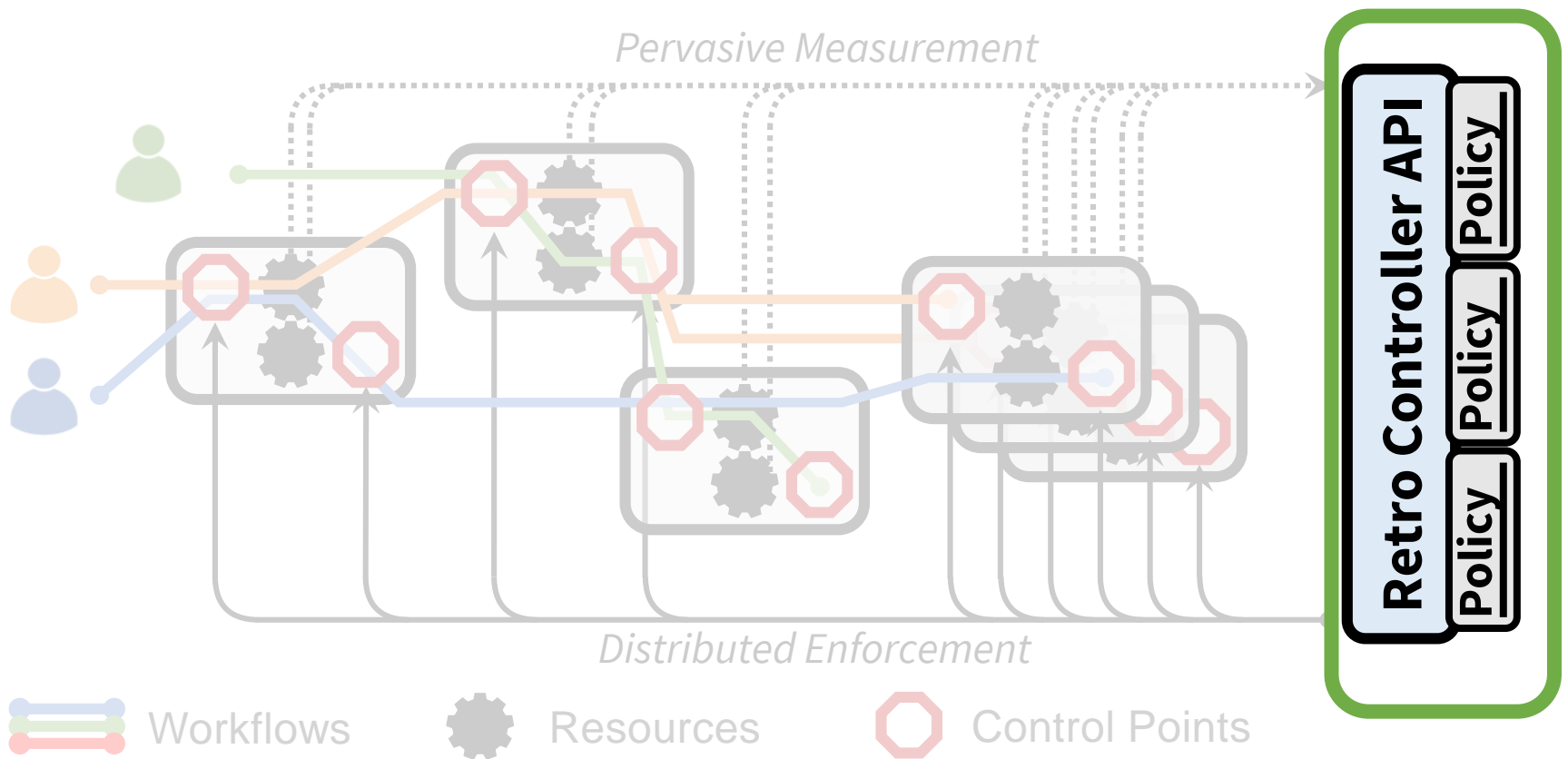
## 2. Centralized Controller

Global, abstracted view of the system

Policies run in continuous control loop



1. Pervasive Measurement  
Aggregated locally then reported centrally once per second
2. Centralized Controller  
Global, abstracted view of the system  
Policies run in continuous control loop
3. Distributed Enforcement  
Co-ordinates enforcement using distributed token bucket



“Control Plane” for resource management

Global, abstracted view of the system

Easier to write

Reusable



# Example: LatencySLO

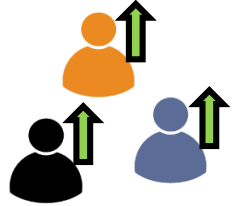


# Example: LatencySLO

H High Priority Workflows



*"200ms average request latency"*

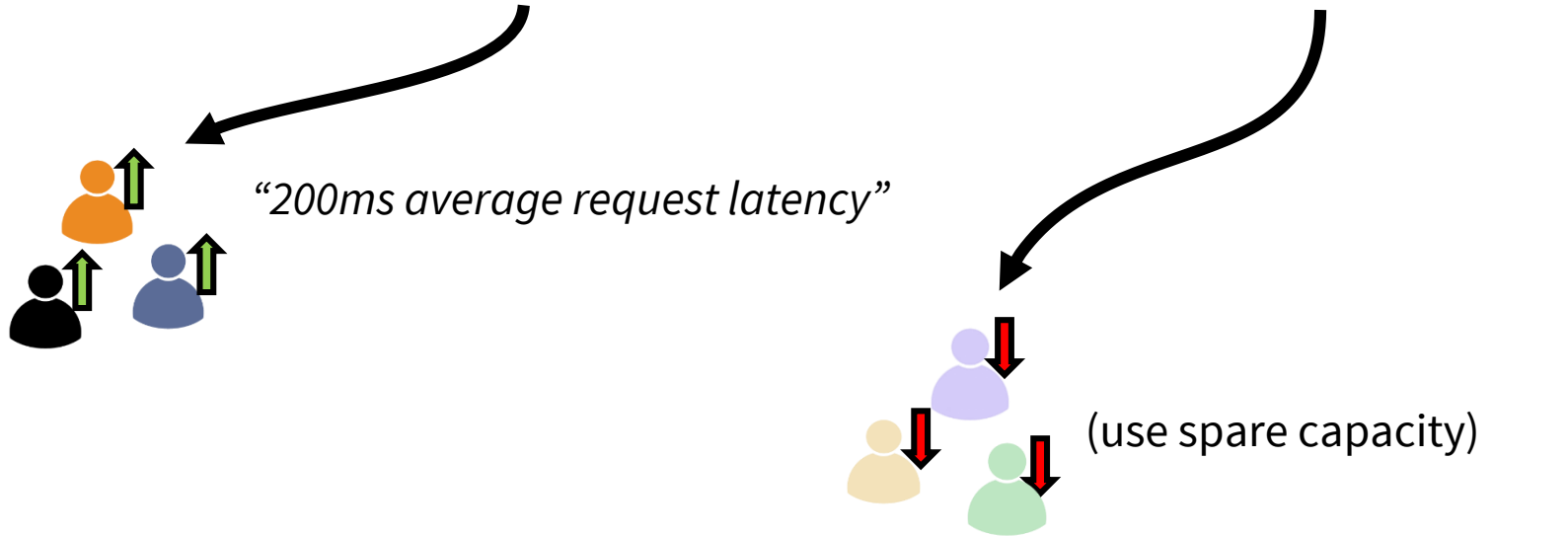


## Policy

# Example: LatencySLO

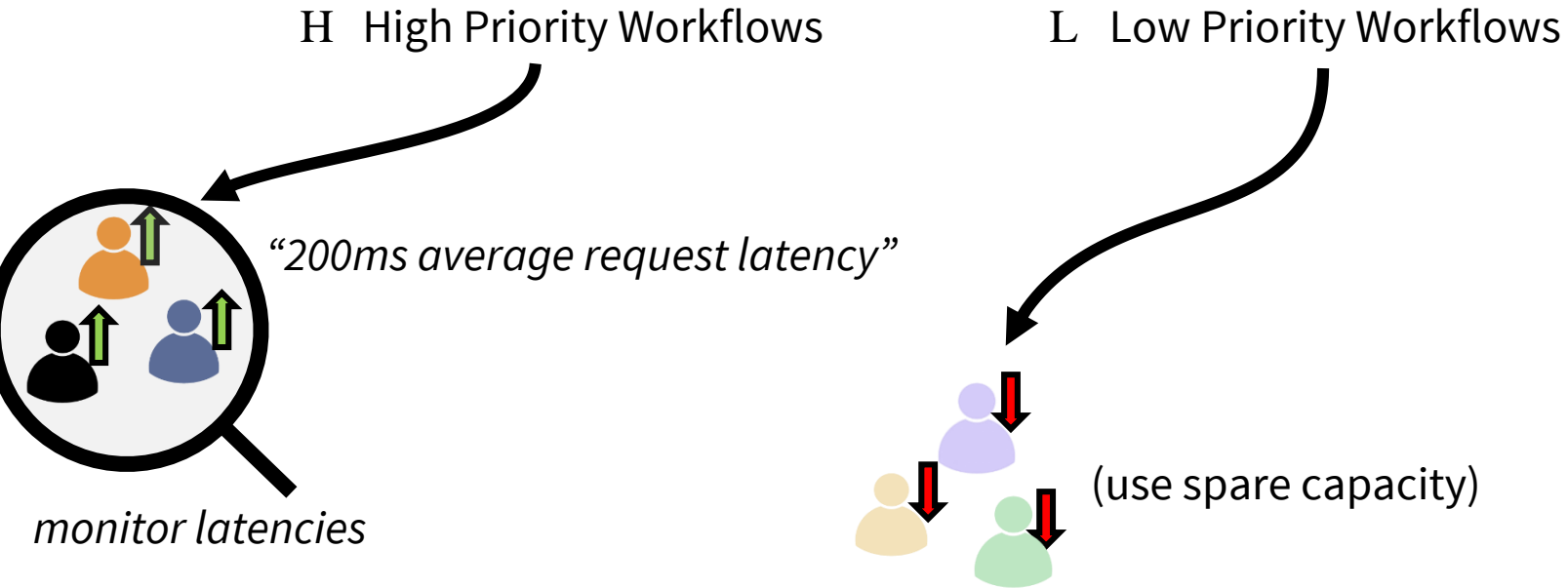
H High Priority Workflows

L Low Priority Workflows



## Policy

# Example: LatencySLO

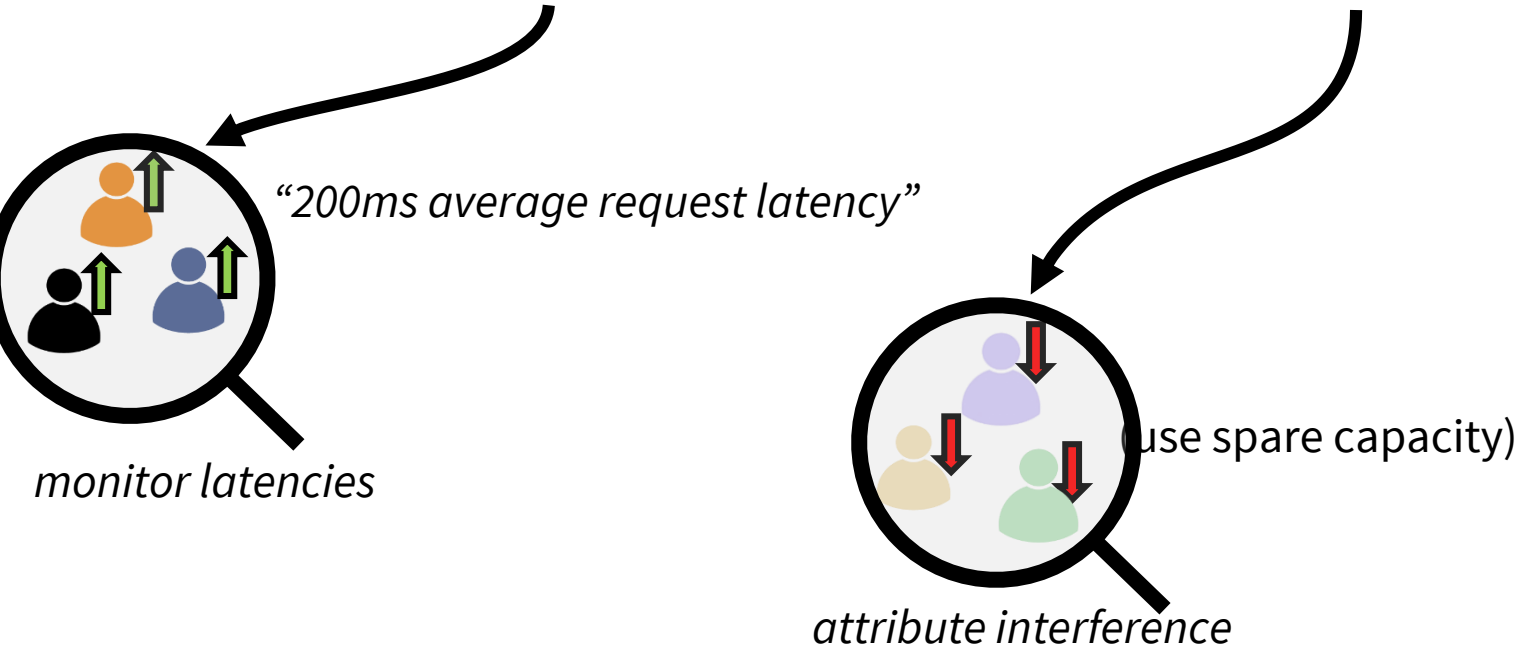


## Policy

# Example: LatencySLO

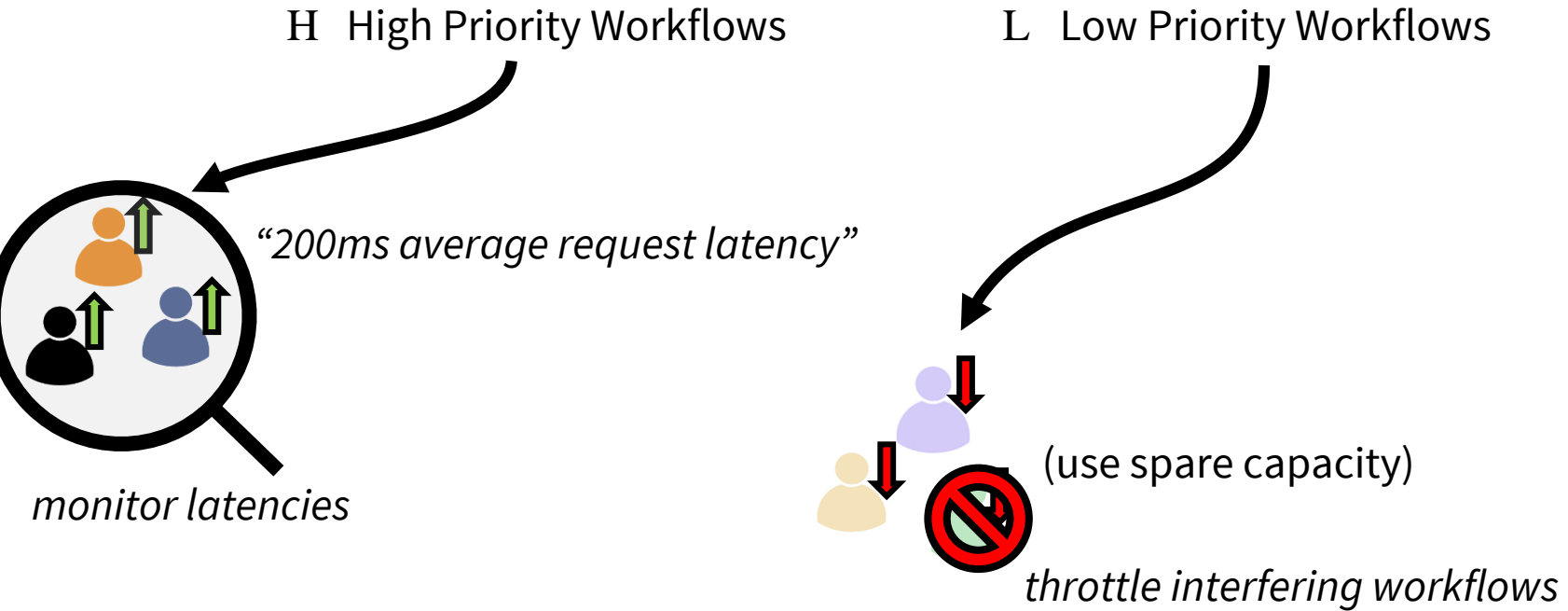
H High Priority Workflows

L Low Priority Workflows



## Policy

# Example: LatencySLO



# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

```
1  foreach candidate in H
2    miss[candidate] = latency(candidate) / guarantee[candidate]
3  W = candidate in H with max miss[candidate]

4  foreach rsrc in resources() // calculate importance of each resource for hipri
5    importance[rsrc] = latency(W, rsrc) * log(slowdown(rsrc))

6  foreach lopri in L // calculate low priority workflow interference
7    interference[lopri] =  $\sum_{rsrc}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)

8  foreach lopri in L // normalize interference
9    interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11   if miss[W] > 1 // throttle
12     scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]
13   else // release
14     scalefactor = 1 +  $\beta$ 

15   foreach cpoint in controlpoints() // apply new rates
16     set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

```
1  foreach candidate in H
2    miss[candidate] = latency(candidate) / guarantee[candidate]
3  W = candidate in H with max miss[candidate]
```

```
4  foreach rsrc in resources() // calculate importance of each resource for hipri
5    importance[rsrc] = latency(W, rsrc) * log(slowdown(rsrc))

6  foreach lopri in L // calculate low priority workflow interference
7    interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

```
8  foreach lopri in L // normalize interference
9    interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11   if miss[W] > 1 // throttle
12     scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]
13   else // release
14     scalefactor = 1 +  $\beta$ 

15 foreach cpoint in controlpoints() // apply new rates
16   set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

```
1  foreach candidate in H  
2  miss[candidate] = latency(candidate, resources)  
3  W = candidate in H with max miss[candidate]
```

**Select the high priority workflow  $W$  with worst performance**

```
4  foreach rsrc in resources() // calculate importance of each resource for hipri  
5  importance[rsrc] = latency(W, rsrc) * log(slowdown(rsrc))  
  
6  foreach lopri in L // calculate low priority workflow interference  
7  interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

```
8  foreach lopri in L // normalize interference  
9  interference[lopri] /=  $\sum_k$  interference[k]  
  
10 foreach lopri in L  
11  if miss[W] > 1 // throttle  
12    scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]  
13  else // release  
14    scalefactor = 1 +  $\beta$   
  
15  foreach cpoint in controlpoints() // apply new rates  
16  set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```



# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

```
1  foreach candidate in H
2  miss[candidate] = latency(candidate, resources)
3  W = candidate in H with max miss[candidate]
```

**Select the high priority workflow  $W$  with worst performance**

```
4  foreach rsrc in resources() // calculate importance of each resource for hipri
5  importance[rsrc] = latency(W, rsrc) * log(load(rsrc))
6  foreach lopri in L // calculate low priority workflow interference
7  interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

**Weight low priority workflows by their interference with  $W$**

```
8  foreach lopri in L // normalize interference
9  interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11  if miss[W] > 1 // throttle
12    scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]
13  else // release
14    scalefactor = 1 +  $\beta$ 

15 foreach cpoint in controlpoints() // apply new rates
16  set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

```
1 foreach candidate in H
2   miss[candidate] = latency(candidate, resources)
3 W = candidate in H with max miss[candidate]
```

**Select the high priority workflow  $W$  with worst performance**

```
4 foreach rsrc in resources() // calculate importance of each resource for hipri
5   importance[rsrc] = latency(W, rsrc) * log(load(rsrc))
6 foreach lopri in L // calculate low priority workflow interference
7   interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

**Weight low priority workflows by their interference with  $W$**

```
8 foreach lopri in L // normalize interference
9   interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11   scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]
12   else // release
13     scalefactor = 1 +  $\beta$ 

15 foreach cpoint in controlpoints() // apply new rates
16   set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

**Throttle low priority workflows proportionally to their weight**

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

**Select the high priority workflow  $W$  with worst performance**

```
1  foreach candidate in H
2    miss[candidate] = latency(candidate) / guarantee[candidate]
3   $W$  = candidate in H with max miss[candidate]
```

```
4  foreach rsrc in resources() // calculate importance of each resource for hipri
5    importance[rsrc] = latency( $W$ , rsrc) * log(slowdown(rsrc))
6  foreach lopri in L // calculate low priority workflow interference
7    interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

**Weight low priority workflows by their interference with  $W$**

```
8  foreach lopri in L // normalize interference
9    interference[lopri] /=  $\sum_k$  interference[k]
```

```
10 foreach lopri in L
11   if miss[ $W$ ] > 1 // throttle
12   else // release
13     scalefactor = 1 +  $\beta$ 
```

**Throttle low priority workflows proportionally to their weight**

```
15 foreach cpoint in controlpoints() // apply new rates
16   set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

**Select the high priority workflow  $W$  with worst performance**

```

1  foreach candidate in H
2    miss[candidate] = latency(candidate) / guarantee[candidate]
3  W = candidate in H with max miss[candidate]
```

**Weight low priority workflows by their interference with  $W$**

```

4  foreach rsrc in resources() // calculate importance of each resource for hipri
5    importance[rsrc] = latency(W, rsrc) * log(slowdown(rsrc))

6  foreach lopri in L // calculate low priority workflow interference
7    interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

```

8  foreach lopri in L // normalize interference
9    interference[lopri] /=  $\sum_k$  interference[k]
```

**Throttle low priority workflows proportionally to their weight**

```

10 foreach lopri in L
11   if miss[W] > 1 // throttle
12   else // release
13     scalefactor = 1 +  $\beta$ 

15   foreach cpoint in controlpoints() // apply new rates
16     set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

**Select the high priority workflow  $W$  with worst performance**

```
1 foreach candidate in H
2   miss[candidate] = latency(candidate) / guarantee[candidate]
3 W = candidate in H with max miss[candidate]
```

**Weight low priority workflows by their interference with  $W$**

```
4 foreach rsrc in resources() // calculate importance of each resource for hipri
5   importance[rsrc] = latency(W, rsrc) * log(slowdown(rsrc))

6 foreach lopri in L // calculate low priority workflow interference
7   interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)
```

**Throttle low priority workflows proportionally to their weight**

```
8 foreach lopri in L // normalize interference
9   interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11   if miss[W] > 1 // throttle
12     scalefactor = 1 -  $\alpha$  * (miss[W] - 1) * interference[lopri]
13   else // release
14     scalefactor = 1 +  $\beta$ 

15 foreach cpoint in controlpoints() // apply new rates
16   set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))
```

# Example: LatencySLO

H High Priority Workflows

L Low Priority Workflows

**Select the high priority workflow  $W$  with worst performance**

```

1  foreach candidate in H
2    miss[candidate] = latency(candidate) / guarantee[candidate]
3   $W$  = candidate in H with max miss[candidate]

```

**Weight low priority workflows by their interference with  $W$**

```

4  foreach rsrc in resources() // calculate importance of each resource for hipri
5    importance[rsrc] = latency( $W$ , rsrc) * log(slowdown(rsrc))

6  foreach lopri in L // calculate low priority workflow interference
7    interference[lopri] =  $\sum_{\text{rsrc}}$  importance[rsrc] * load(lopri, rsrc) / load(rsrc)

```

**Throttle low priority workflows proportionally to their weight**

```

8  foreach lopri in L // normalize interference
9    interference[lopri] /=  $\sum_k$  interference[k]

10 foreach lopri in L
11   if miss[ $W$ ] > 1 // throttle
12     scalefactor =  $1 - \alpha$  * (miss[ $W$ ] - 1) * interference[lopri]
13   else // release
14     scalefactor =  $1 + \beta$ 

15   foreach cpoint in controlpoints() // apply new rates
16     set_rate(cpoint, lopri, scalefactor * get_rate(cpoint, lopri))

```

# Other types of policy...

# Other types of policy...



## Bottleneck Fairness

Detect most overloaded resource

Fair-share resource between tenants using it



# Other types of policy...



## Bottleneck Fairness

Detect most overloaded resource

Fair-share resource between tenants using it



## Dominant Resource Fairness

Estimate demands and capacities from measurements

# Other types of policy...



## Bottleneck Fairness

Detect most overloaded resource

Fair-share resource between tenants using it



## Dominant Resource Fairness

Estimate demands and capacities from measurements

Concise

Any resources can be bottleneck (policy doesn't care)

Not system specific

# *Retro* ***Evaluation***

# Instrumentation

# Instrumentation

Retro implementation in Java

Instrumentation Library

Central controller implementation

# Instrumentation

## Retro implementation in Java

Instrumentation Library

Central controller implementation

## To enable Retro



Propagate Workflow ID within application (like X-Trace, Dapper)



Instrument resources with wrapper classes

# Instrumentation

## Retro implementation in Java

Instrumentation Library

Central controller implementation

## To enable Retro



Propagate Workflow ID within application (like X-Trace, Dapper)



Instrument resources with wrapper classes

## Overheads

# Instrumentation

## Retro implementation in Java

Instrumentation Library

Central controller implementation

## To enable Retro



Propagate Workflow ID within application (like X-Trace, Dapper)



Instrument resources with wrapper classes

## Overheads

Resource instrumentation automatic using AspectJ



# Instrumentation

## Retro implementation in Java

Instrumentation Library

Central controller implementation

## To enable Retro



Propagate Workflow ID within application (like X-Trace, Dapper)



Instrument resources with wrapper classes

## Overheads

Resource instrumentation automatic using AspectJ

Overall 50-200 lines per system to modify RPCs

# Instrumentation

## Retro implementation in Java

Instrumentation Library

Central controller implementation

## To enable Retro



Propagate Workflow ID within application (like X-Trace, Dapper)



Instrument resources with wrapper classes

## Overheads

Resource instrumentation automatic using AspectJ

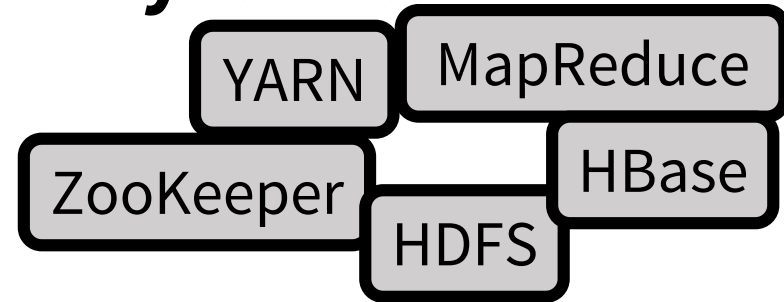
Overall 50-200 lines per system to modify RPCs

Retro overhead: max 1-2% on throughput, latency

# Experiments

# Experiments

## Systems



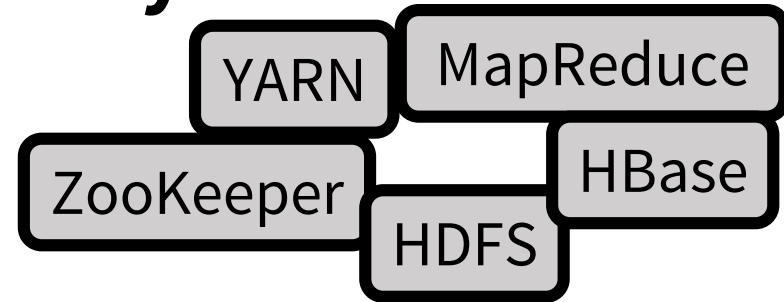
# Experiments

## Workflows



MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Systems



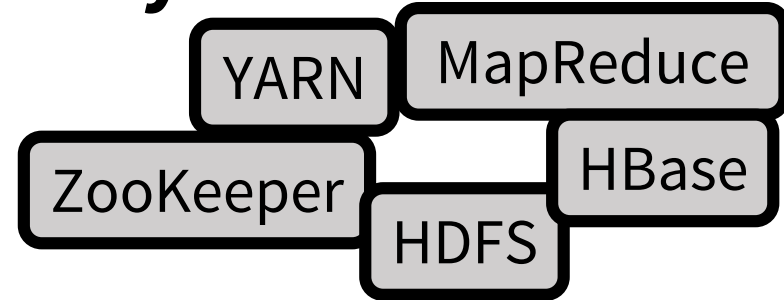
# Experiments

## Workflows



MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Systems



## Resources



CPU, Disk, Network (All systems)  
Locks, Queues (HDFS, HBase)

# Experiments

## Workflows



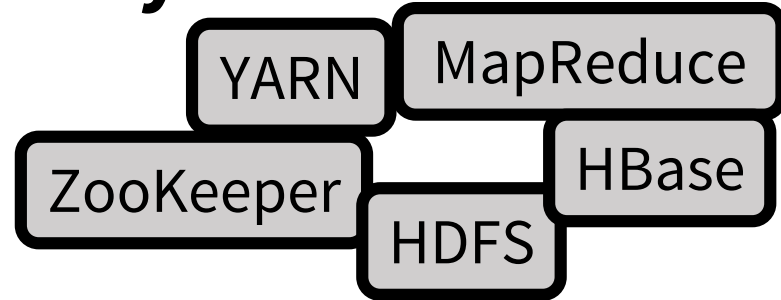
MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Resources



CPU, Disk, Network (All systems)  
Locks, Queues (HDFS, HBase)

## Systems



## Policies



# Experiments

## Workflows



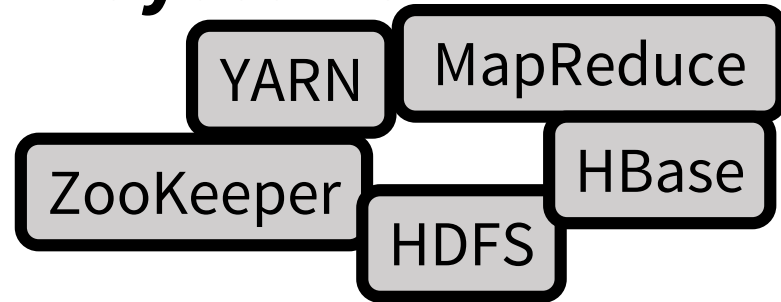
MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Resources



CPU, Disk, Network (All systems)  
Locks, Queues (HDFS, HBase)

## Systems



## Policies



Policies for a mixture of systems, workflows, and resources

Results on clusters up to 200 nodes

See paper for full experiment results



# Experiments

## Workflows



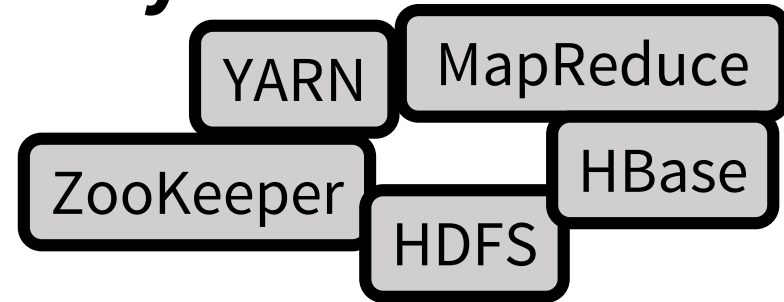
MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Resources



CPU, Disk, Network (All systems)  
Locks, Queues (HDFS, HBase)

## Systems



## Policies



Policies for a mixture of systems, workflows, and resources

Results on clusters up to 200 nodes

See paper for full experiment results

**This talk** LatencySLO policy results

# Experiments

## Workflows



MapReduce Jobs (HiBench)  
HBase (YCSB)  
HDFS clients  
Background Data Replication  
Background Heartbeats

## Resources



CPU, Disk, Network (All systems)  
Locks, Queues (HDFS, HBase)

## Systems



## Policies



Policies for a mixture of systems, workflows, and resources



Results on clusters up to 200 nodes

See paper for full experiment results



**This talk** LatencySLO policy results



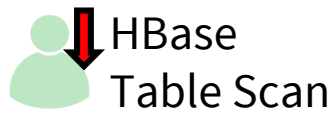
HDFS read 8k



HBase read 1 row



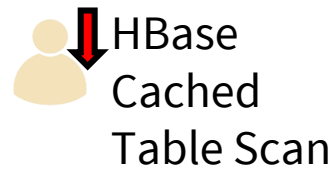
HBase read 1  
cached row



HBase  
Table Scan



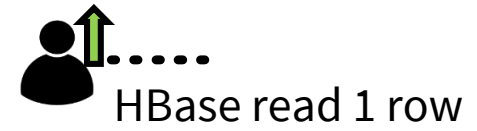
HDFS  
mkdir



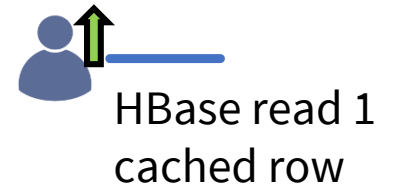
HBase  
Cached  
Table Scan



HDFS read 8k




HBase read 1 row




HBase read 1  
cached row


SLO-Normalized Latency


 HBase  
Table Scan

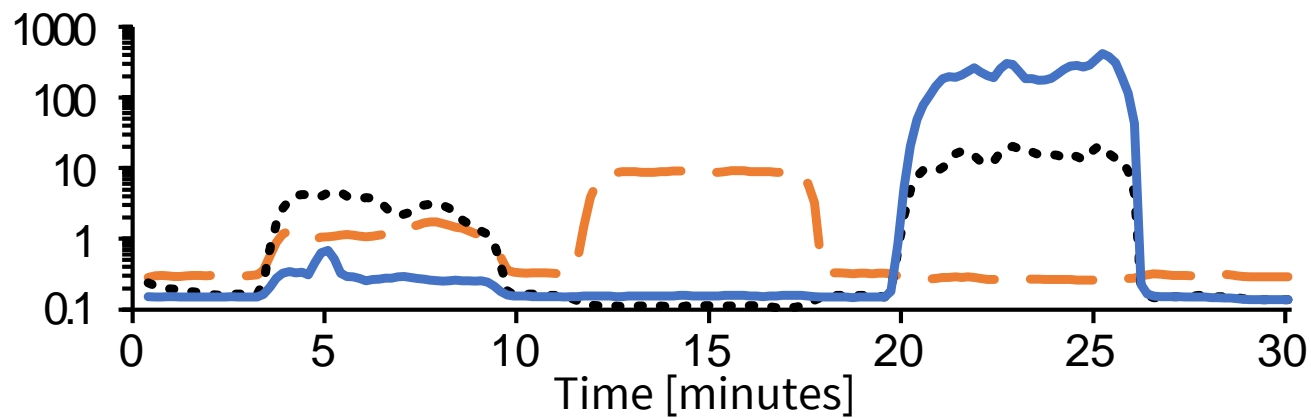
 HDFS  
mkdir

 HBase  
Cached  
Table Scan

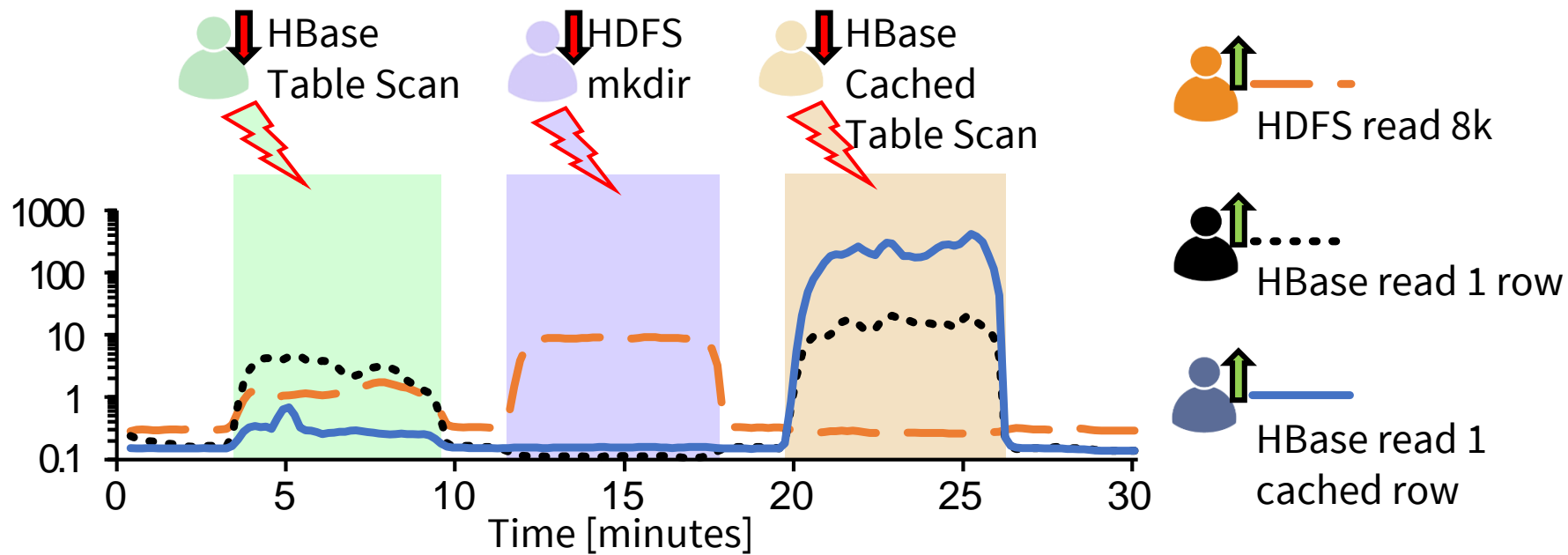
 HDFS read 8k

 HBase read 1 row

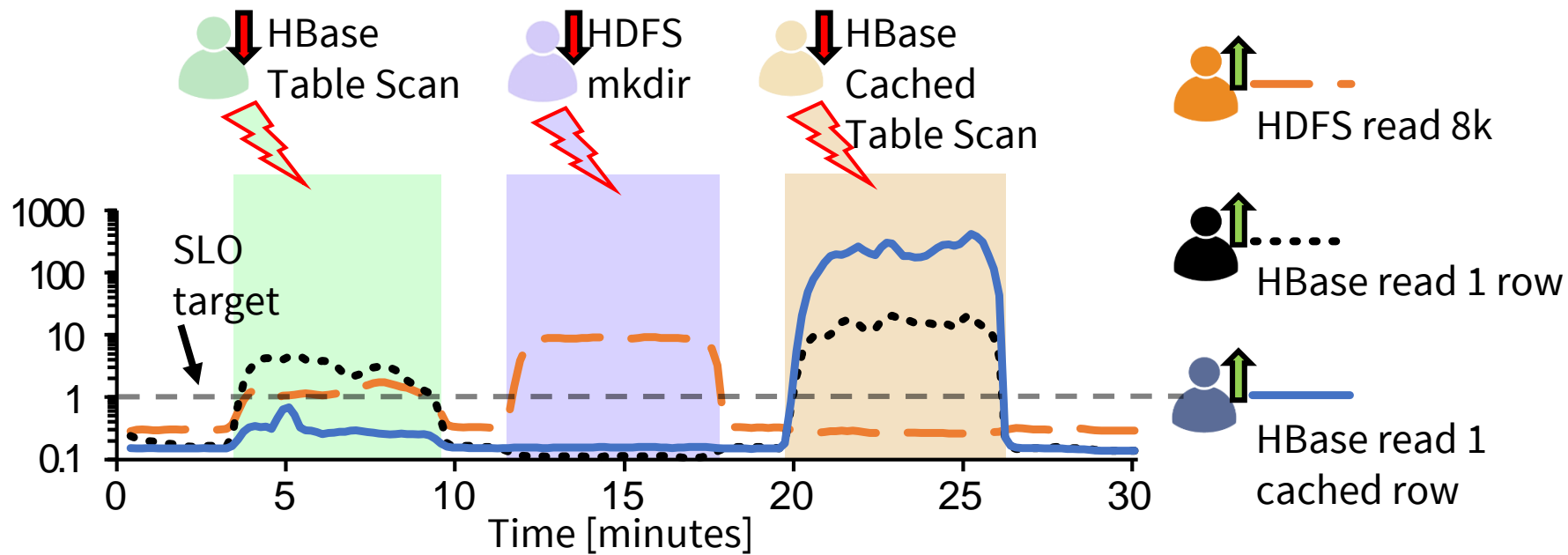
 HBase read 1  
cached row



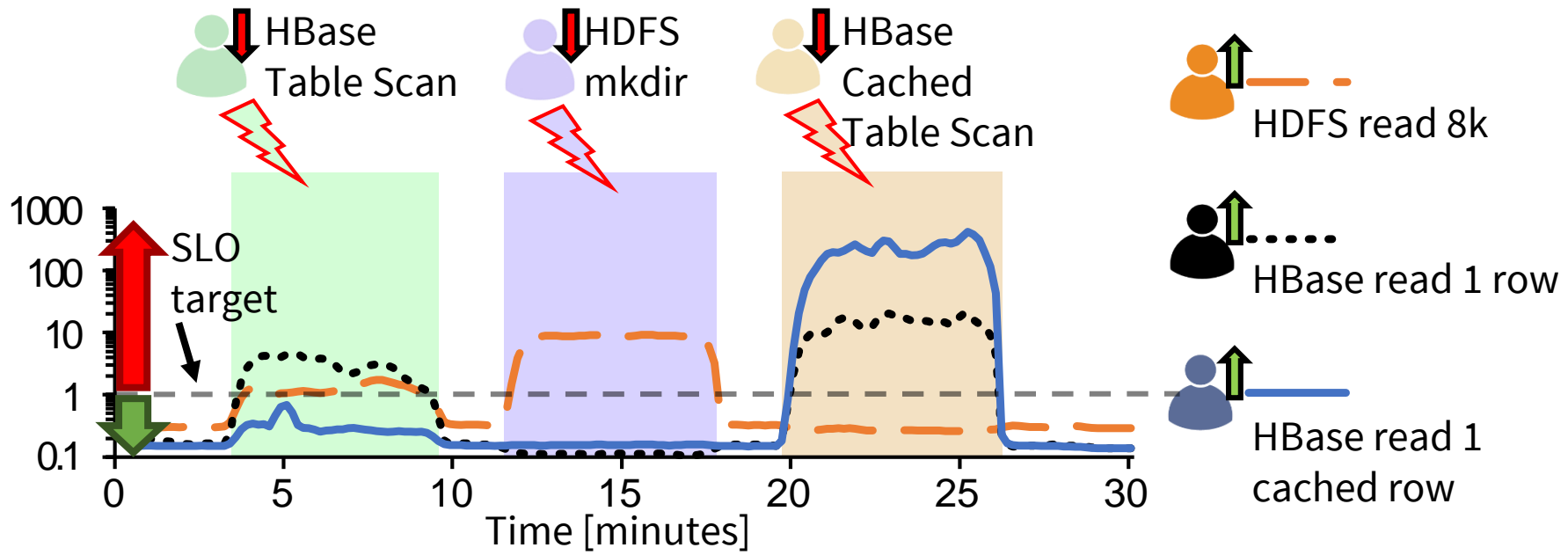
SLO-Normalized Latency



SLO-Normalized Latency

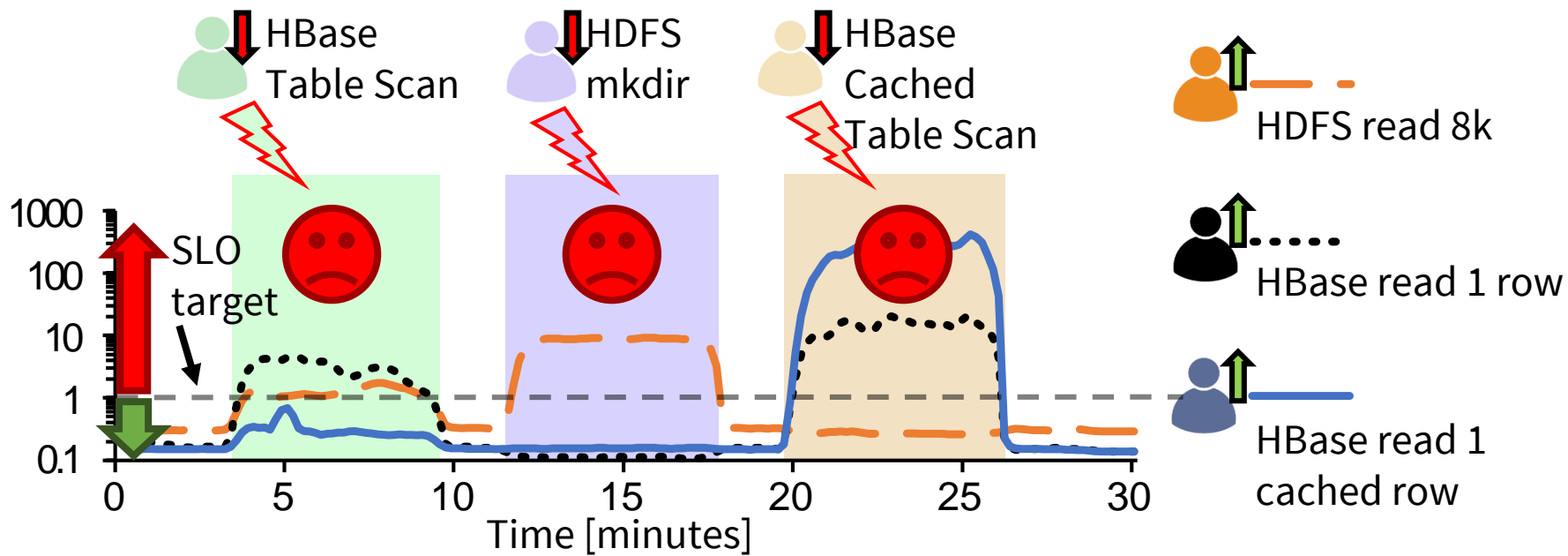


SLO-Normalized Latency





SLO-Normalized Latency



SLO-Normalized Latency

HBase  
Table Scan

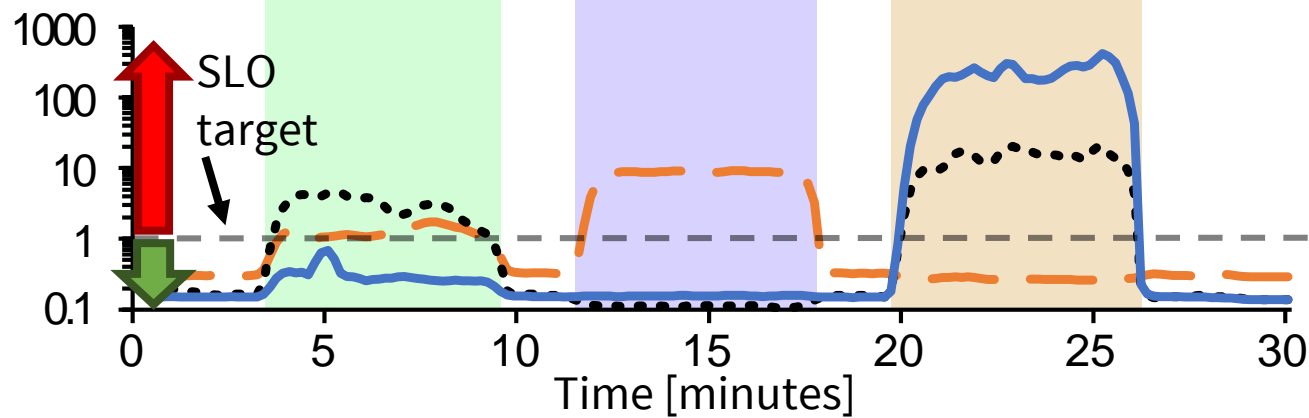
HDFS  
mkdir

HBase  
Cached  
Table Scan

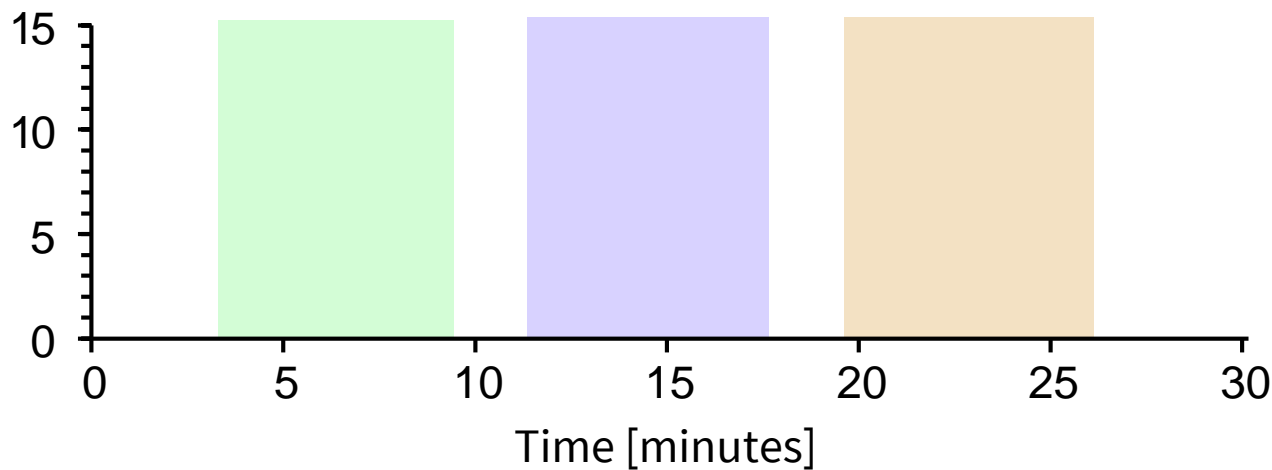
HDFS read 8k

HBase read 1 row

HBase read 1  
cached row



Slowdown



SLO-Normalized Latency

Slowdown

HBase  
Table Scan

HDFS  
mkdir

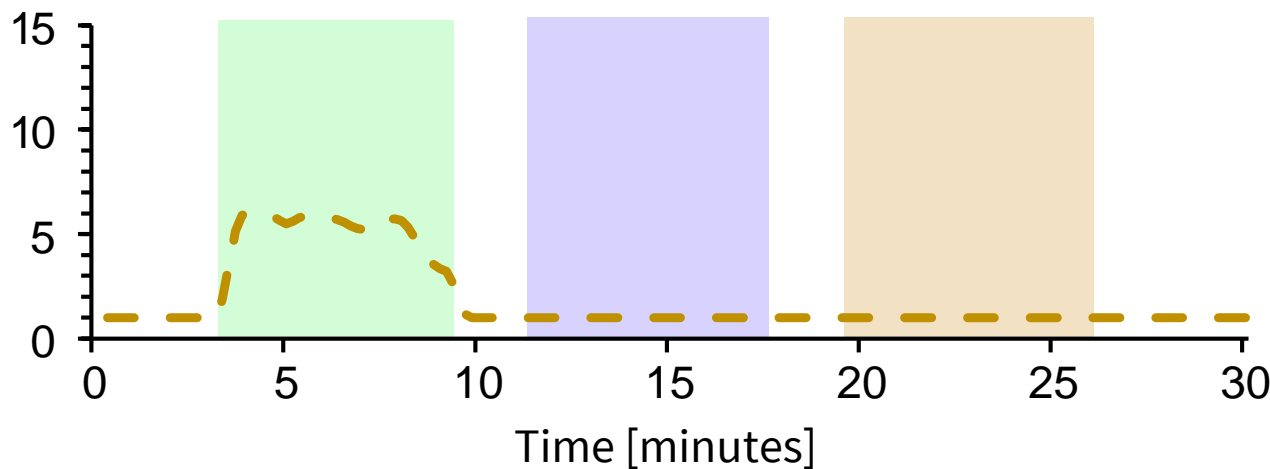
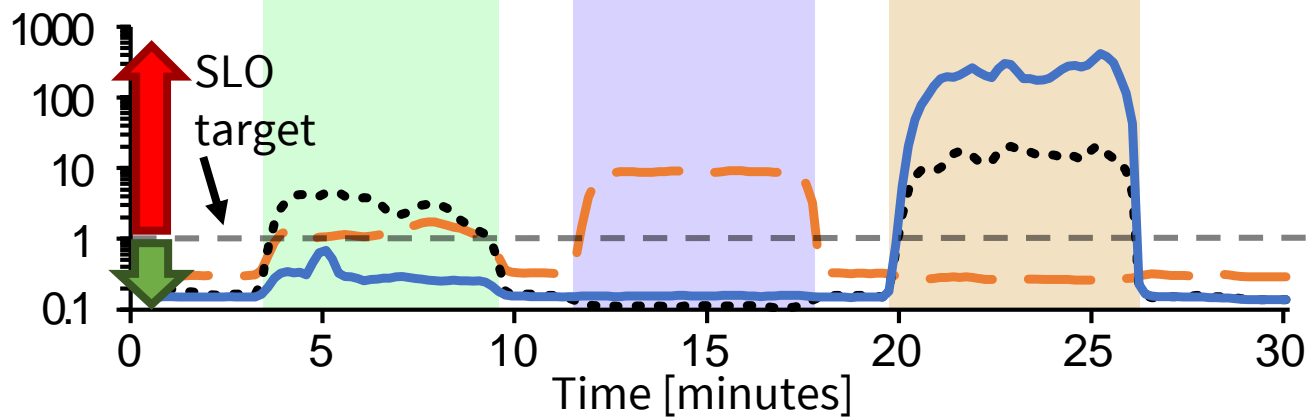
HBase  
Cached  
Table Scan

HDFS read 8k

HBase read 1 row


HBase read 1  
cached row

Disk




SLO-Normalized Latency


Slowdown


 HBase  
Table Scan

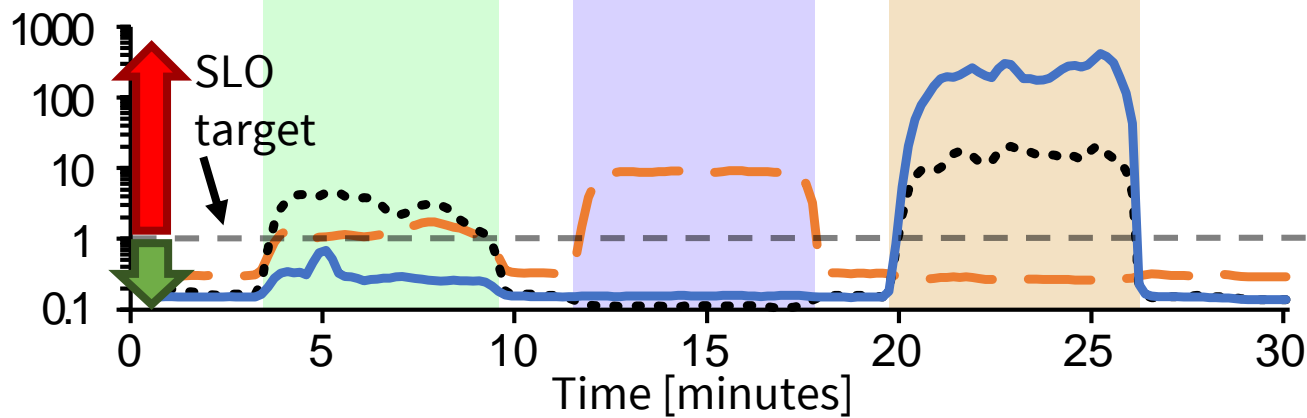
 HDFS  
mkdir


 HBase  
Cached  
Table Scan

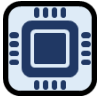
 HDFS read 8k

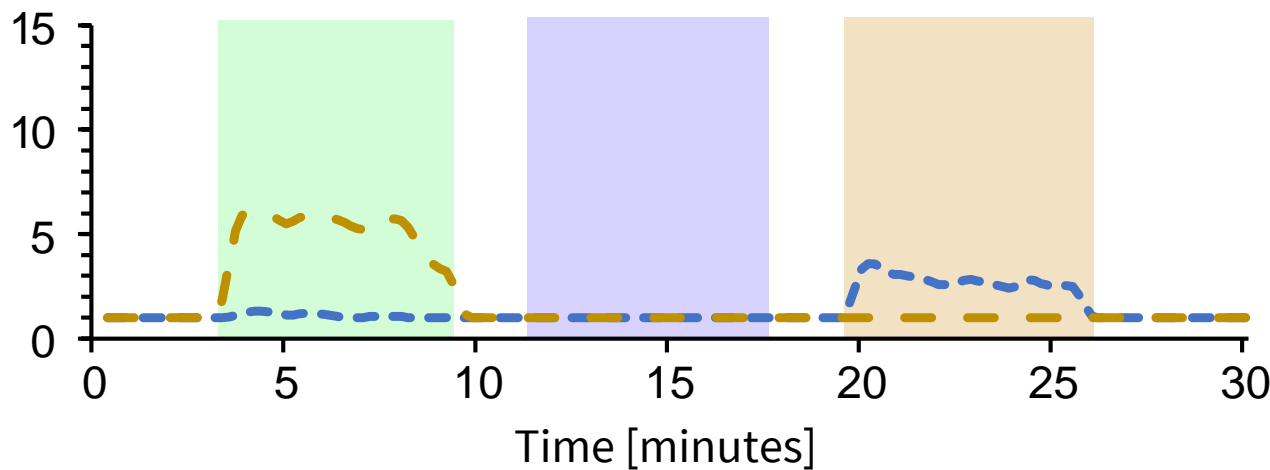
 HBase read 1 row

 HBase read 1  
cached row




 Disk

 CPU



SLO-Normalized Latency


Slowdown


 HBase  
Table Scan

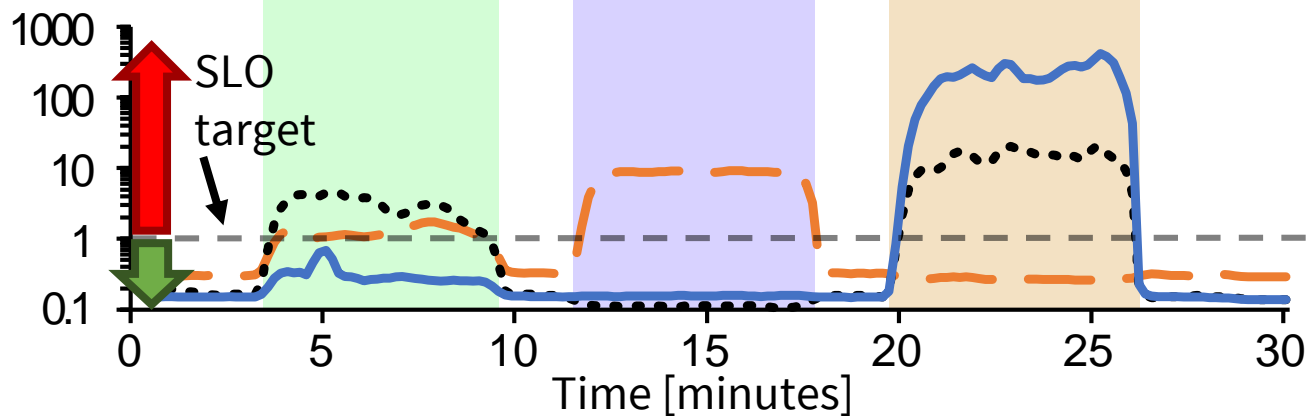
 HDFS  
mkdir

 HBase  
Cached  
Table Scan


 HDFS read 8k

 HBase read 1 row

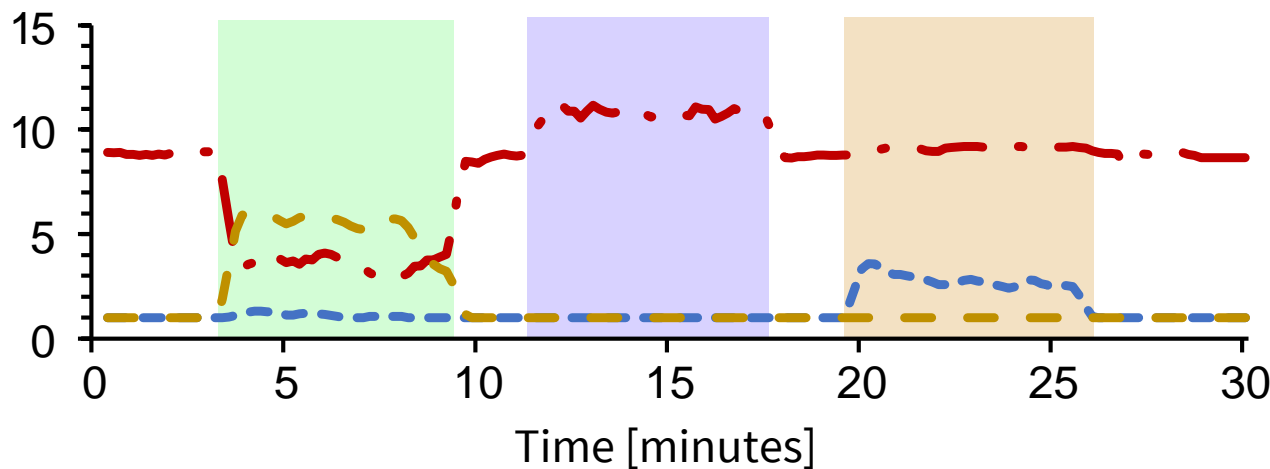
 HBase read 1  
cached row



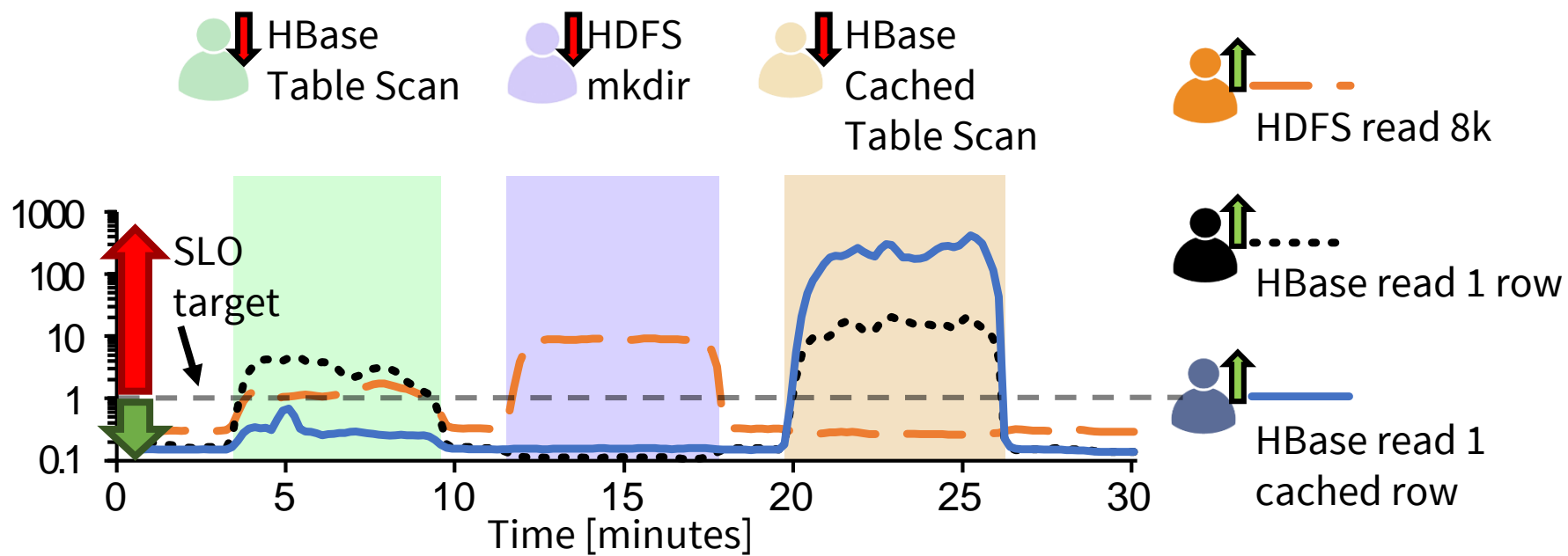
 Disk

 CPU

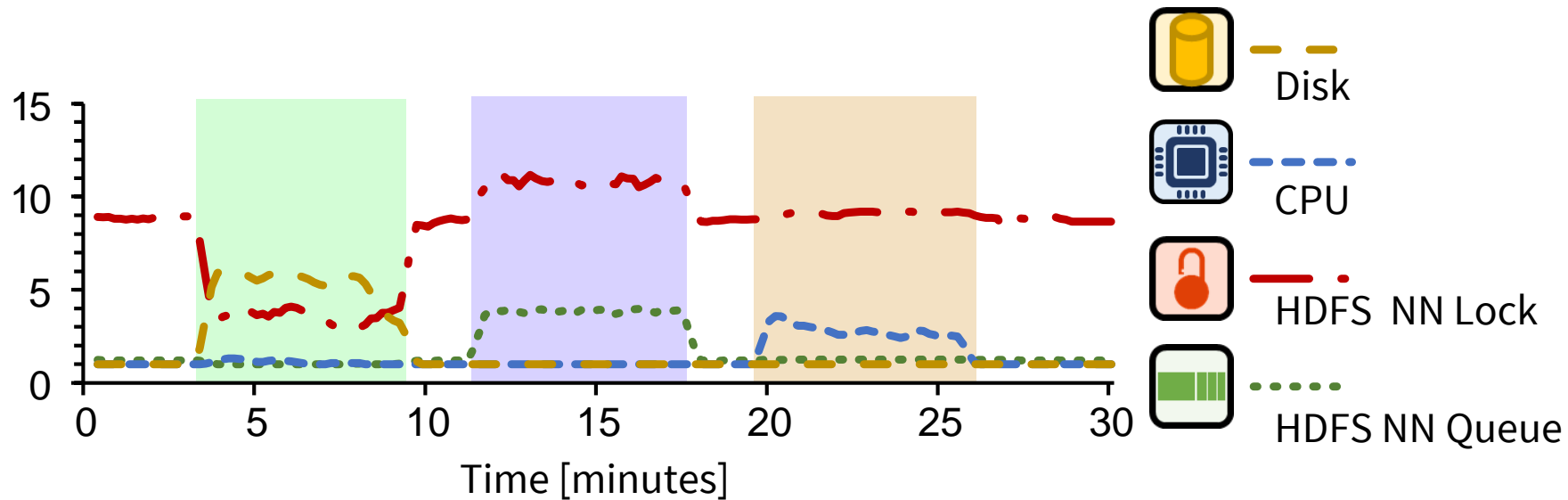
 HDFS NN Lock



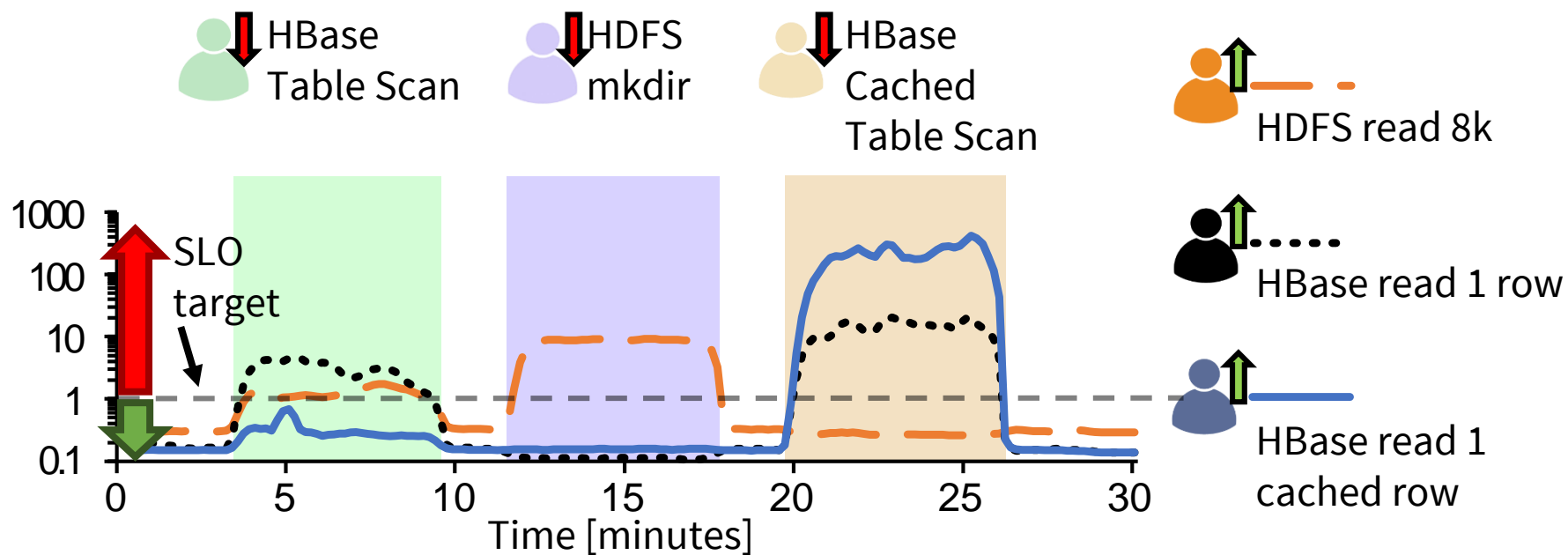
SLO-Normalized Latency



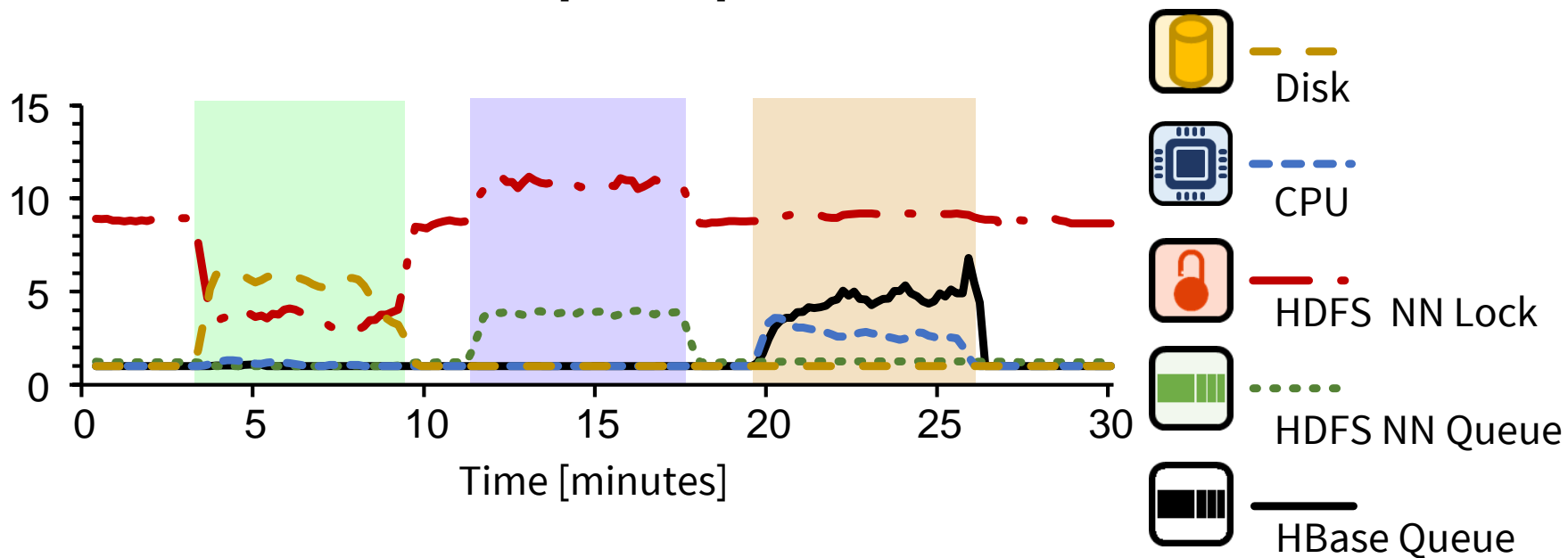
Slowdown



SLO-Normalized Latency



Slowdown



SLO-Normalized Latency

SLO-Normalized Latency

↓ HBase  
Table Scan

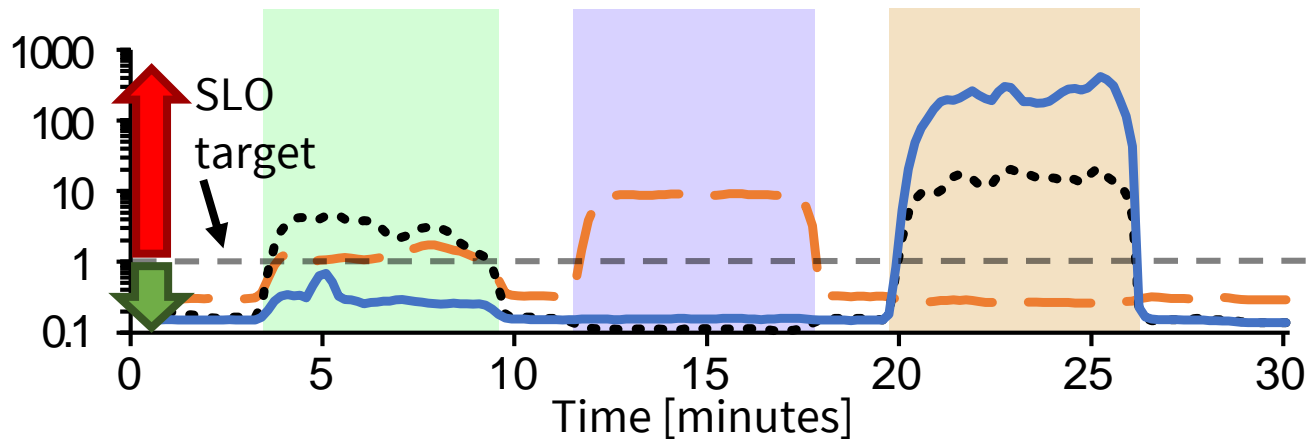
↓ HDFS  
mkdir

↓ HBase  
Cached  
Table Scan

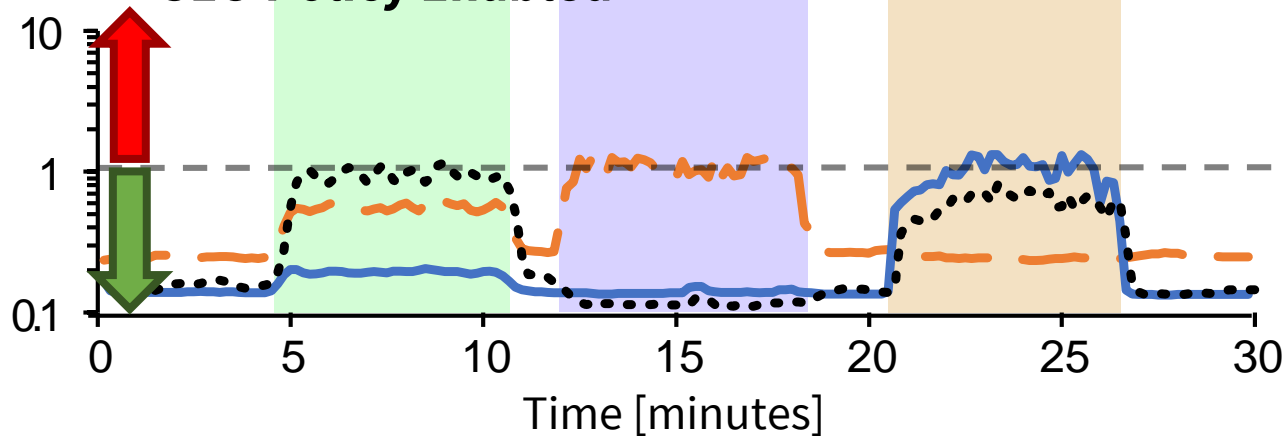
↑ HDFS read 8k

↑ HBase read 1 row

↑ HBase read 1  
cached row




**+ SLO Policy Enabled**






SLO-Normalized Latency


 HBase  
Table Scan

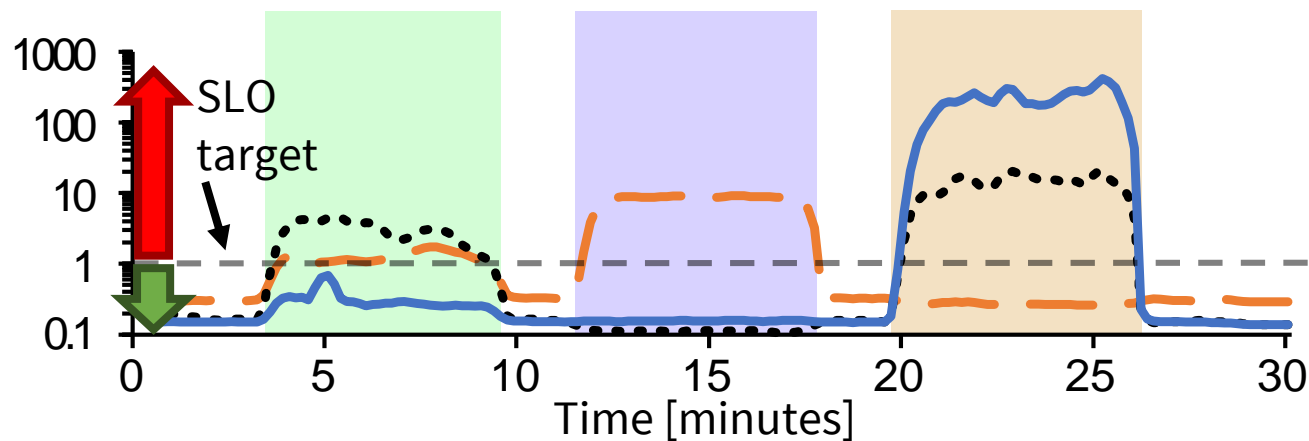
 HDFS  
mkdir

 HBase  
Cached  
Table Scan

 HDFS read 8k

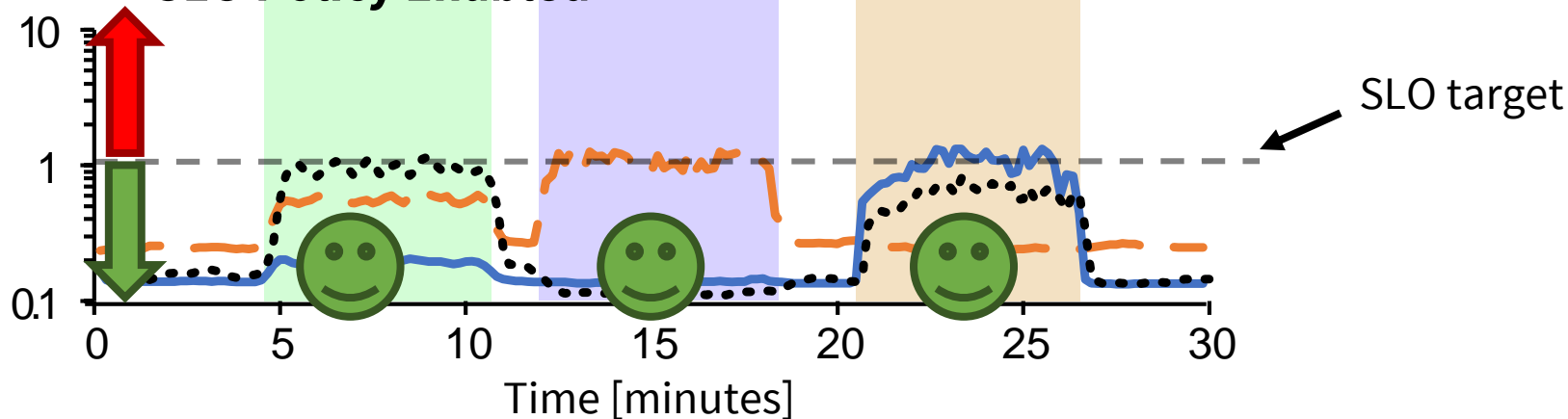
 HBase read 1 row

 HBase read 1  
cached row

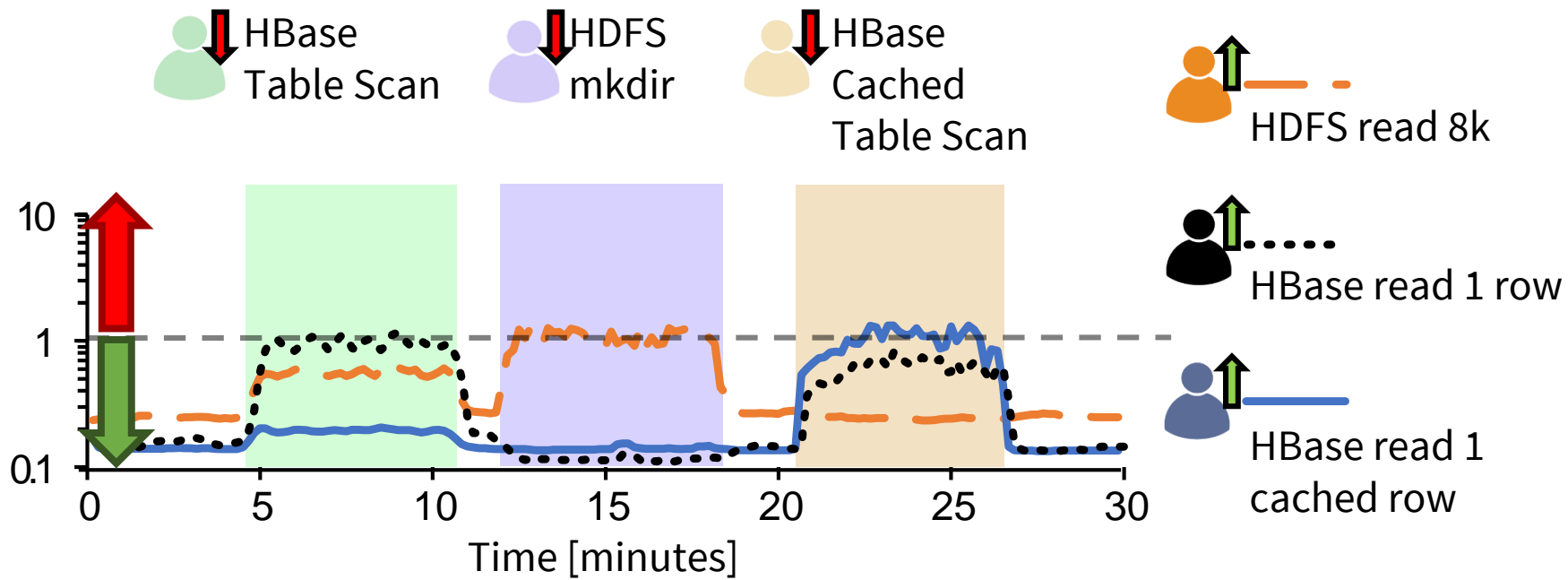


SLO-Normalized Latency

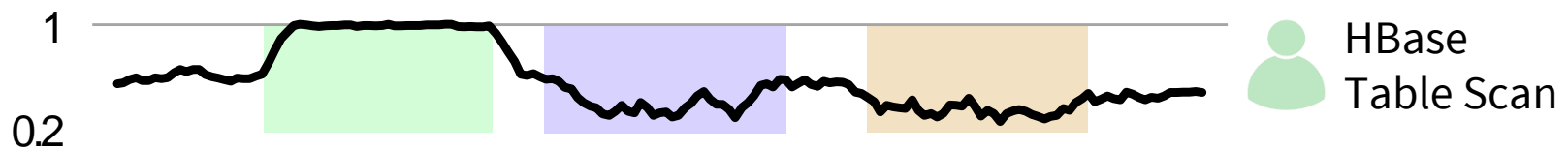
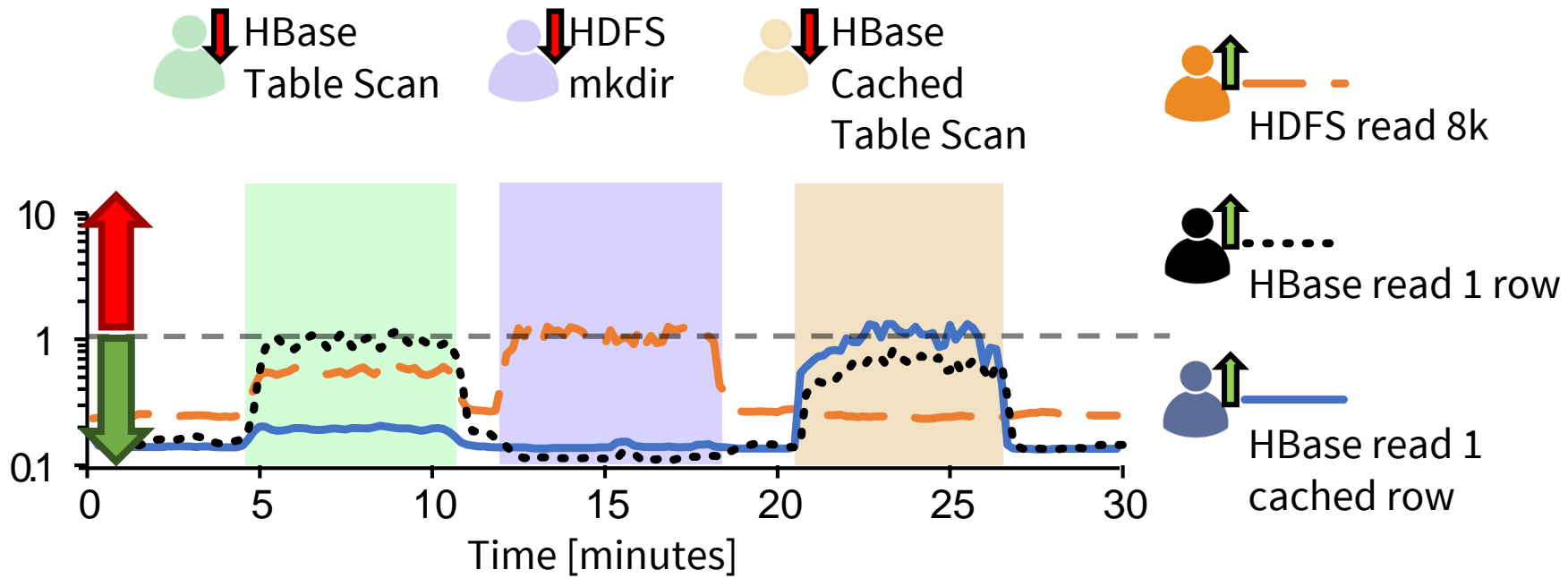
**+ SLO Policy Enabled**



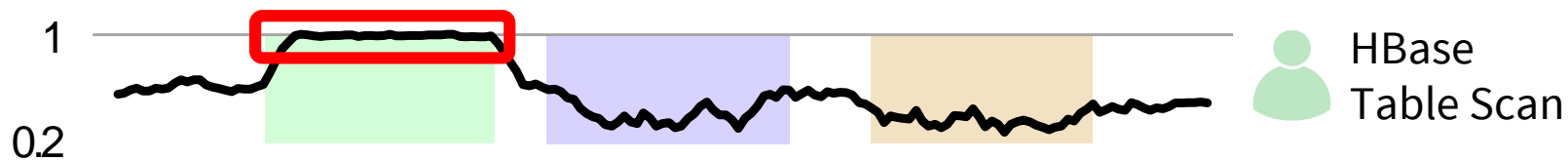
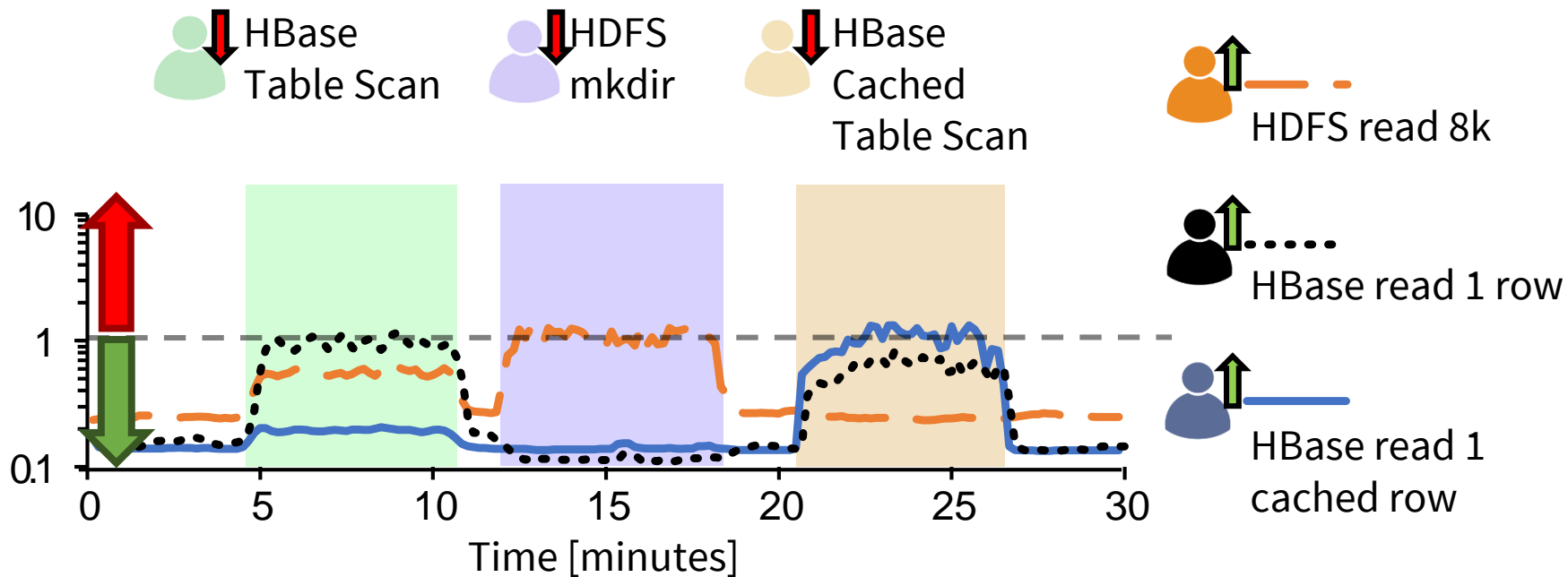
SLO-Normalized Latency



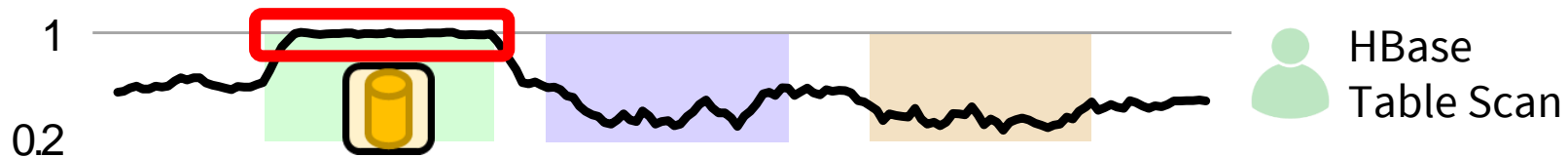
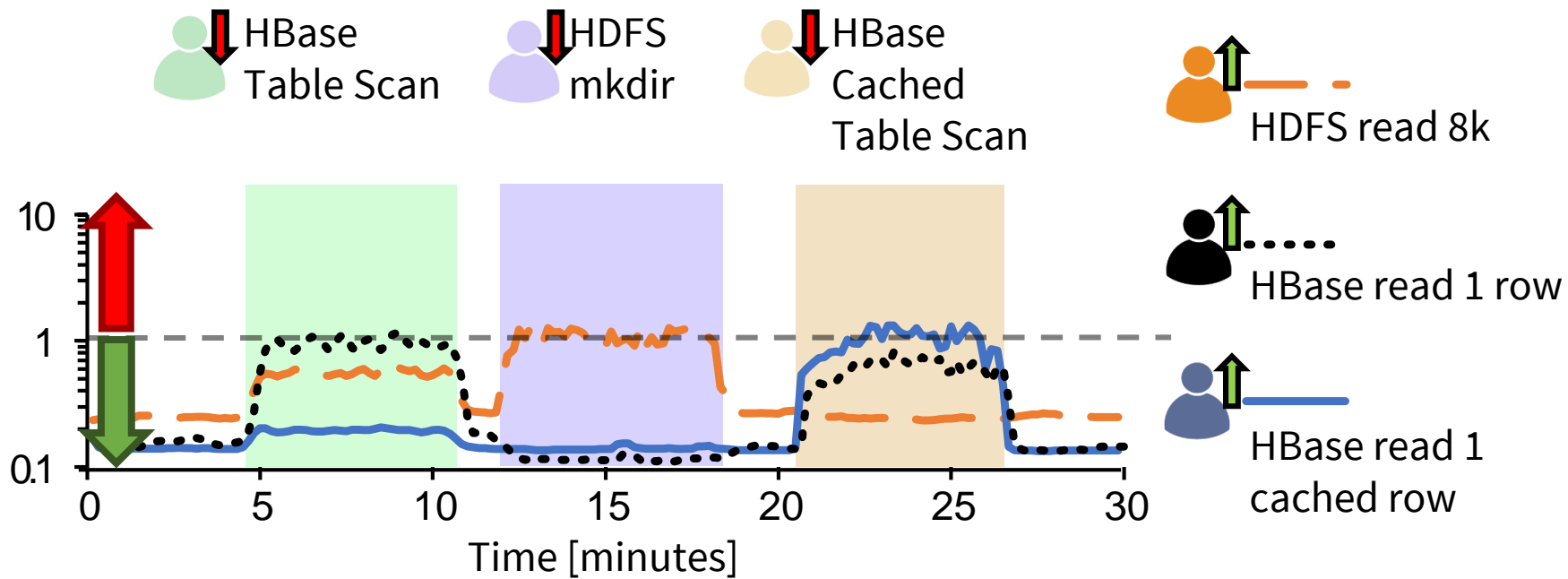
SLO-Normalized Latency



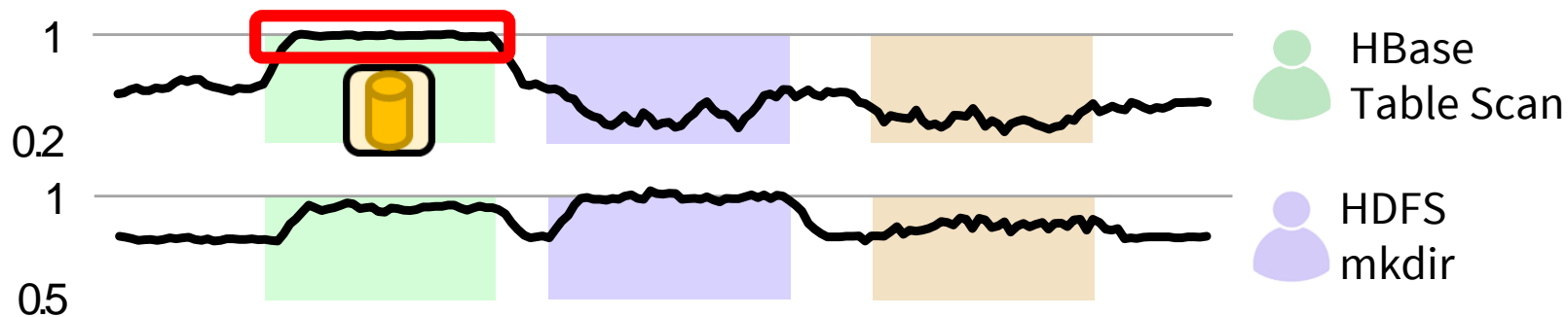
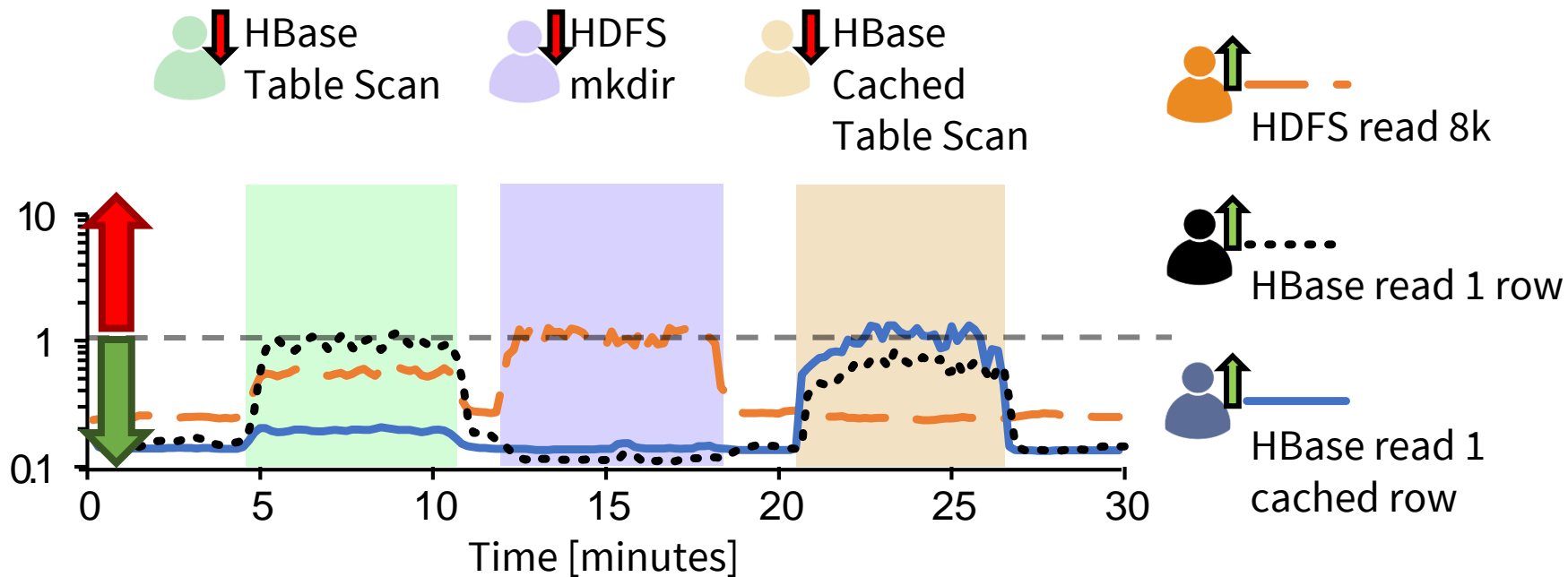
SLO-Normalized Latency



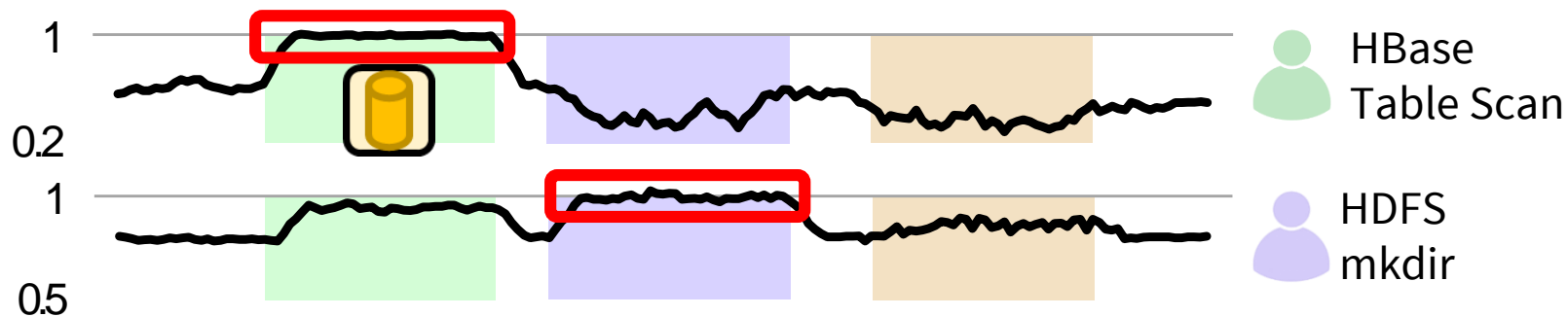
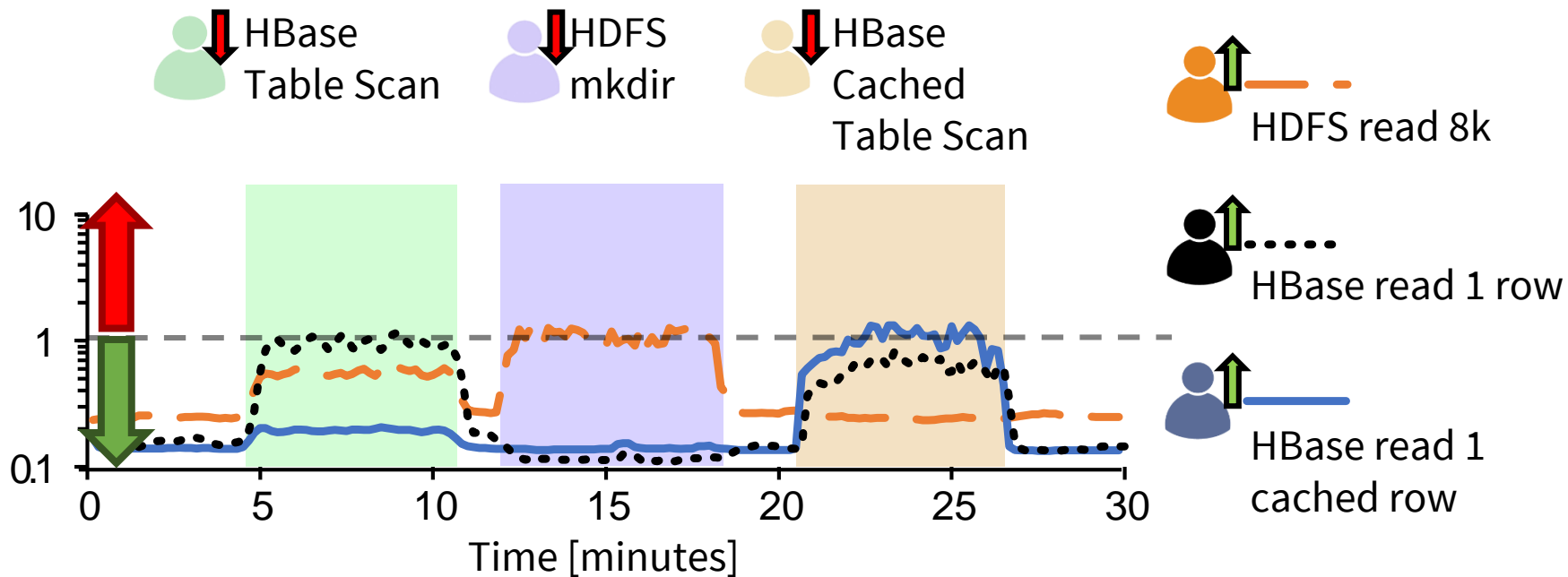
SLO-Normalized Latency



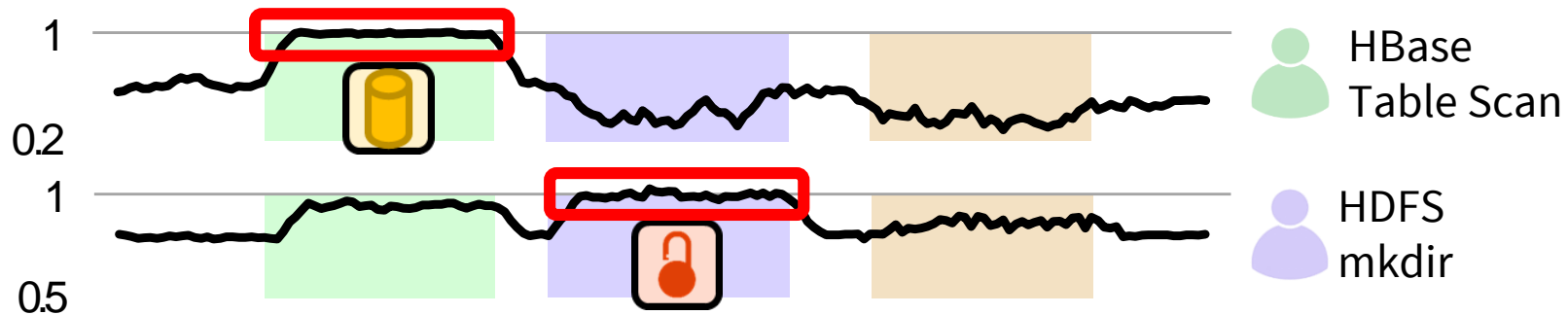
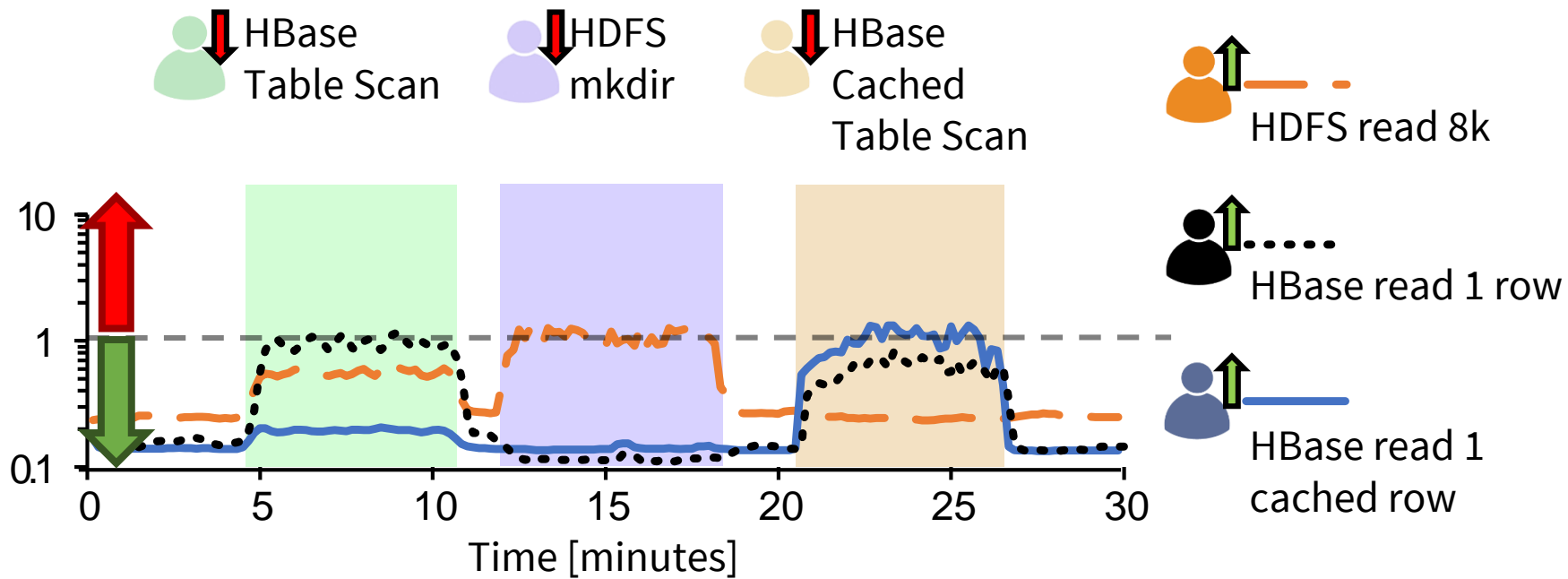
SLO-Normalized Latency



SLO-Normalized Latency

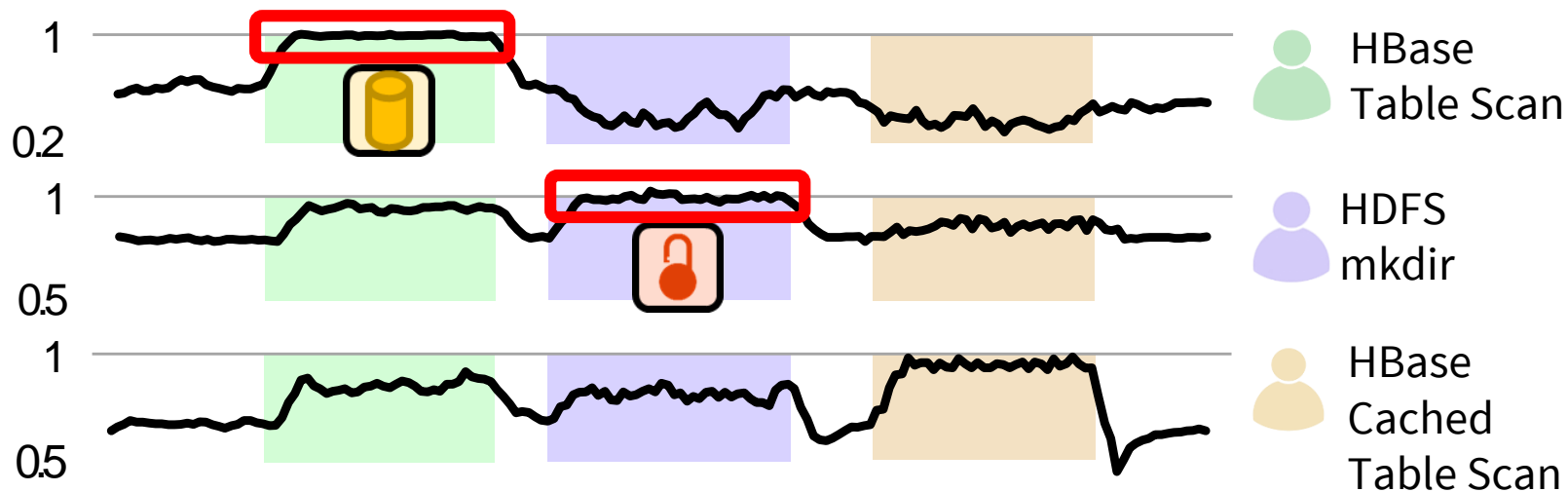
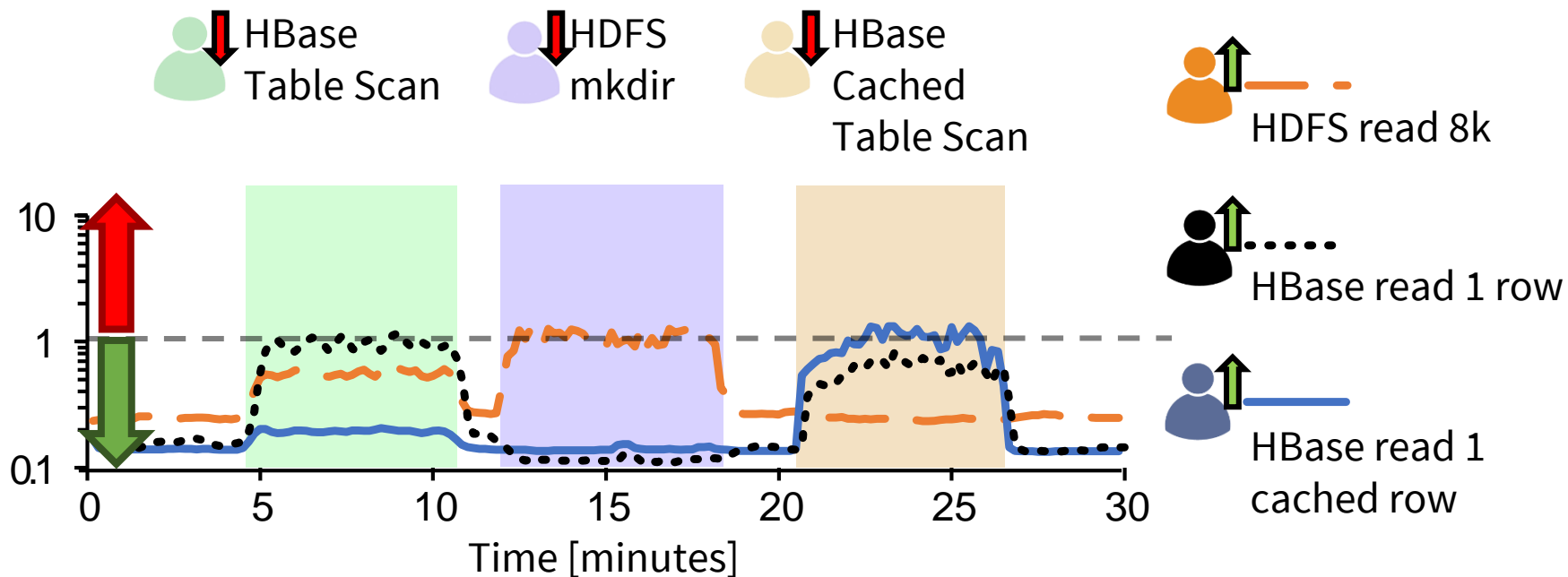


SLO-Normalized Latency

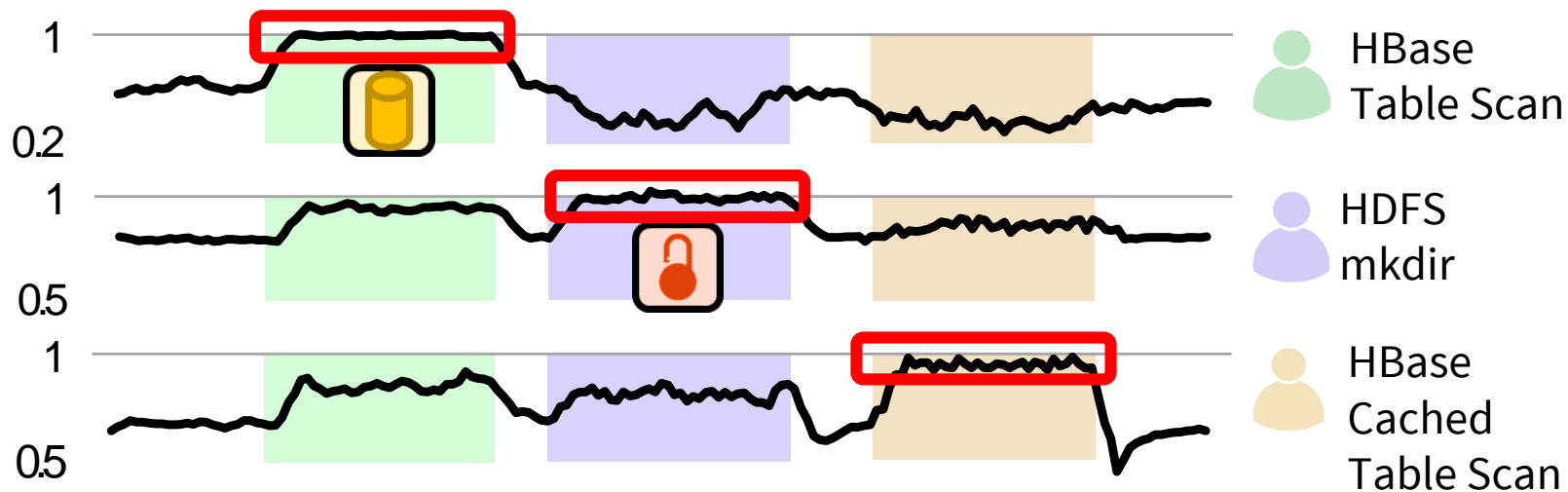
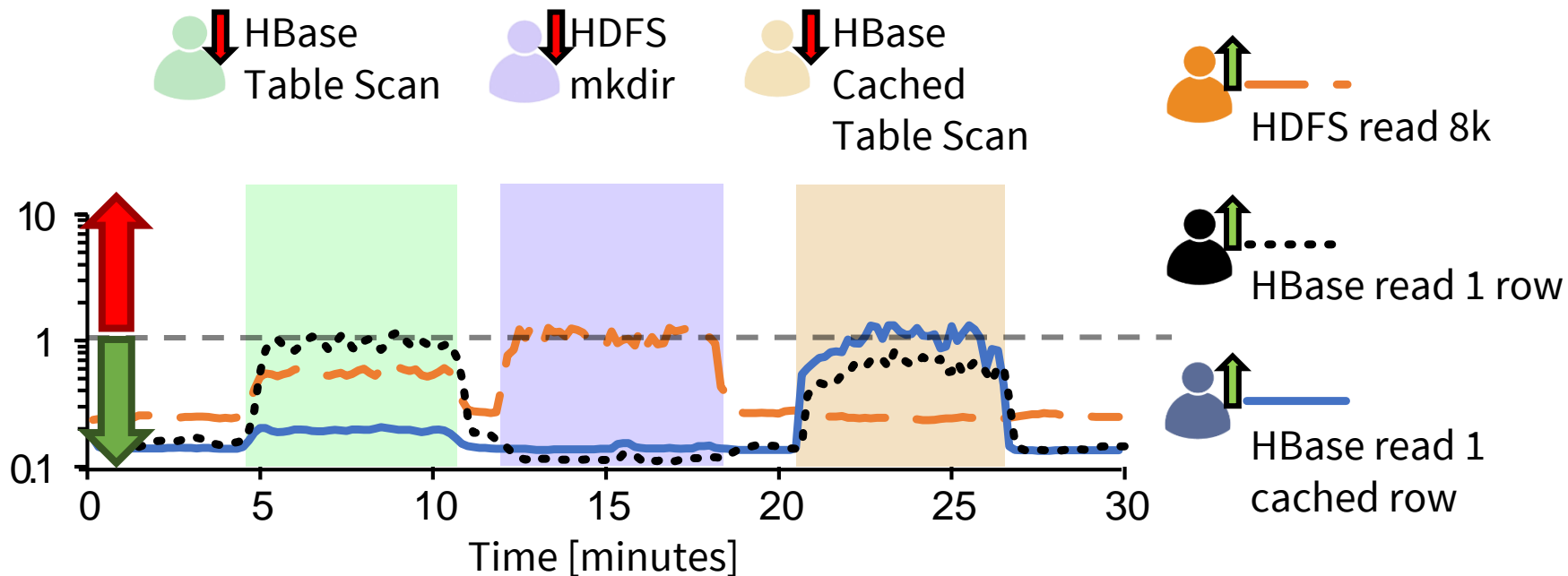




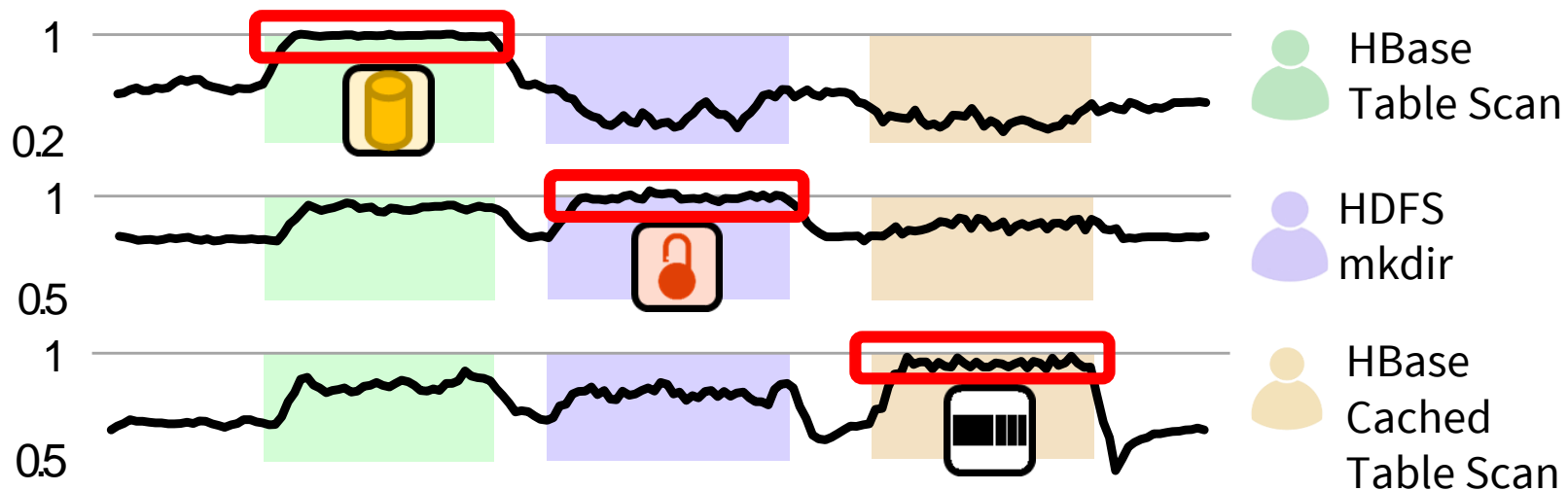
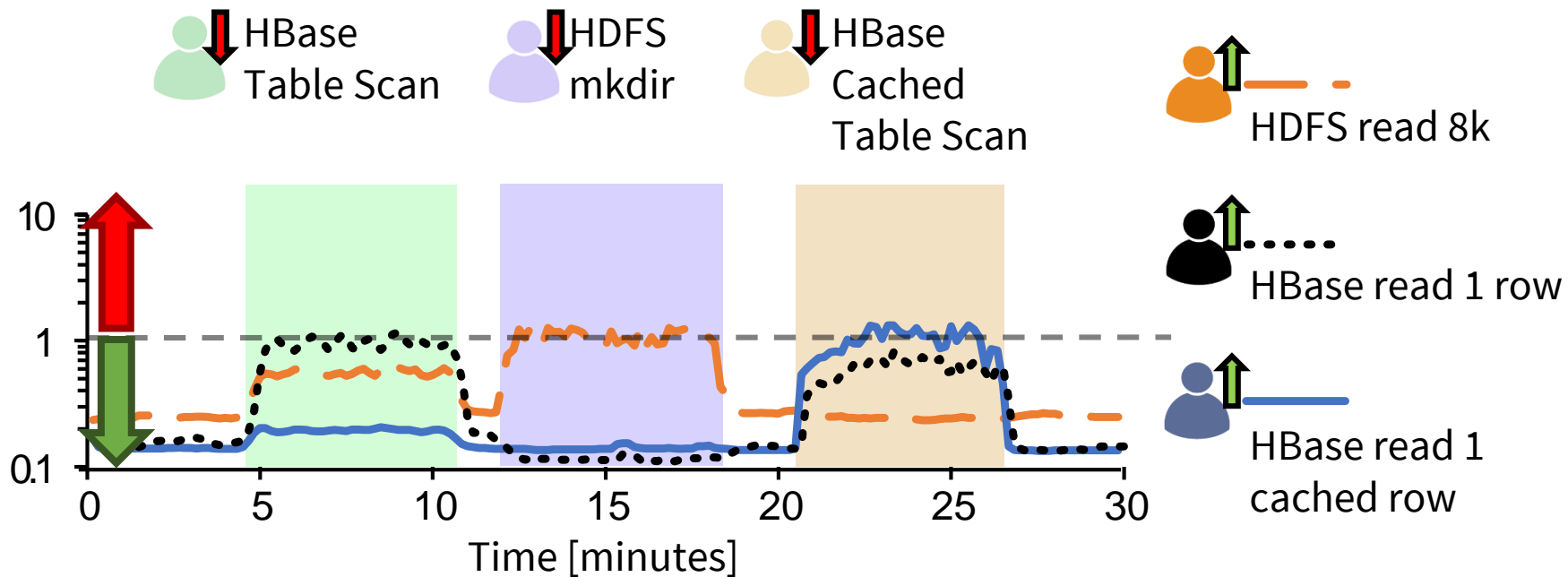
SLO-Normalized Latency



SLO-Normalized Latency

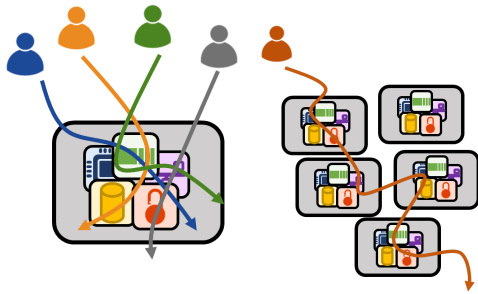


SLO-Normalized Latency

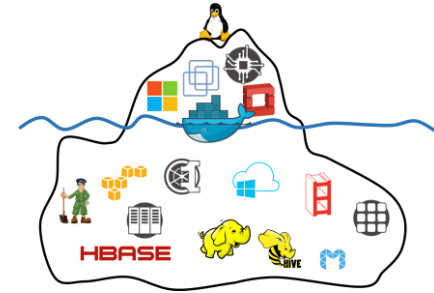


# *Retra* **Conclusion**

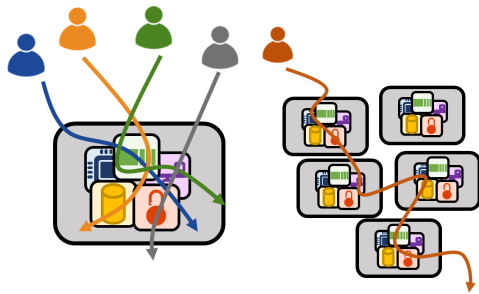
# Retro Conclusion



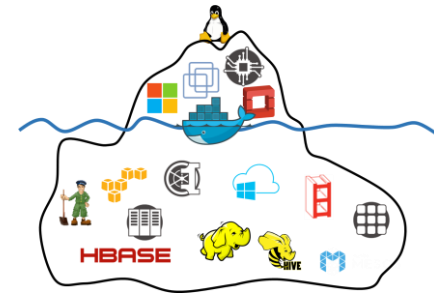
Resource management for *shared*  
distributed systems



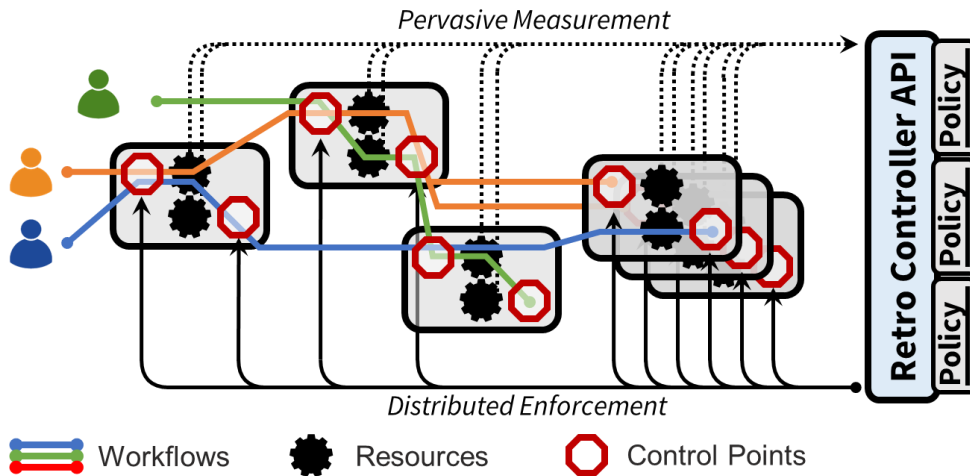
# Retro Conclusion



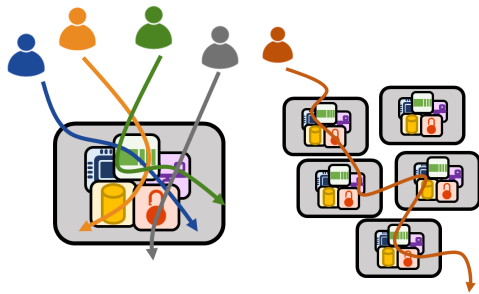
Resource management for *shared* distributed systems



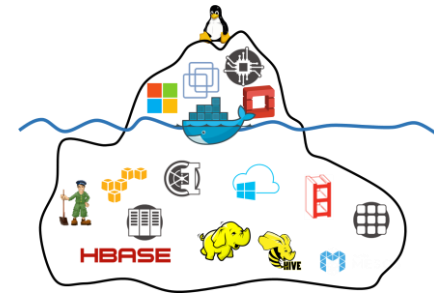
*Centralized* resource management



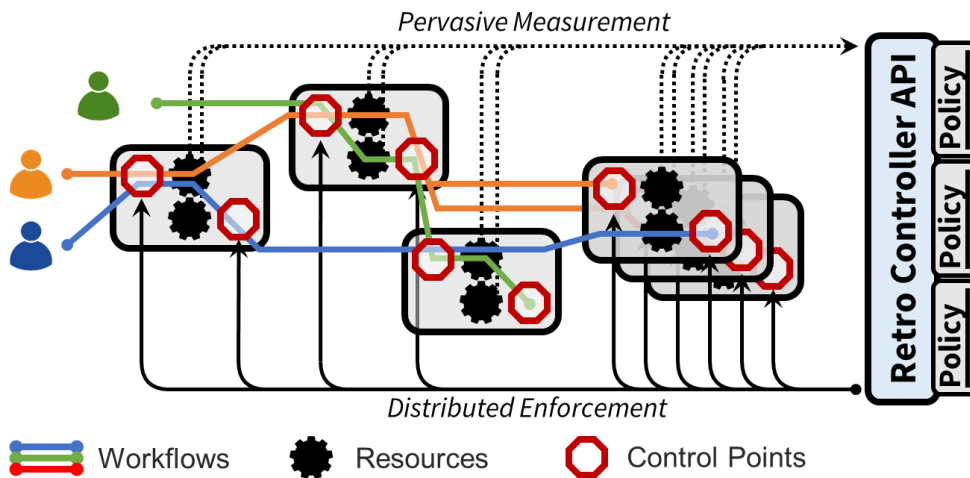
# Retro Conclusion



Resource management for *shared* distributed systems

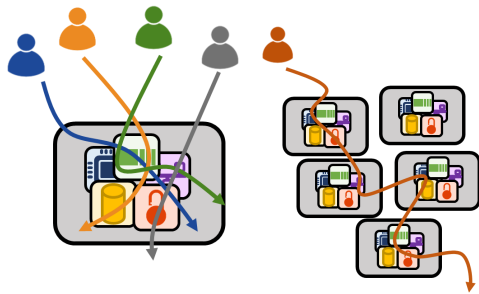


*Centralized* resource management

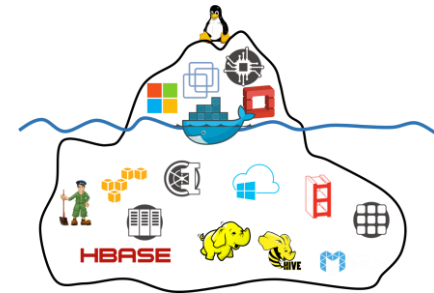


*Comprehensive:* resources, processes, tenants, background tasks

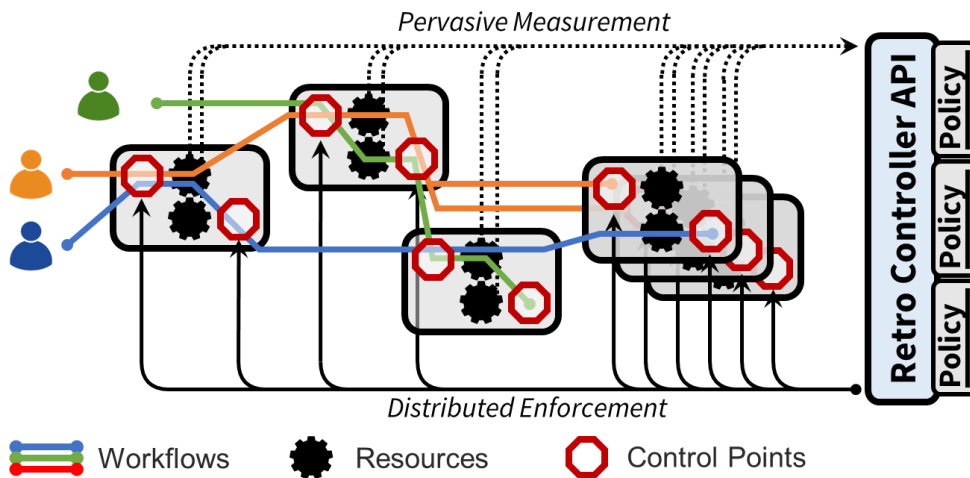
# Retro Conclusion



Resource management for *shared* distributed systems



*Centralized* resource management



*Comprehensive*: resources, processes, tenants, background tasks

Abstractions for writing *concise, general-purpose* policies:

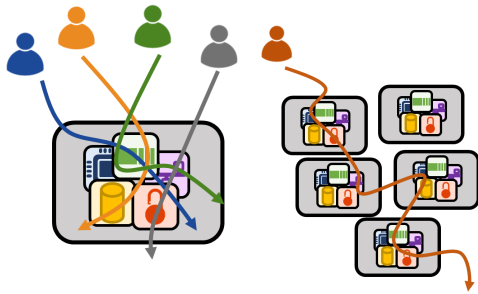
Workflows

Resources (slowdown, load)

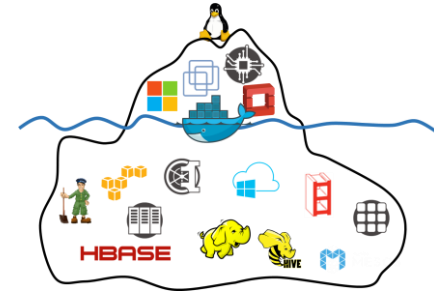
Control points



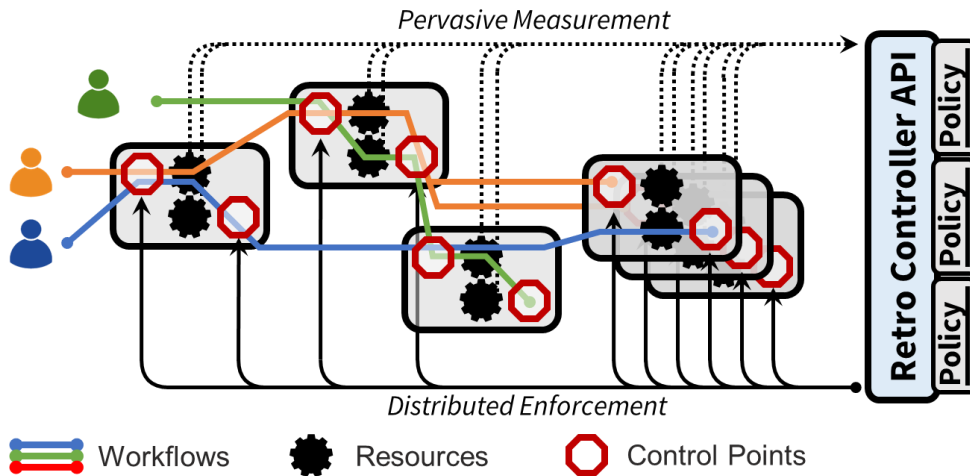
# Retro Conclusion



Resource management for *shared* distributed systems



*Centralized* resource management



*Comprehensive*: resources, processes, tenants, background tasks

Abstractions for writing *concise, general-purpose* policies:

Workflows

Resources (slowdown, load)

Control points

Microsoft

Research



BROWN UNIVERSITY

<http://cs.brown.edu/~jcmace>