

MittOS

Supporting Millisecond Tail Tolerance with Fast
Rejecting SLO-Aware OS Interface

Mingzhe Hao, Huaicheng Li, Michael Hao Tong,
Chrisma Pakha, Riza O. Suminto, Cesar A. Stuardo,
Andrew A. Chien, and Haryadi S. Gunawi



THE UNIVERSITY OF
CHICAGO





Millisecond Matters!

**AMAZON: “EVERY 100MS OF
LATENCY COSTS 1% IN SALES”**

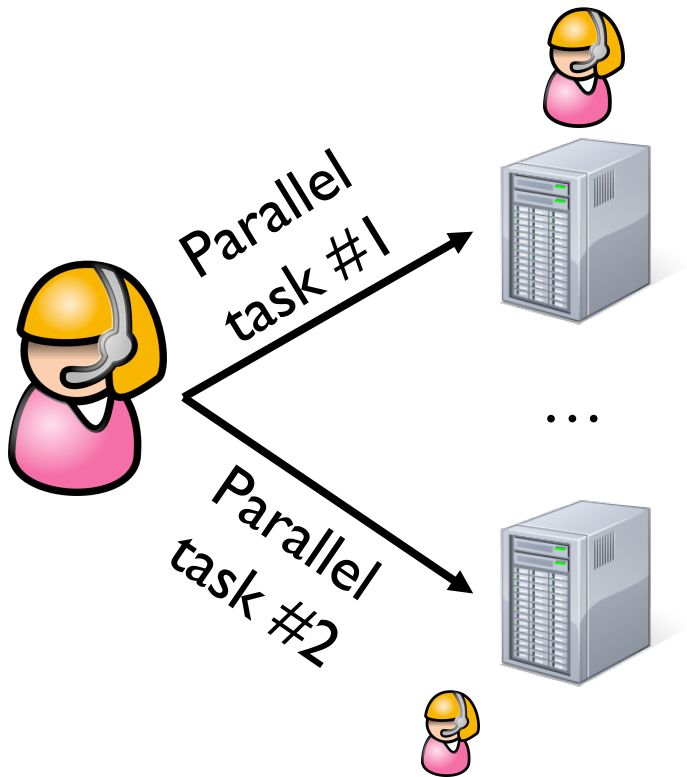
**TABB GROUP: “BROKER COULD LOSE AS MUCH
AS \$4 MILLION IN REVENUES PER MILLISECOND
IF ITS ELECTRONIC TRADING PLATFORM WAS
ONLY 5MS BEHIND THE COMPETITION”**

**GOOGLE: “EXTRA 500MS IN SEARCH PAGE
GENERATION TIME DROPPED TRAFFIC BY
20%”**

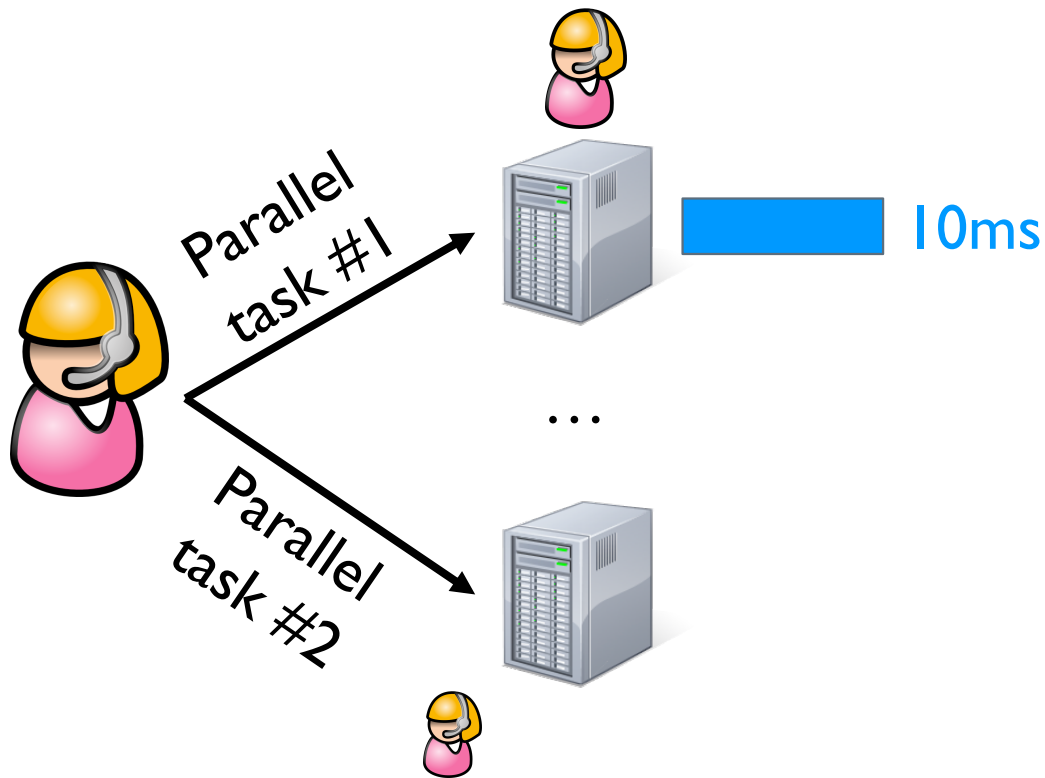


Millisecond Tail Latency

Millisecond Tail Latency

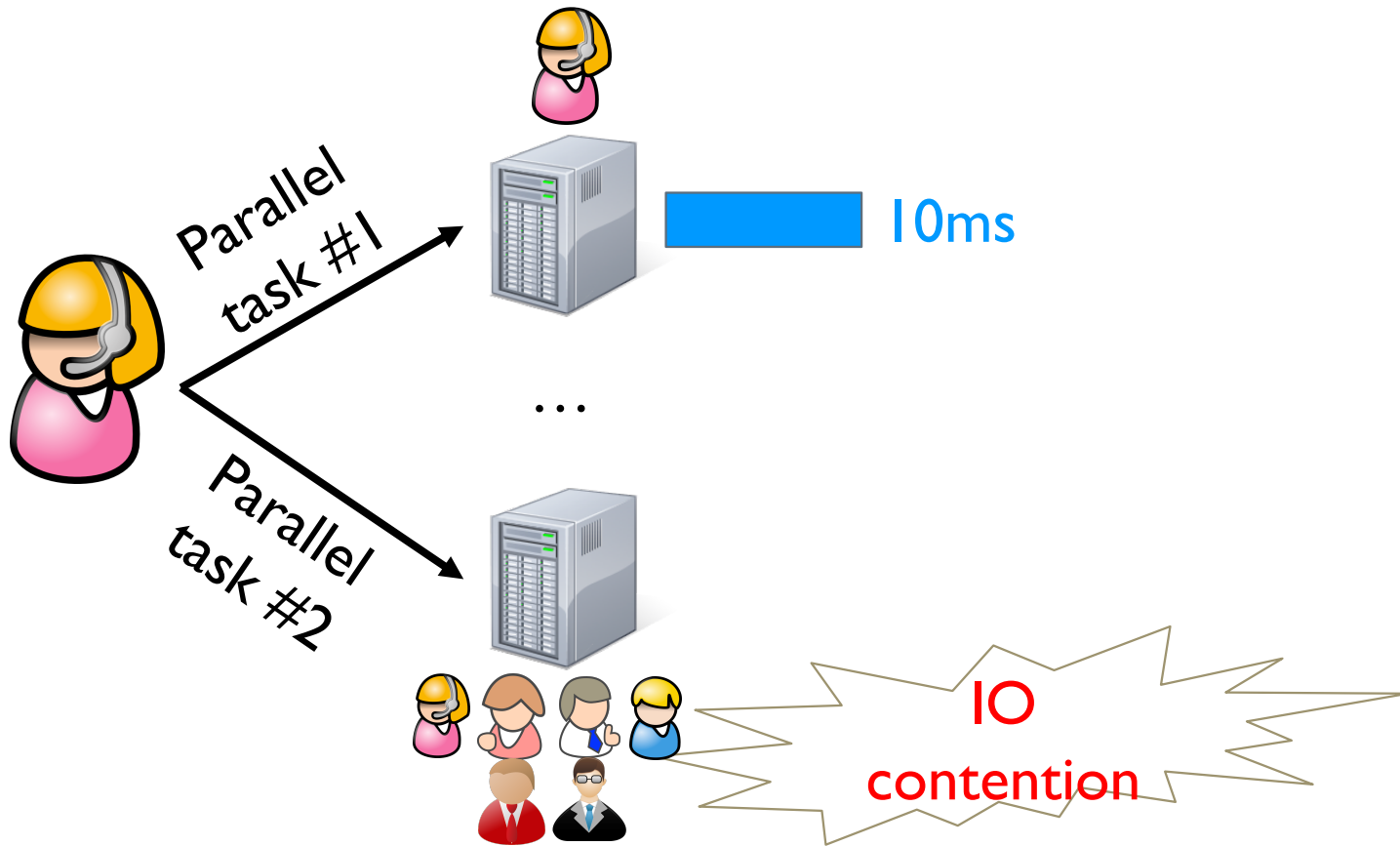


Millisecond Tail Latency



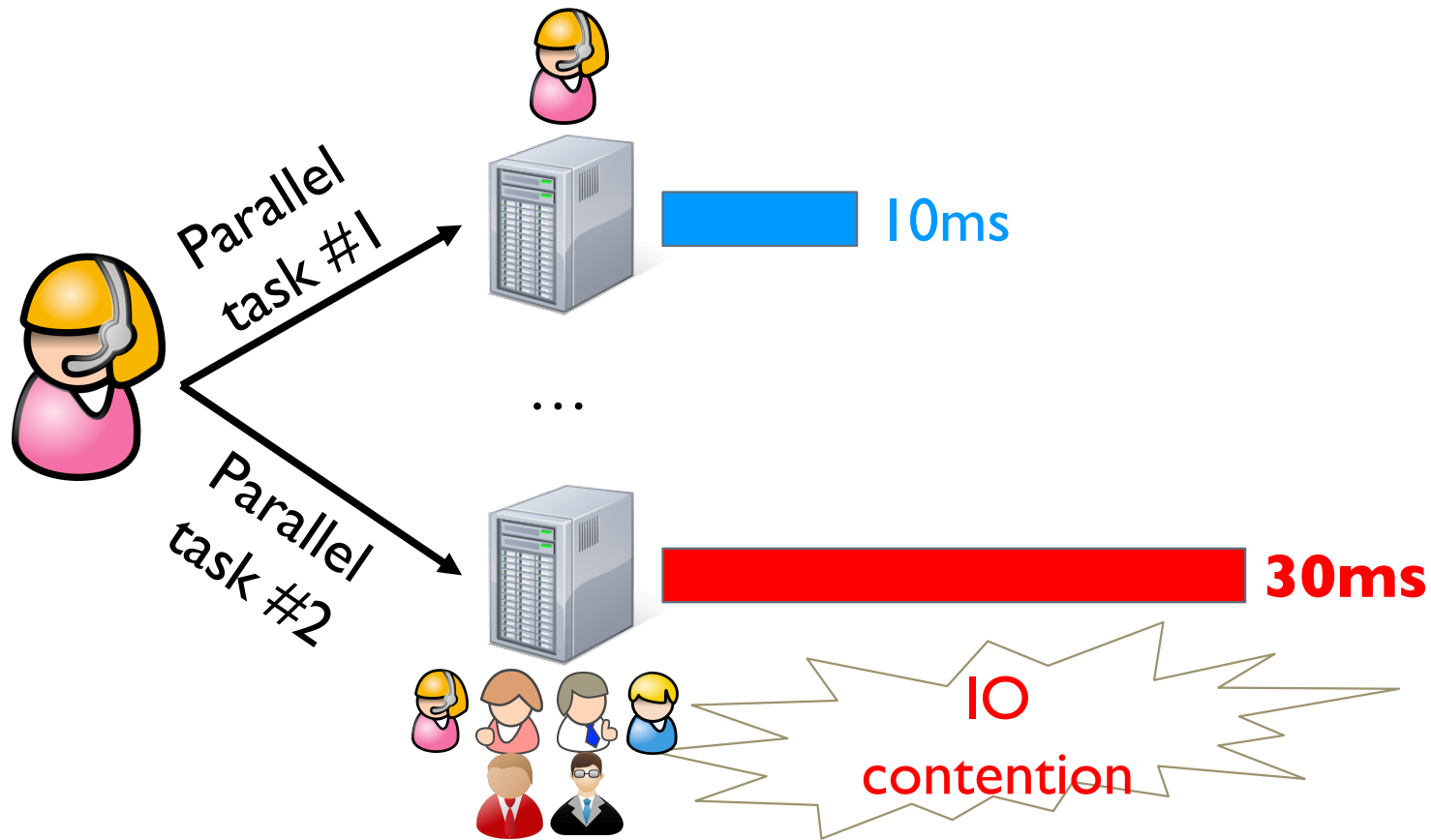


Millisecond Tail Latency

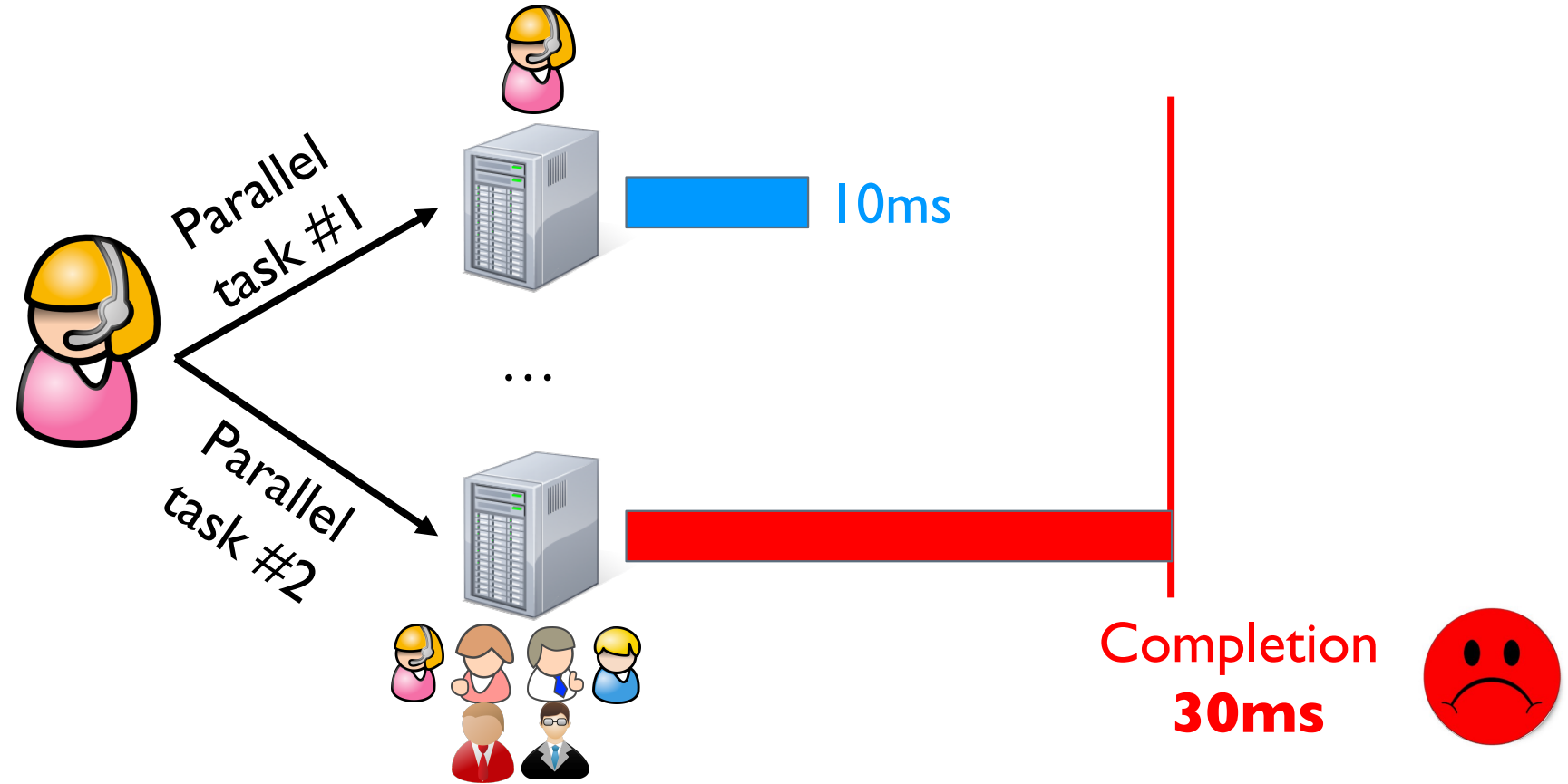




Millisecond Tail Latency



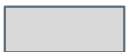
Millisecond Tail Latency





Current Tail-Tolerant Mechanisms

1. Speculation





Current Tail-Tolerant Mechanisms

1. Speculation



Wait
→





Current Tail-Tolerant Mechanisms

1. Speculation



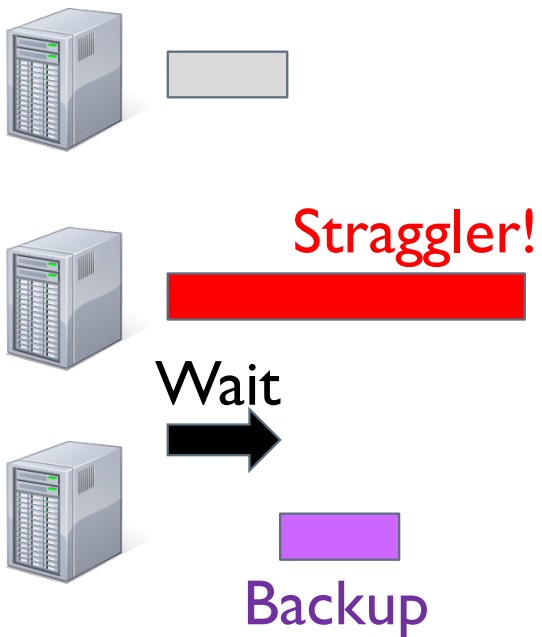
Straggler!





Current Tail-Tolerant Mechanisms

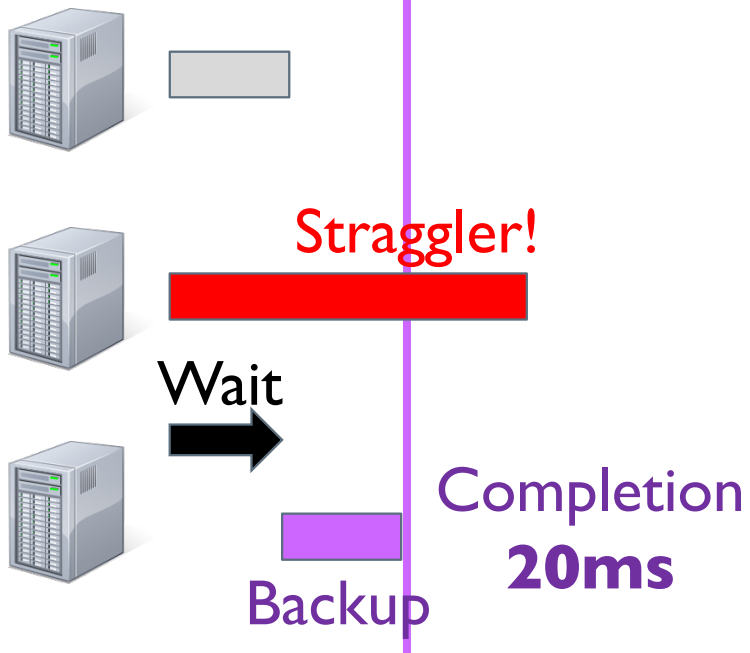
1. Speculation





Current Tail-Tolerant Mechanisms

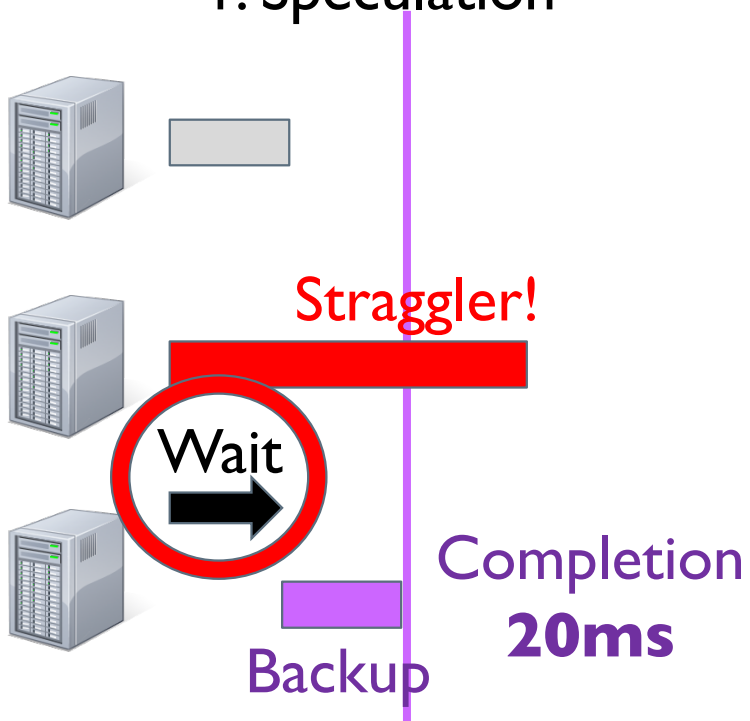
1. Speculation





Current Tail-Tolerant Mechanisms

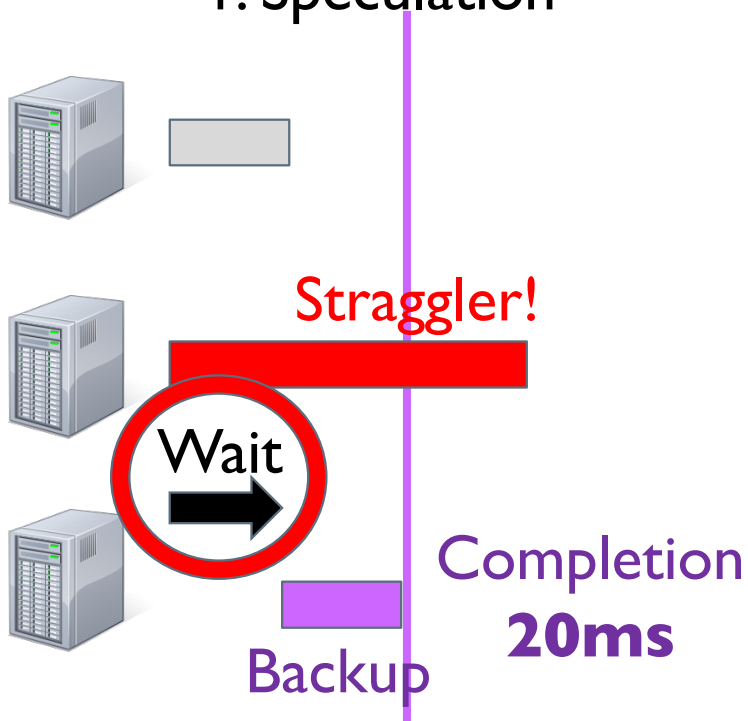
1. Speculation





Current Tail-Tolerant Mechanisms

1. Speculation

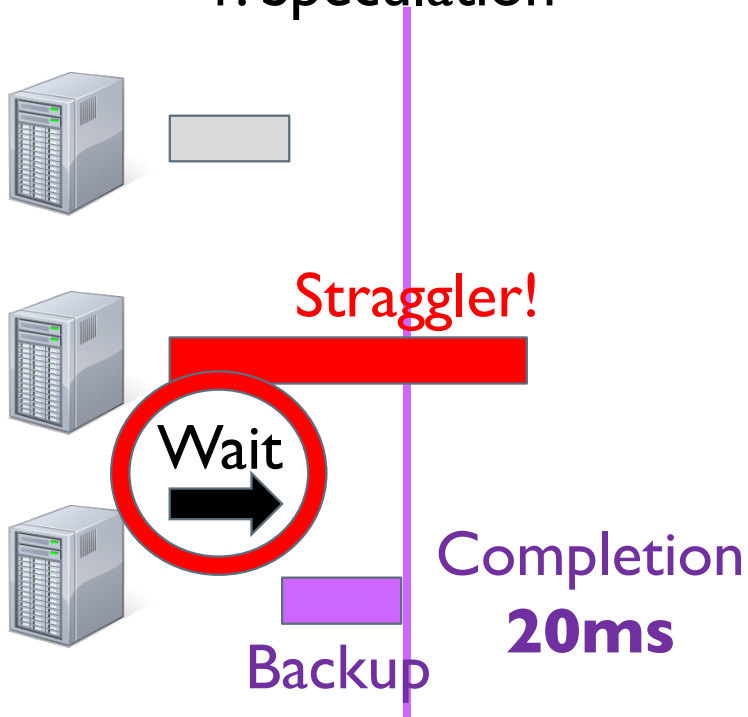


2. Cloning

- Introduces **2x** workload

Current Tail-Tolerant Mechanisms

1. Speculation



2. Cloning

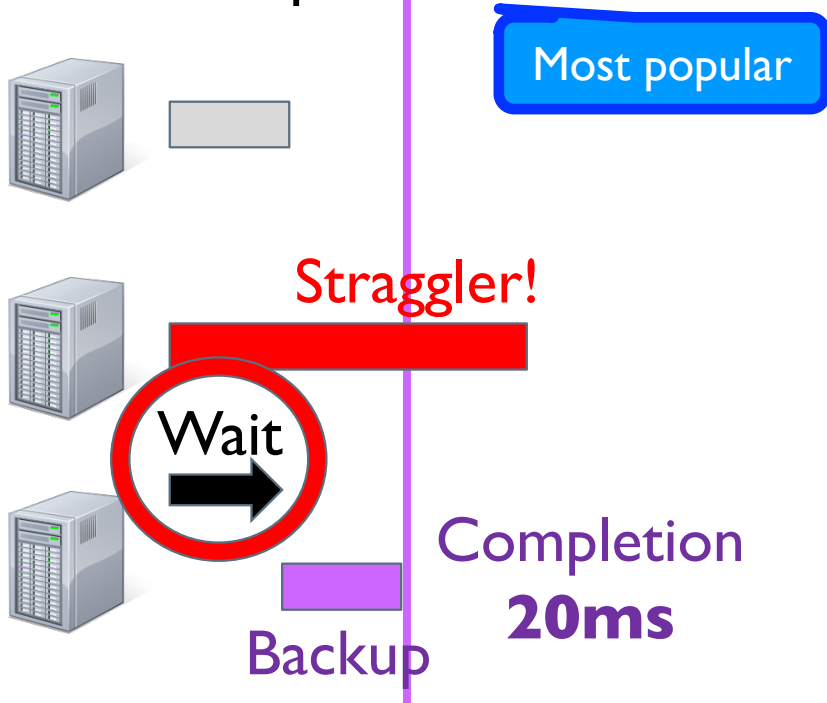
- Introduces **2x** workload

3. Snitching

- Does not work when burstiness fluctuates in ms-level

Current Tail-Tolerant Mechanisms

1. Speculation



2. Cloning

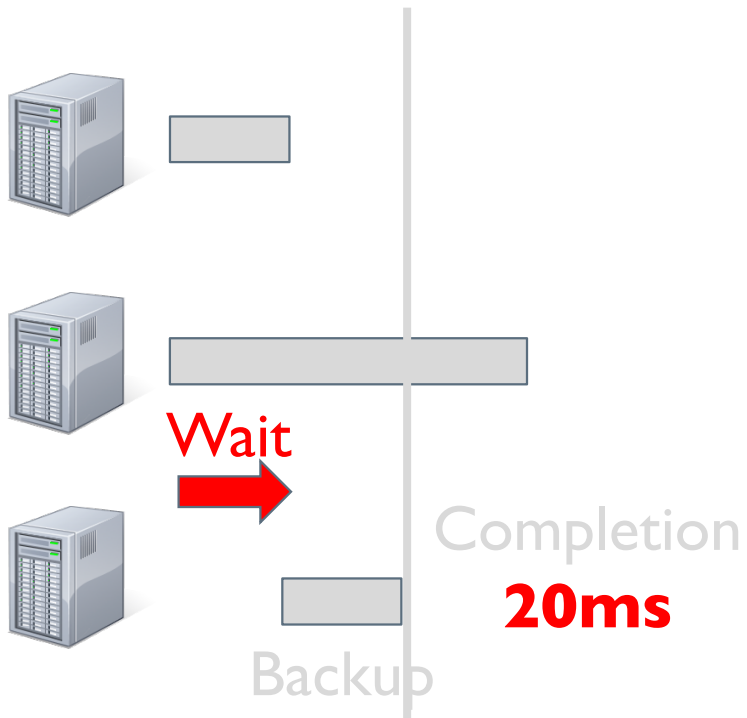
- Introduces **2x** workload

3. Snitching

- Does not work when burstiness fluctuates in ms-level

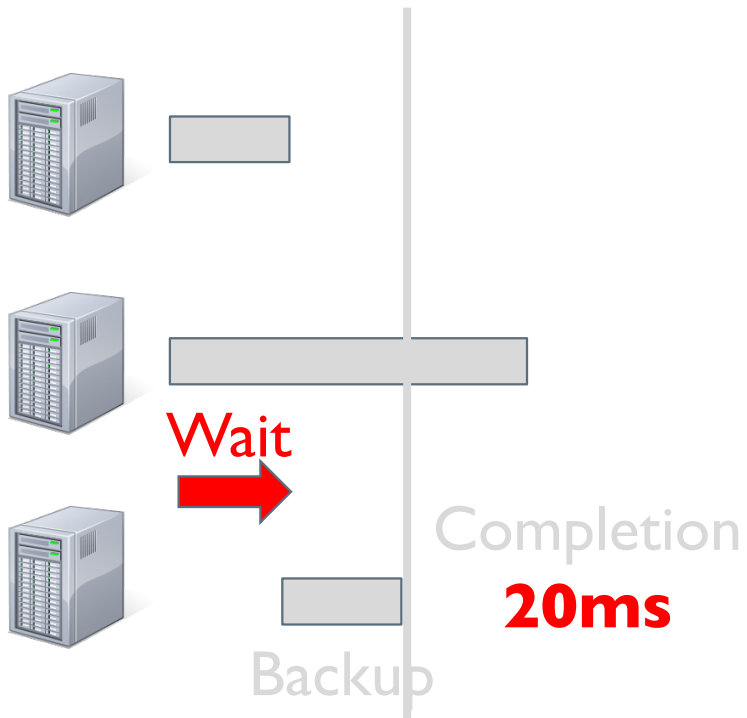


Must Wait!

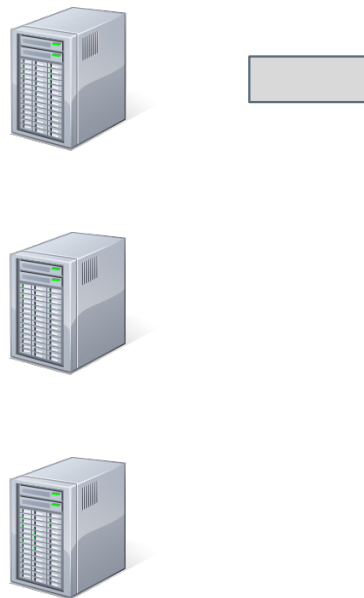




Must Wait!

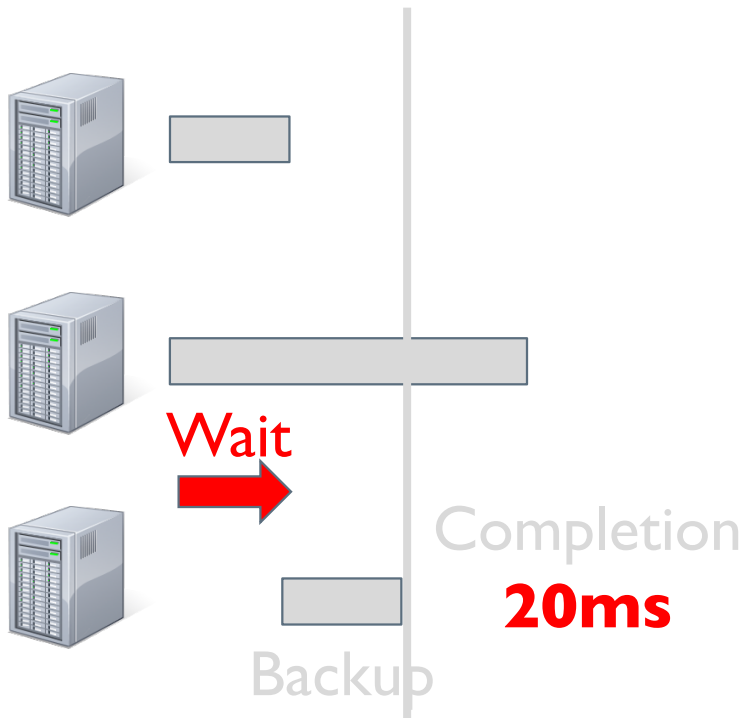


No Wait ?

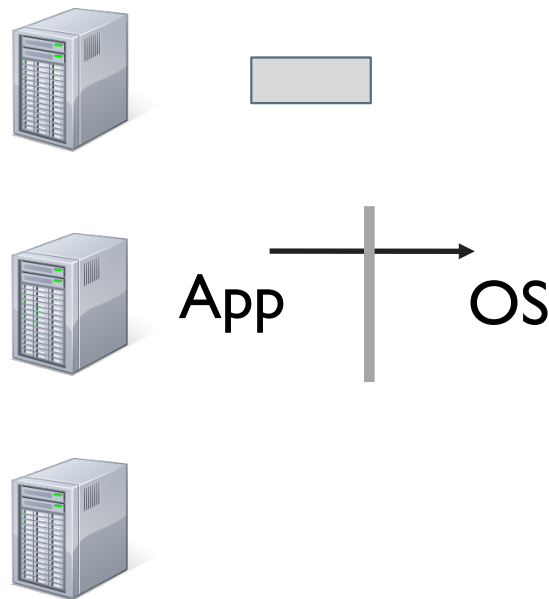




Must Wait!

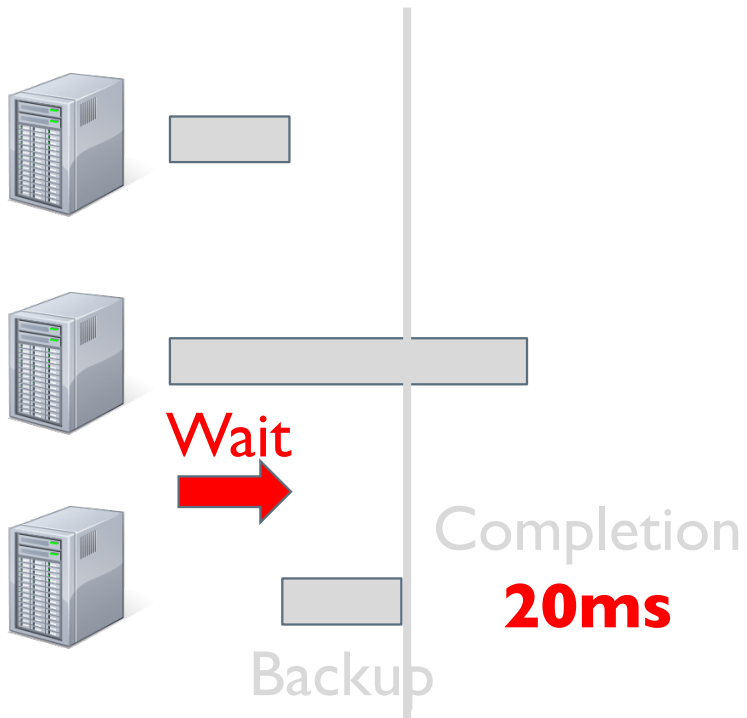


No Wait ?

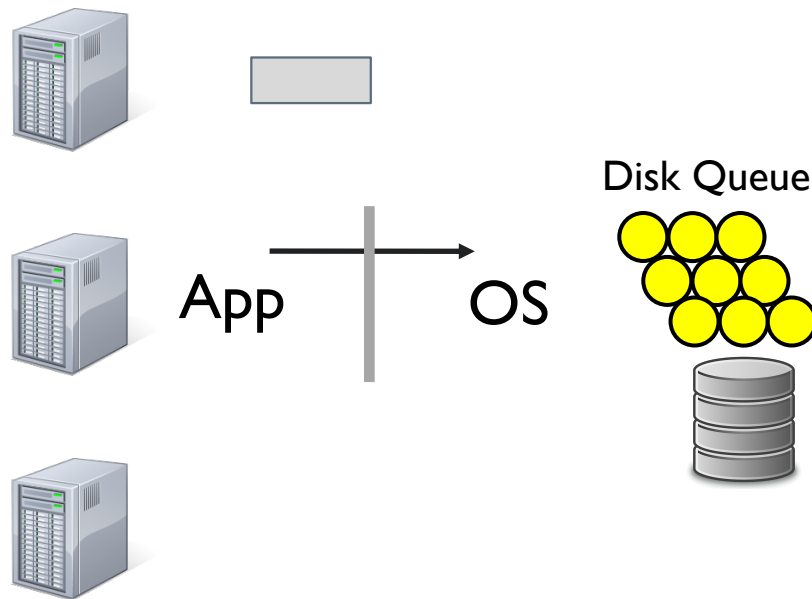




Must Wait!

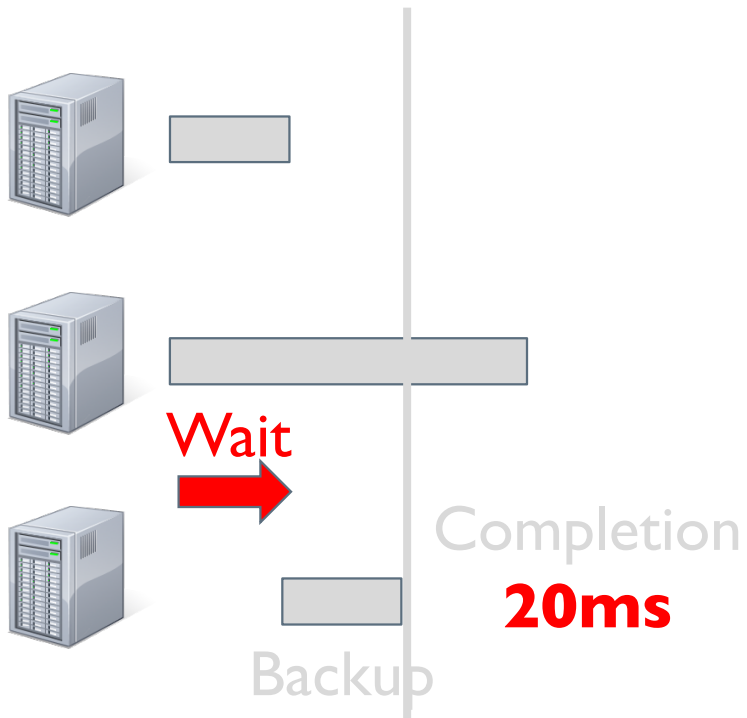


No Wait ?

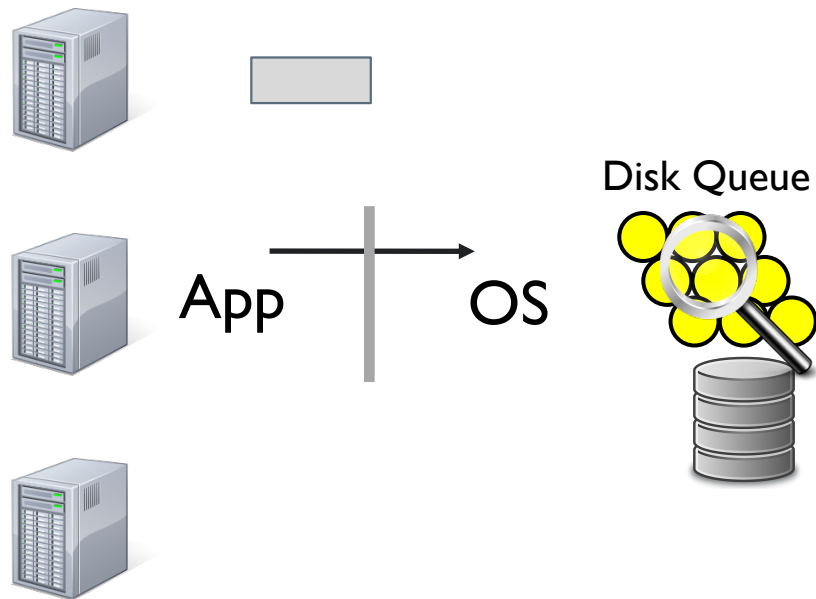




Must Wait!

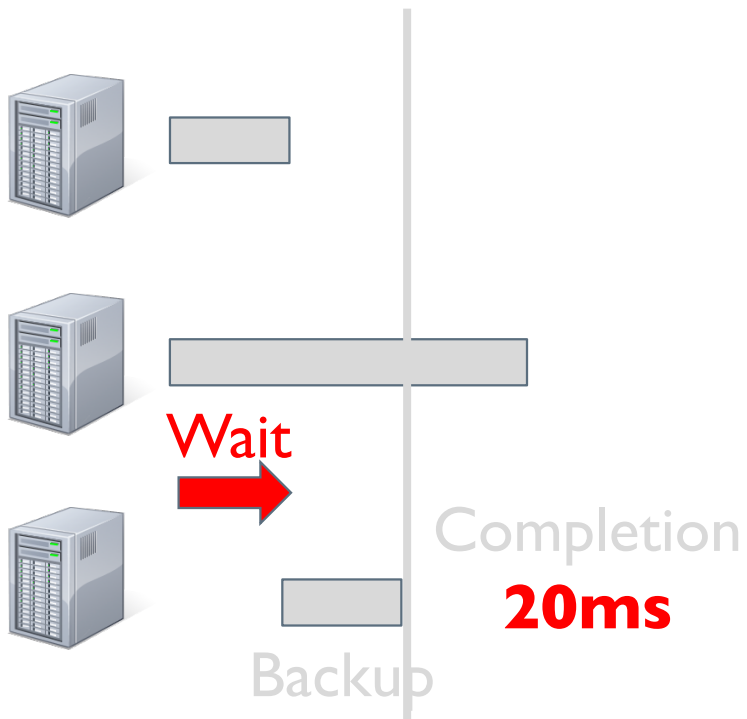


No Wait ?

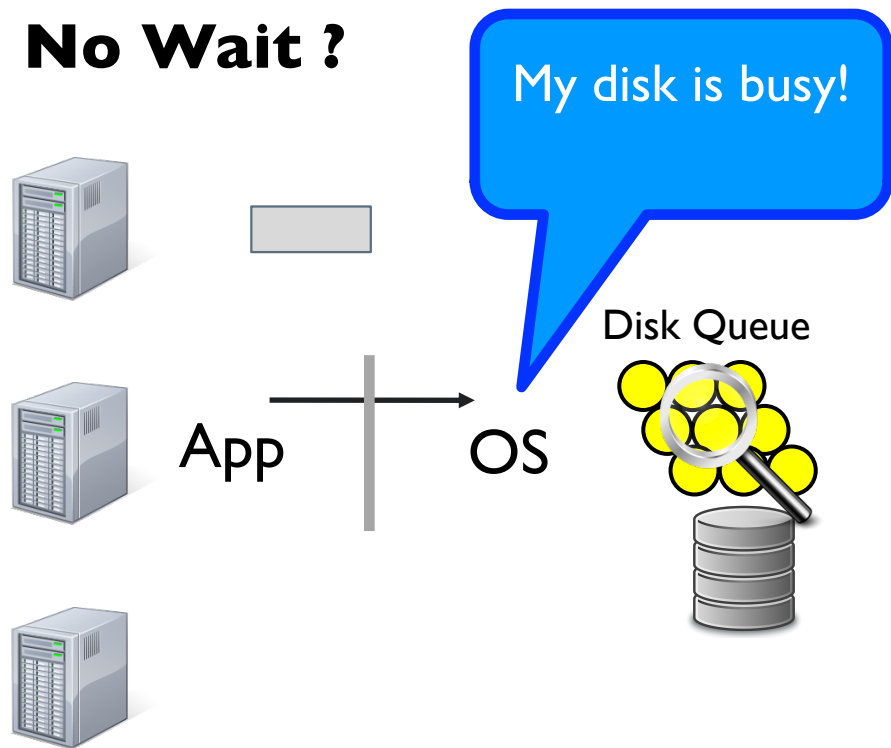




Must Wait!

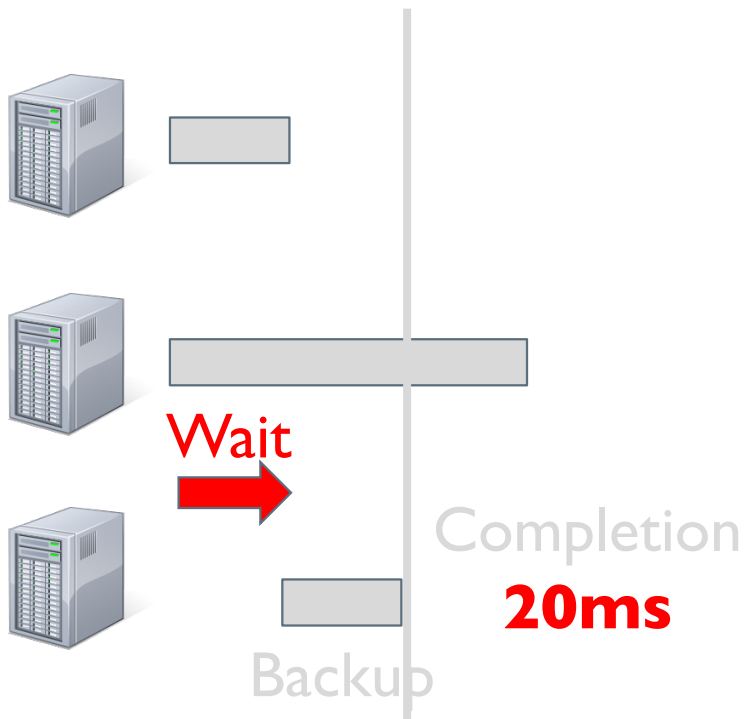


No Wait ?

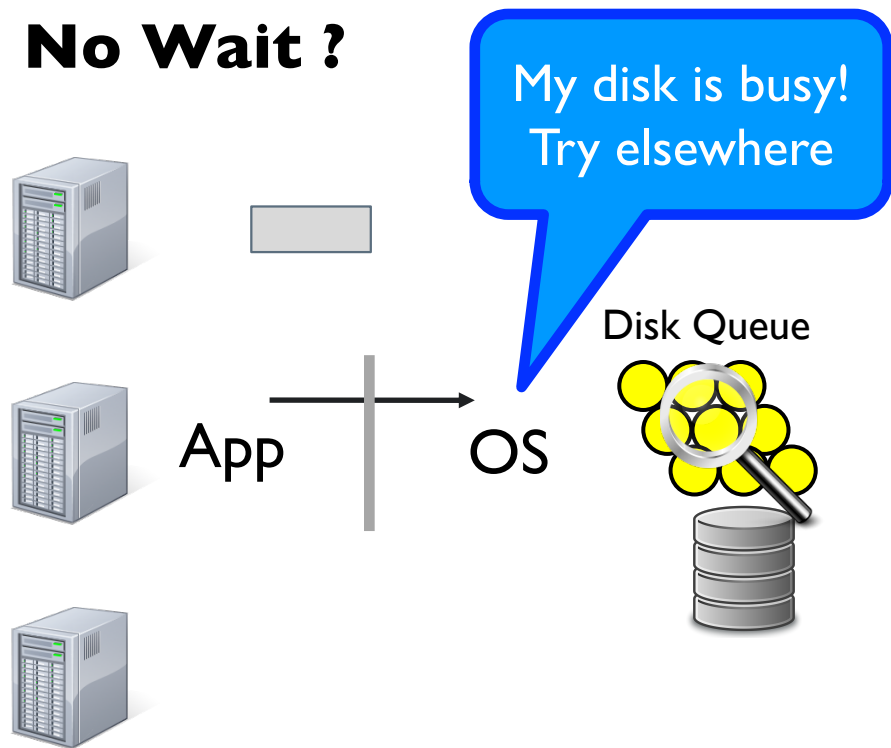




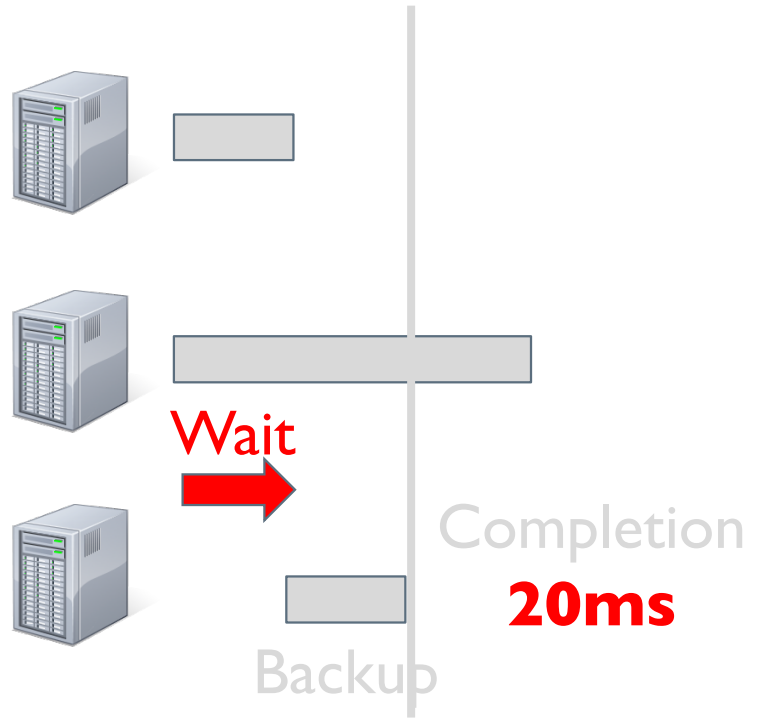
Must Wait!



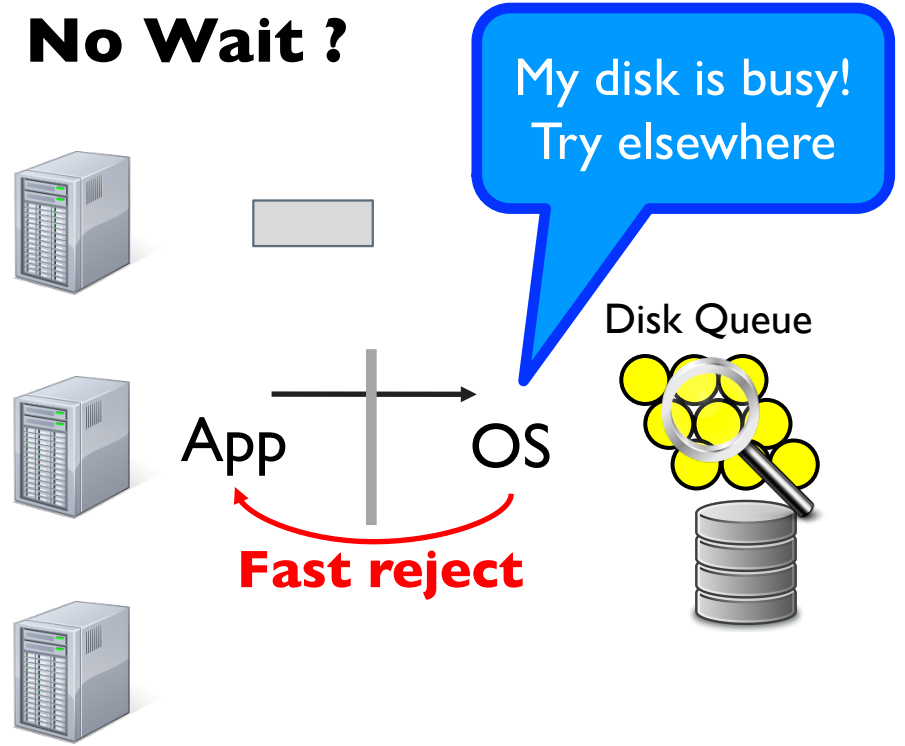
No Wait ?



Must Wait!

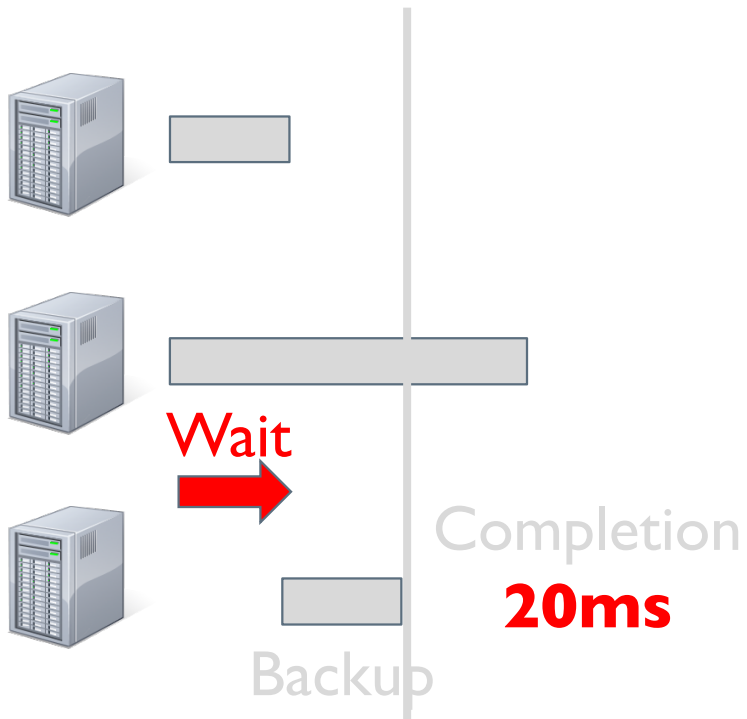


No Wait ?

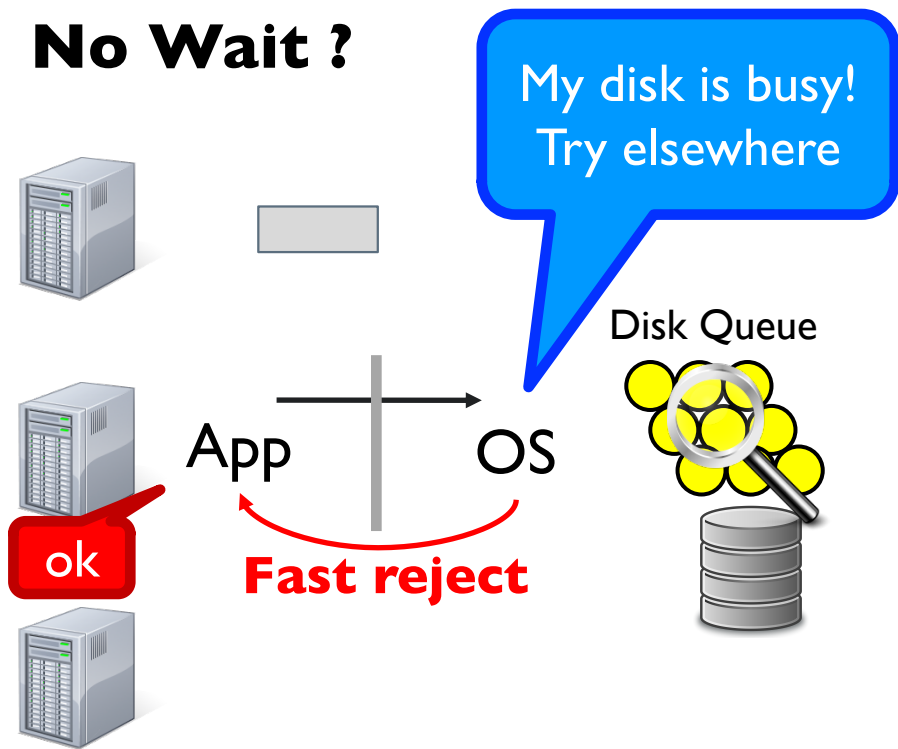




Must Wait!

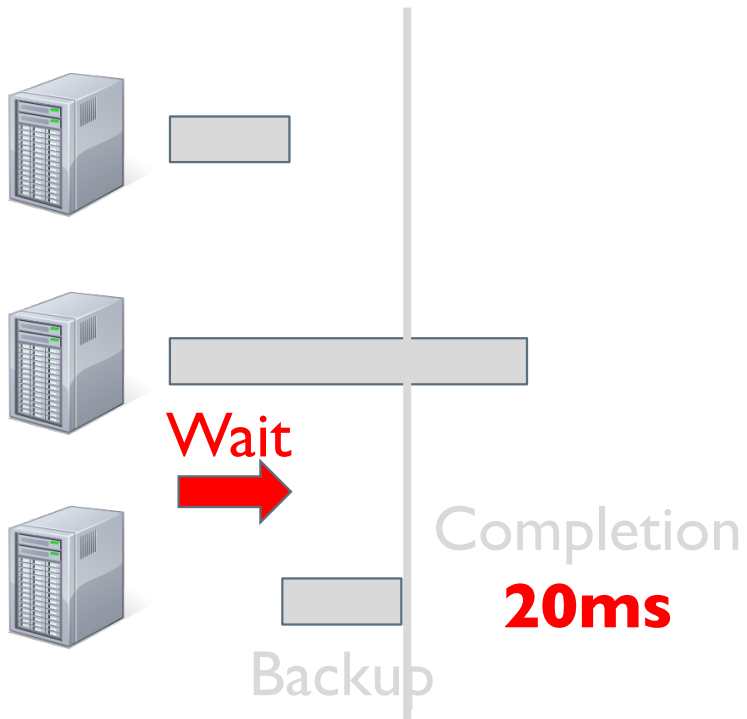


No Wait ?

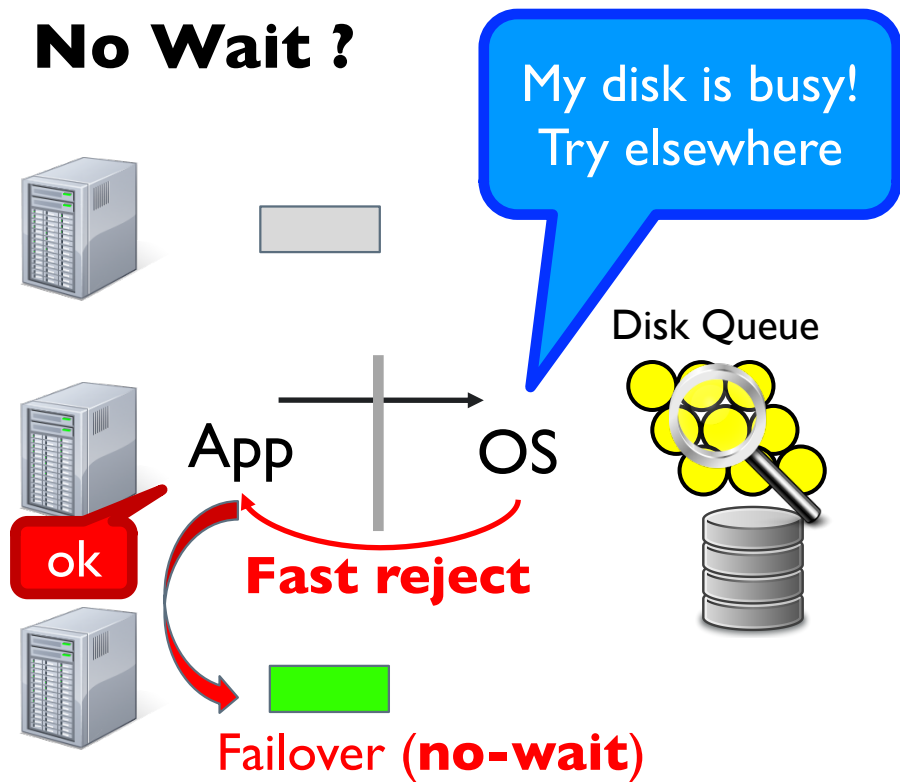




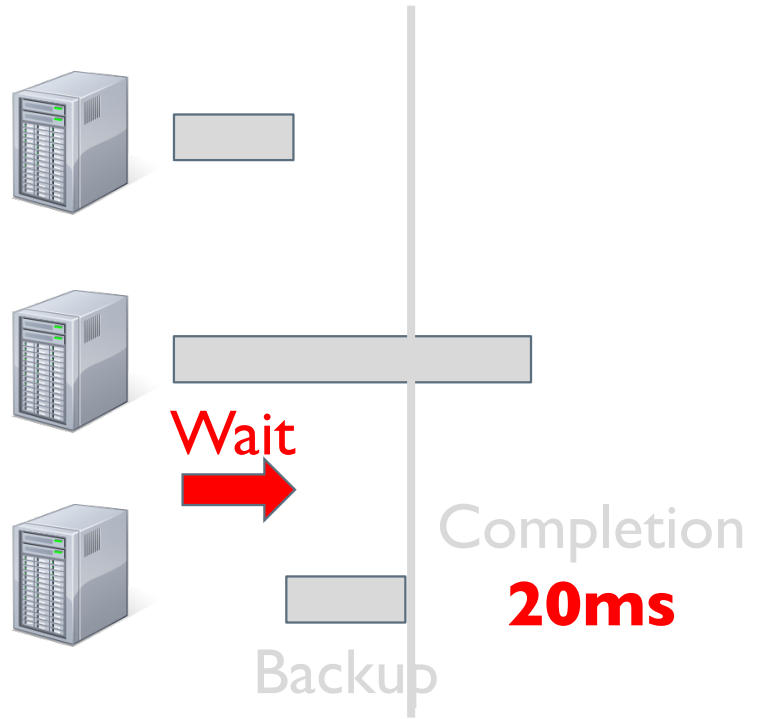
Must Wait!



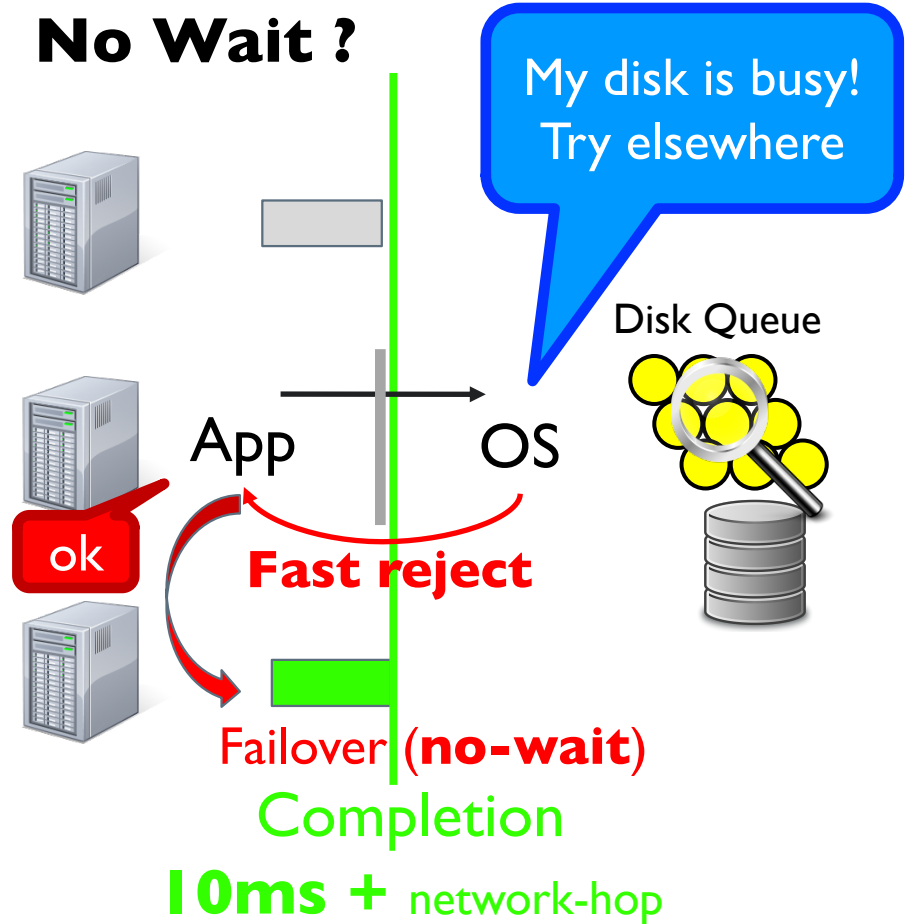
No Wait ?



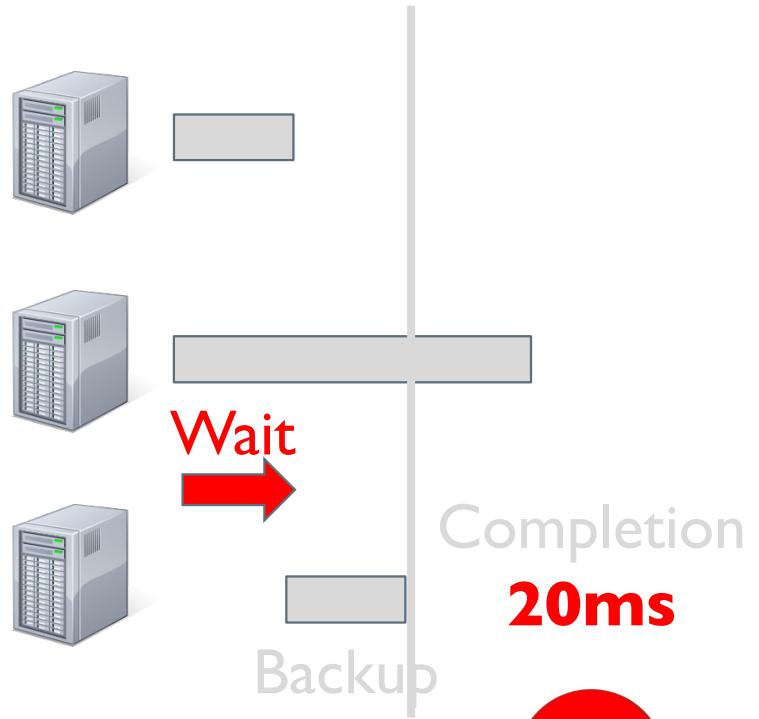
Must Wait!



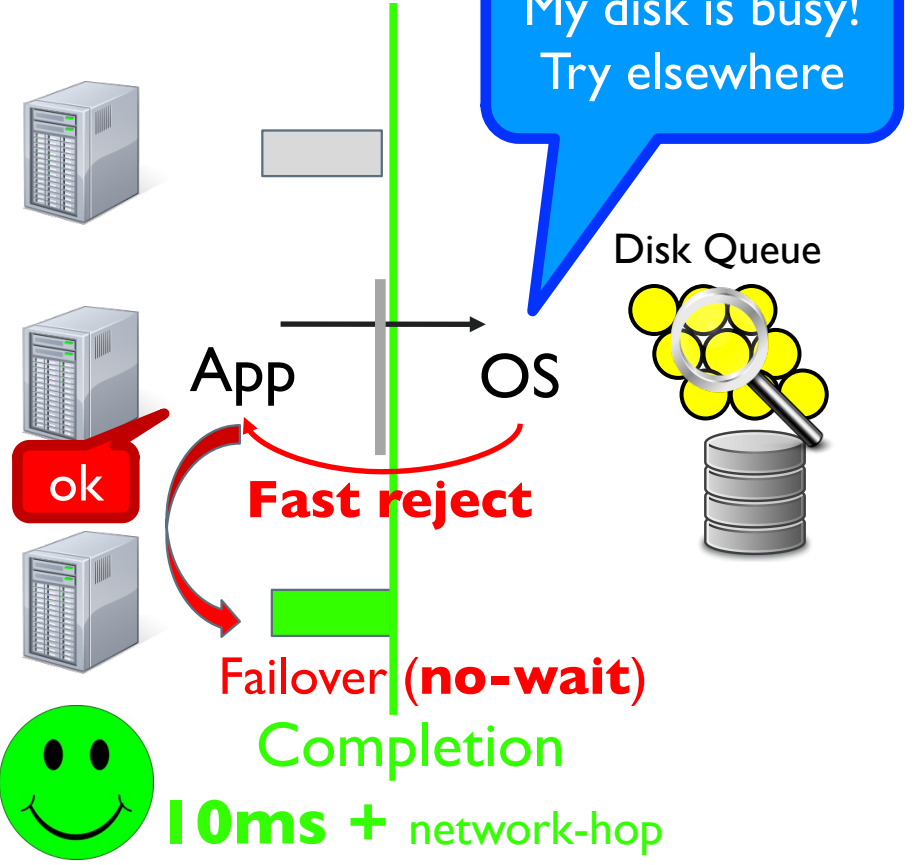
No Wait ?



Must Wait!

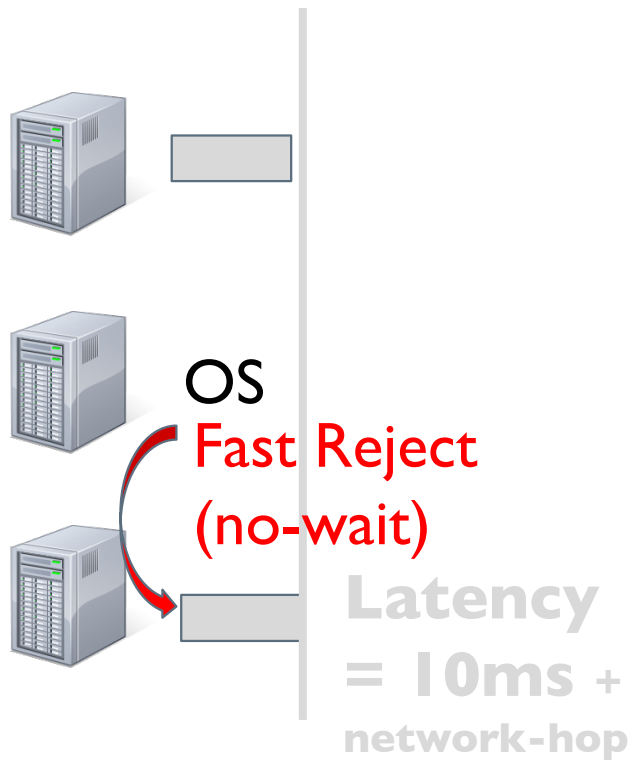


No Wait ?





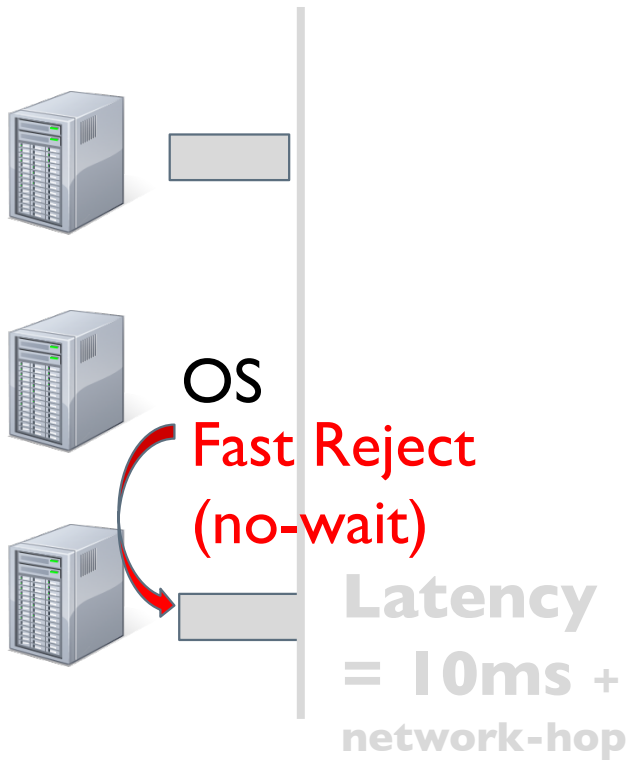
Use-Case



App

OS

Use-Case



I want < 20ms
latency

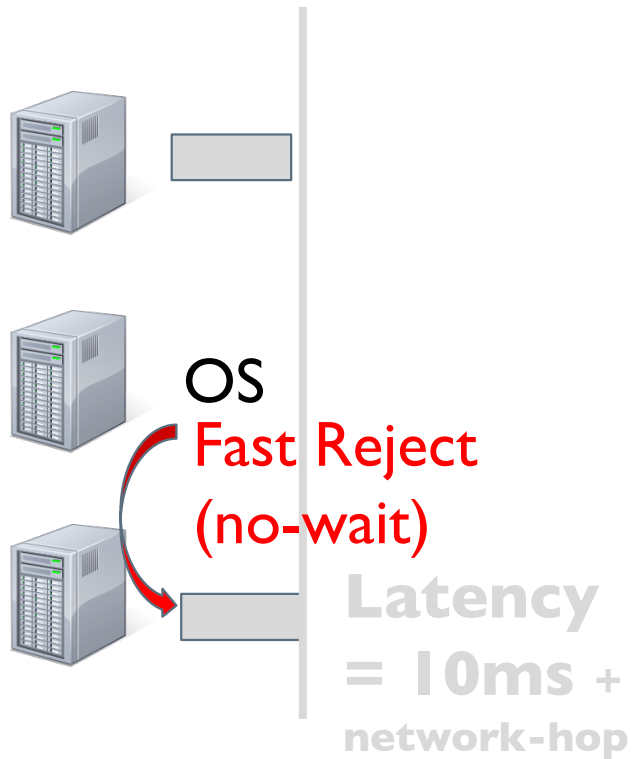
App

OS





Use-Case



I want < 20ms latency

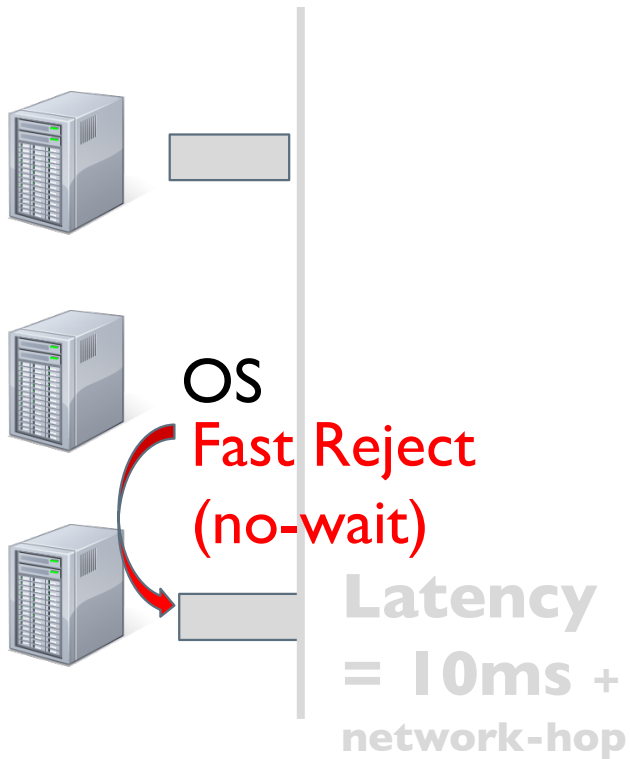
App

OS

1 SLO = 20ms



Use-Case



I want < 20ms latency

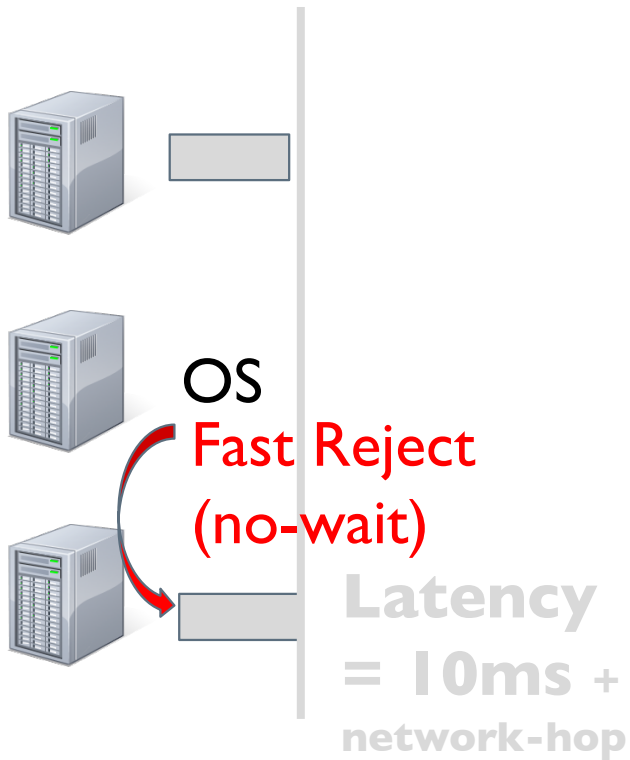
App

OS

- 1 SLO = 20ms
- 2 ret = read(..., SLO)



Use-Case



I want < 20ms
latency

App

OS

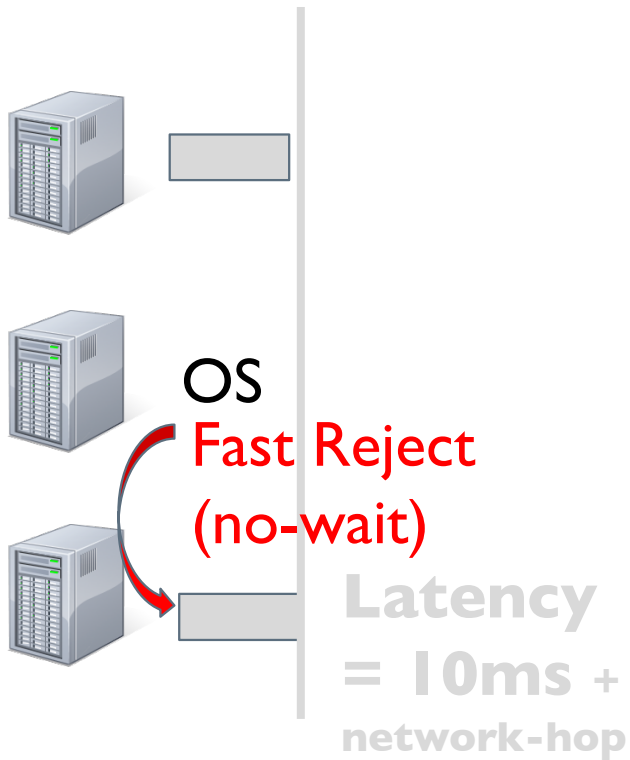
1 SLO = 20ms

2 ret = read(..., SLO) →

3



Use-Case



I want < 20ms latency

App

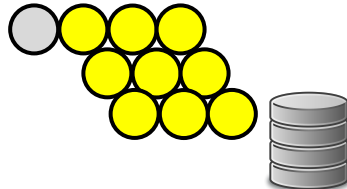
OS

1 SLO = 20ms

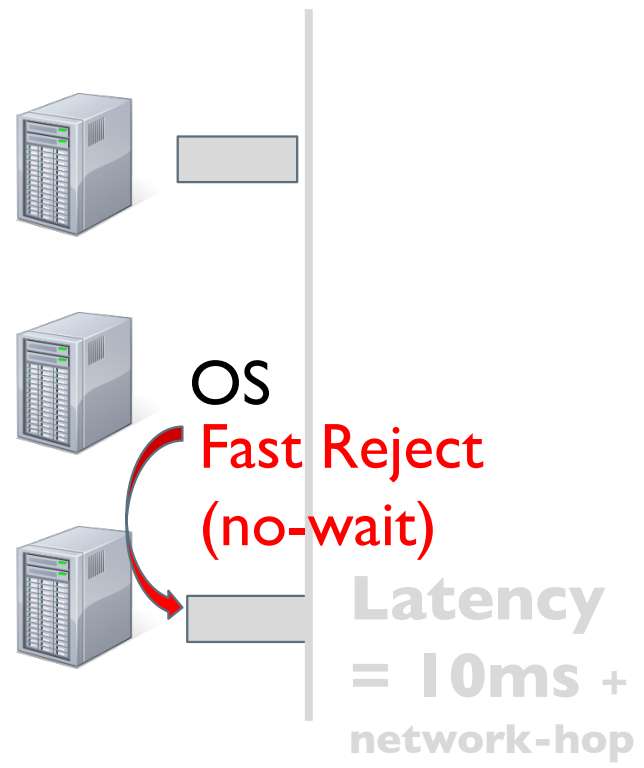
2 ret = read(..., SLO)

3

Disk Queue



Use-Case



I want < 20ms latency

App

OS can see "everything" and tell app when it is busy

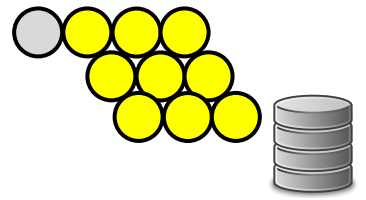
OS

1 SLO = 20ms

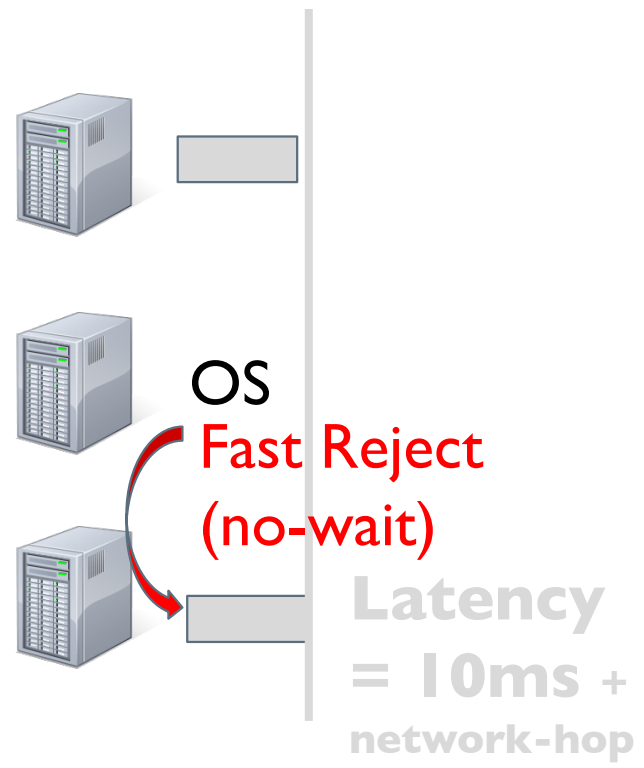
2 ret = read(..., SLO)

3

Disk Queue



Use-Case



I want < 20ms
latency

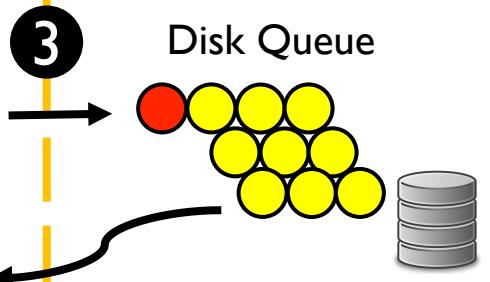
App

OS can see
"everything" and tell
app when it is busy

OS

1 SLO = 20ms

2 ret = read(..., SLO)

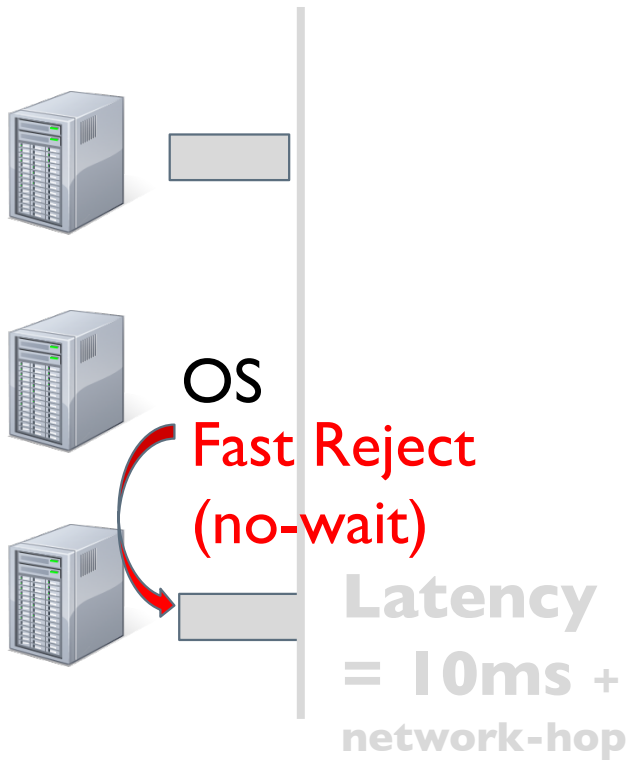


3

4 Reject fast



Use-Case



I want < 20ms latency

App

OS can see "everything" and tell app when it is busy

OS

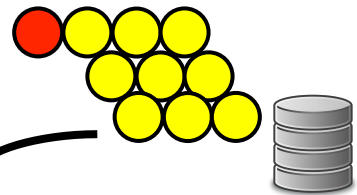
1 SLO = 20ms

2 ret = read(..., SLO)

5 if (ret == Reject)
// failover

3

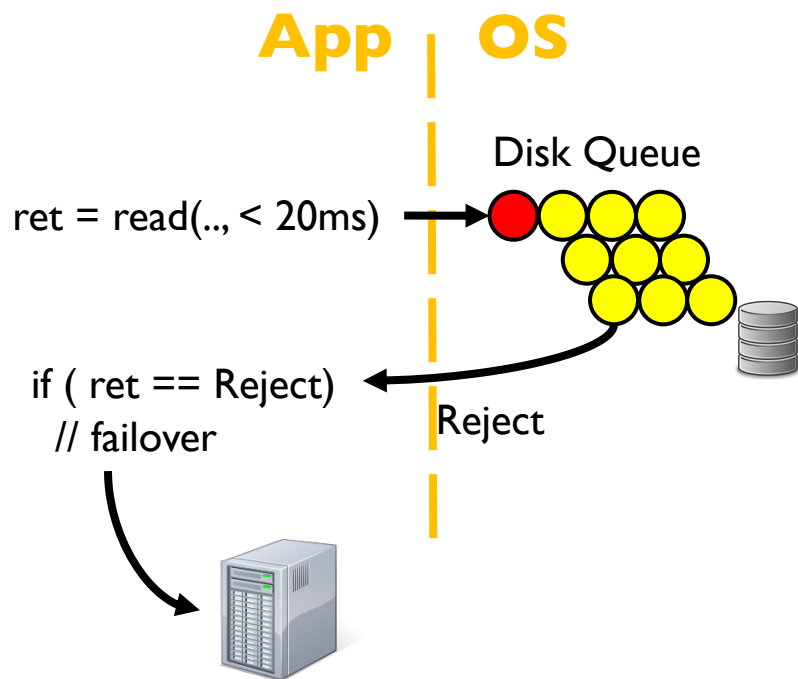
Disk Queue



4 Reject fast

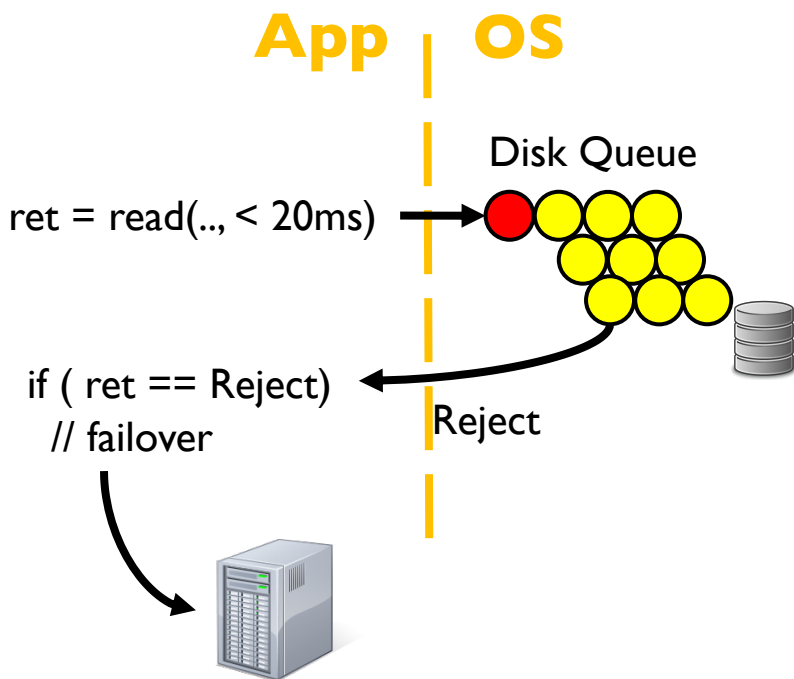


MittOS



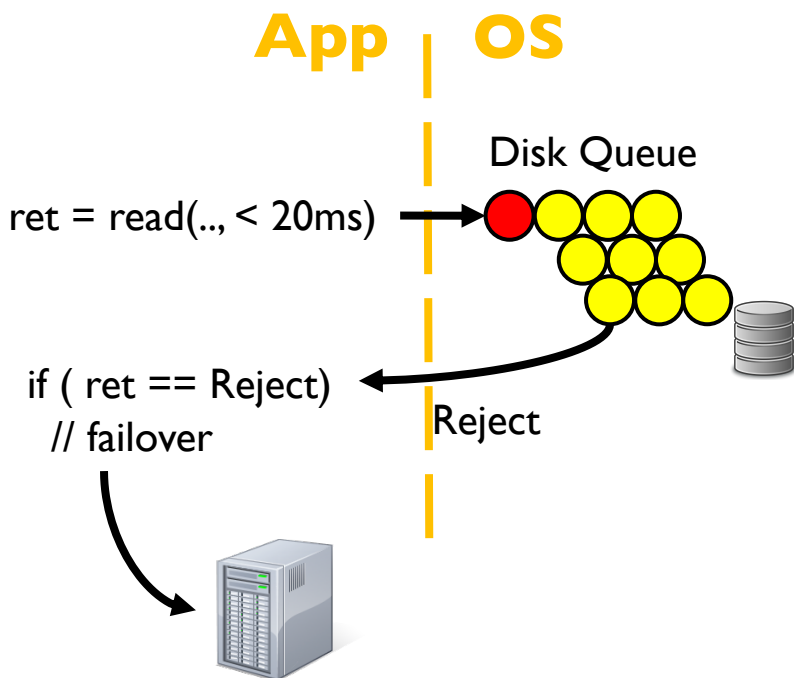
- MittOS **Principles**

MittOS



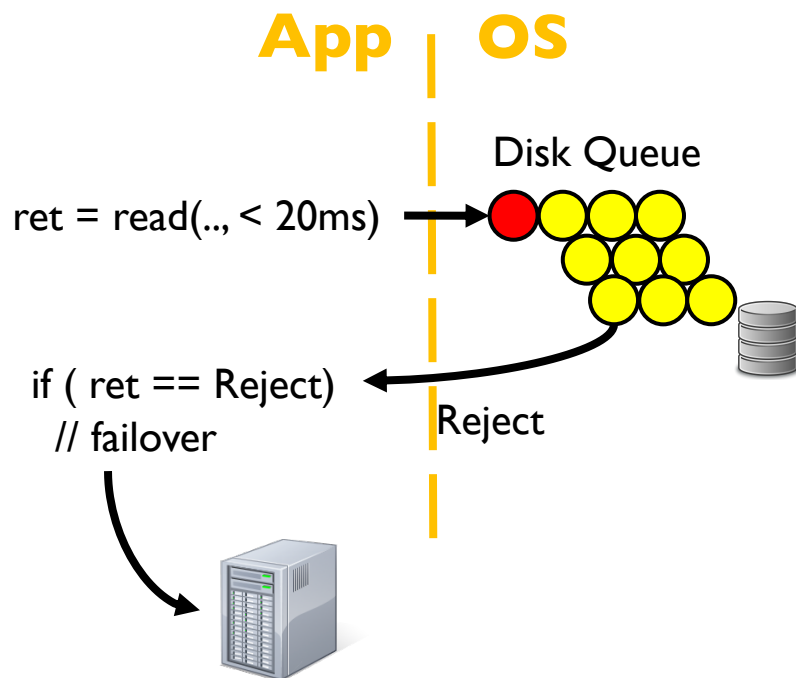
- MittOS **Principles**
 - SLO-aware interface
 - Reject fast

MittOS



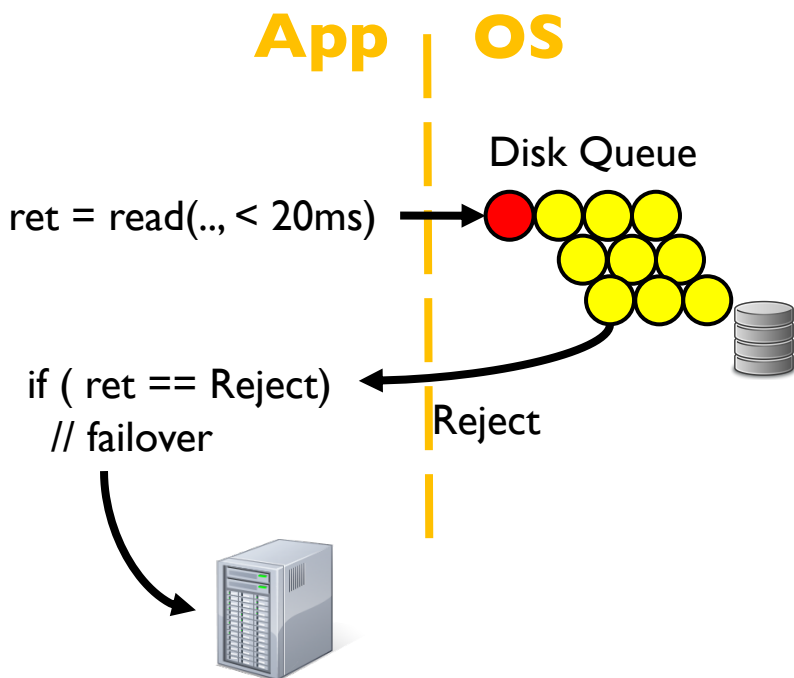
- MittOS **Principles**
 - SLO-aware interface
 - Reject fast
 - *Transparent of busyness*

MittOS



- MittOS **Principles**
 - SLO-aware interface
 - Reject fast
 - *Transparent of busyness*
 - **PC** era: is best effort (cannot reject IOs)

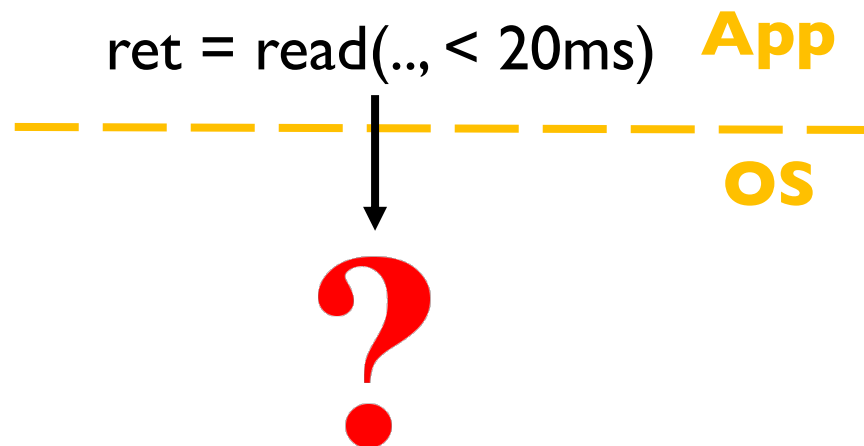
MittOS



- MittOS **Principles**
 - SLO-aware interface
 - Reject fast
 - *Transparent of busyness*
 - **PC** era: is best effort (cannot reject IOs)
 - **DC** era: Less-busy replicas available

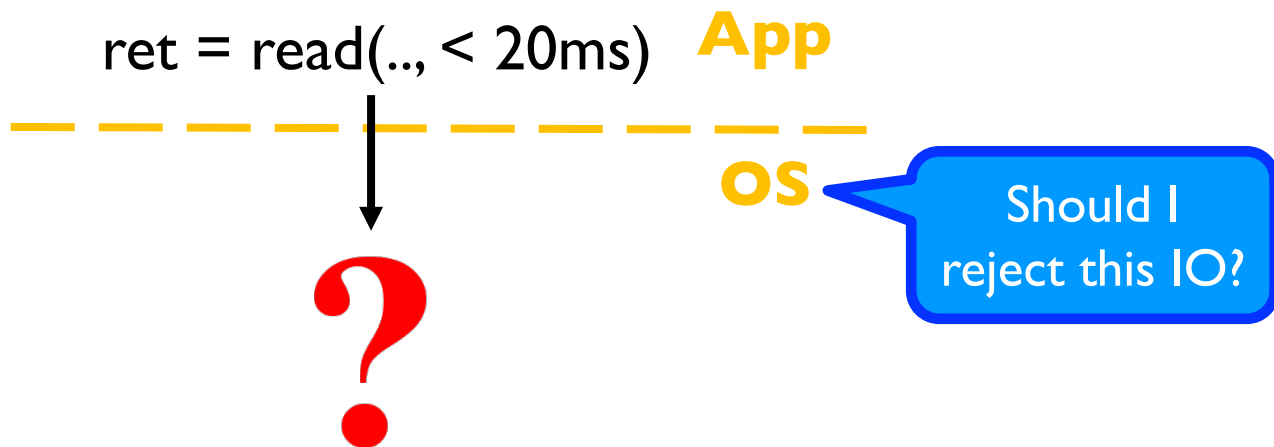


Challenge



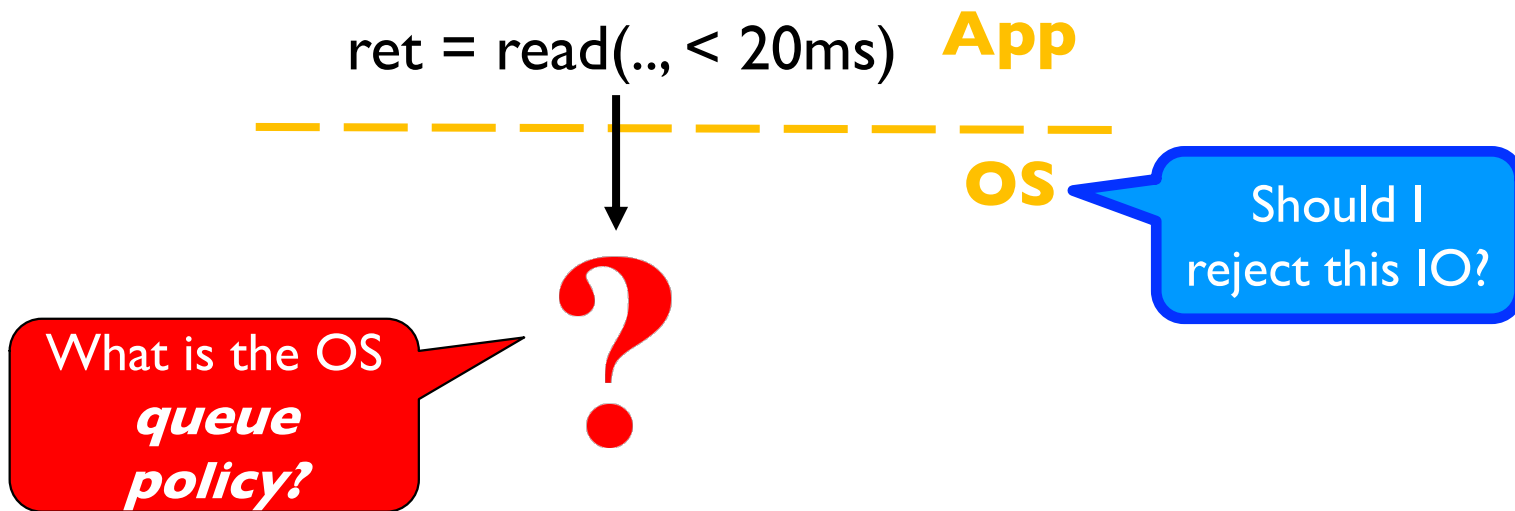


Challenge



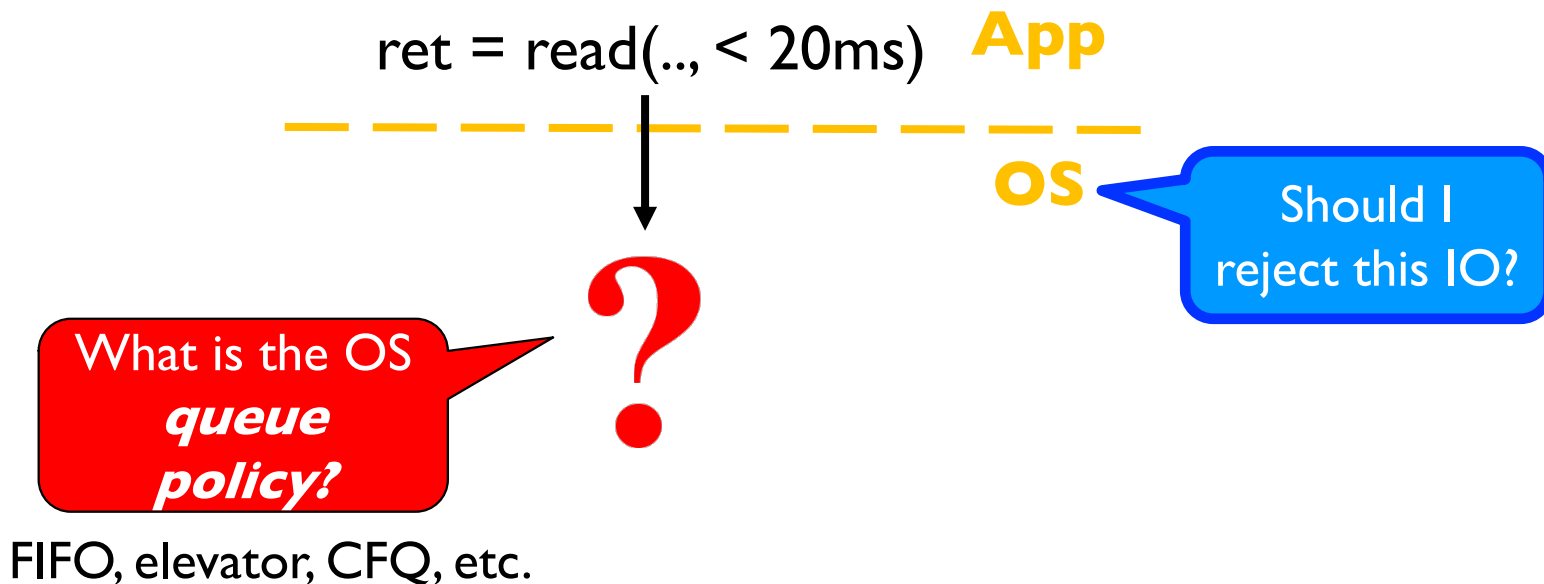


Challenge



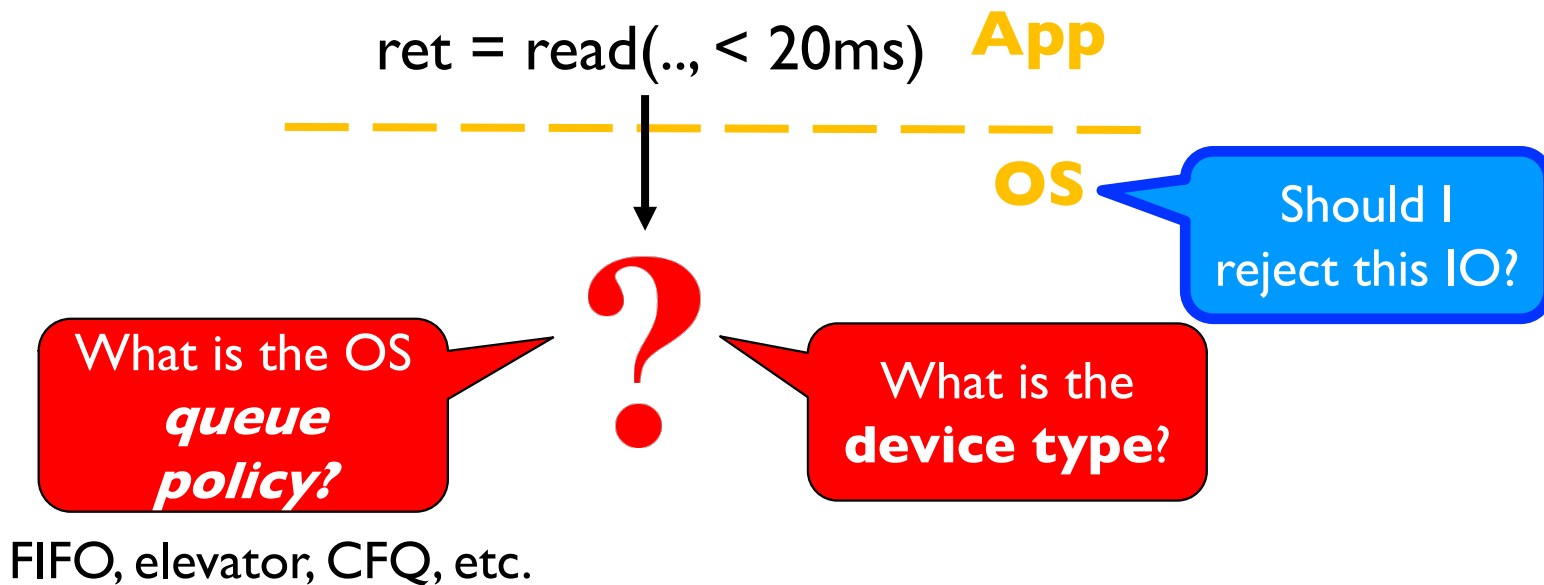


Challenge



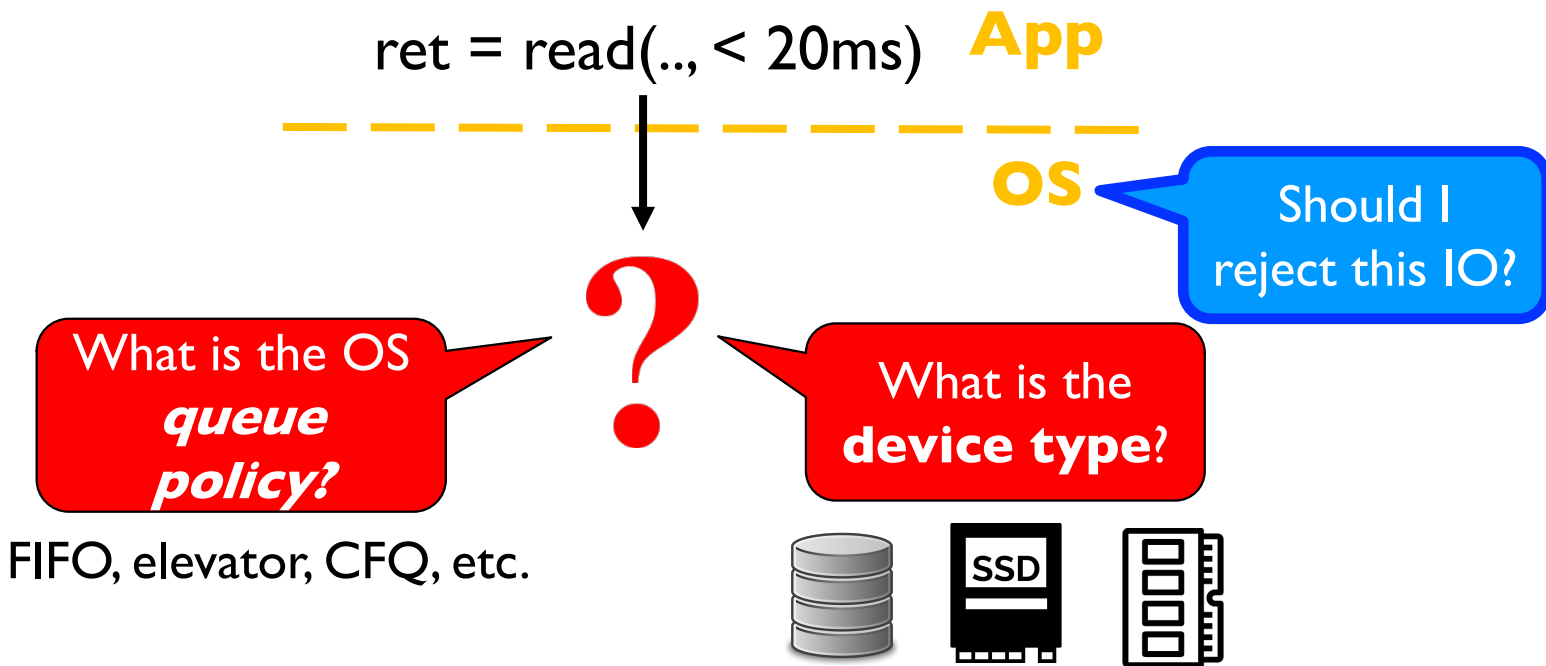


Challenge



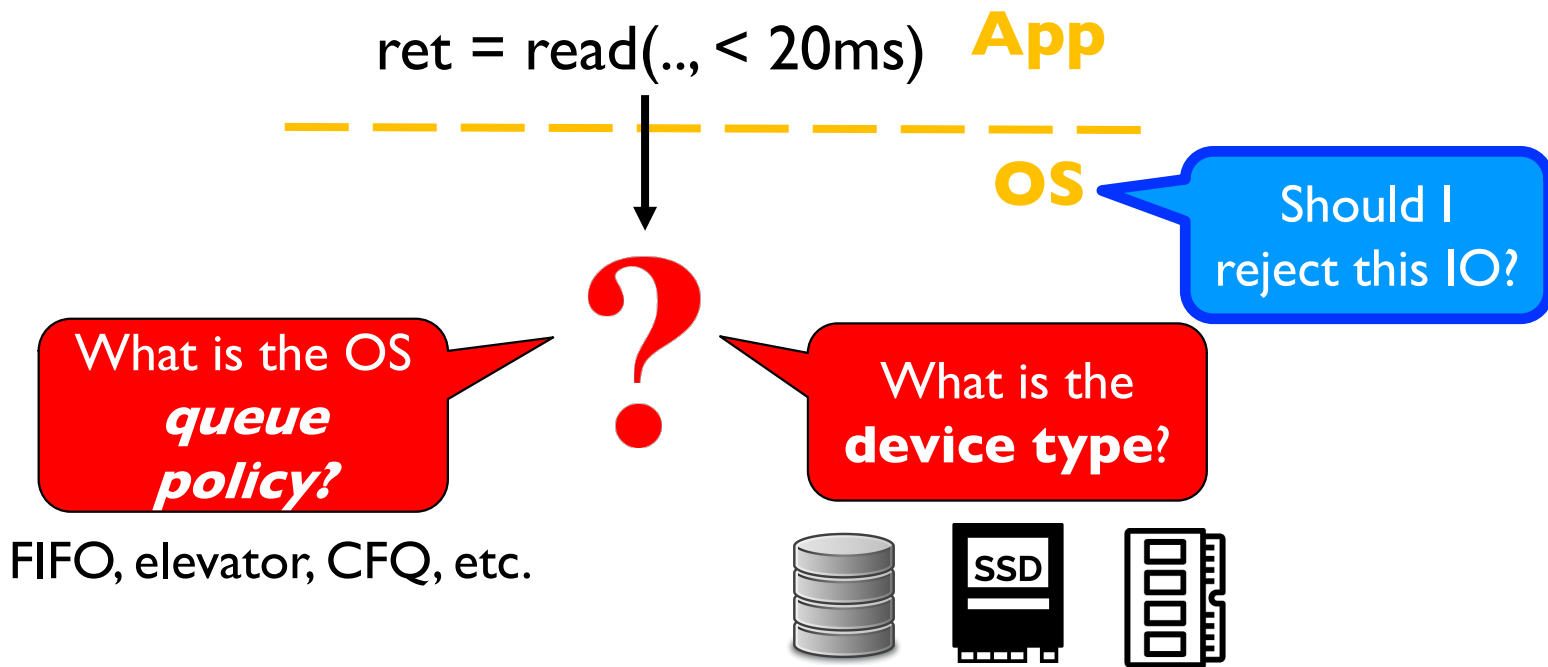


Challenge





Challenge



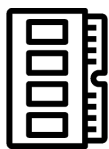
Prediction depends on **queue policy** and **device type**

Contribution

MittOS principle: Support fast
rejecting SLO-aware interface



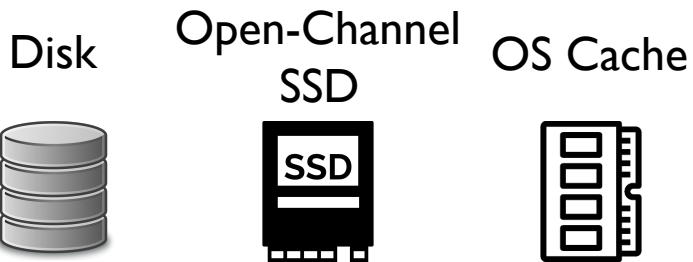
Contribution



MittOS principle: Support fast
rejecting SLO-aware interface

Contribution

MittOS Latency Prediction



MittOS principle: Support fast rejecting SLO-aware interface

Contribution

Fast Reject Interface

MittOS Latency Prediction

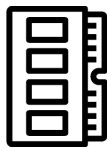
Disk



Open-Channel
SSD



OS Cache



MittOS principle: Support fast rejecting SLO-aware interface

Contribution

MittOS-
powered



MongoDB



Fast Reject Interface

MittOS Latency Prediction

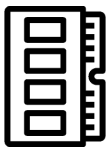
Disk



Open-Channel
SSD



OS Cache



MittOS principle: Support fast
rejecting SLO-aware interface

Contribution **+50 LOC**

MittOS-
powered



MongoDB



Fast Reject Interface

MittOS Latency Prediction

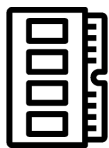
Disk



Open-Channel
SSD



OS Cache



MittOS principle: Support fast rejecting SLO-aware interface

Contribution **+50 LOC**

MittOS-powered



MongoDB



Fast Reject Interface

MittOS Latency Prediction

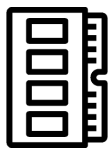
Disk



Open-Channel
SSD



OS Cache



vs. **state of the art**:
hedged requests, cloning,
application timeout, etc.

MittOS principle: Support fast rejecting SLO-aware interface

Contribution **+50 LOC**

MittOS-powered



MongoDB



Fast Reject Interface

MittOS Latency Prediction

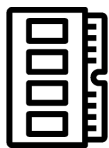
Disk



Open-Channel
SSD



OS Cache



MittOS principle: Support fast rejecting SLO-aware interface

vs. **state of the art:**
hedged requests, cloning,
application timeout, etc.

Cut tail:
**50% latency reduction
above 75 percentile**



Outline

□ Introduction

□ **Design**

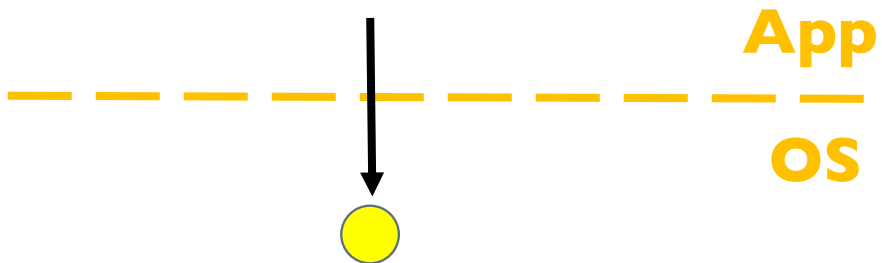
- Challenges
- Solutions

□ Evaluation

□ Conclusion

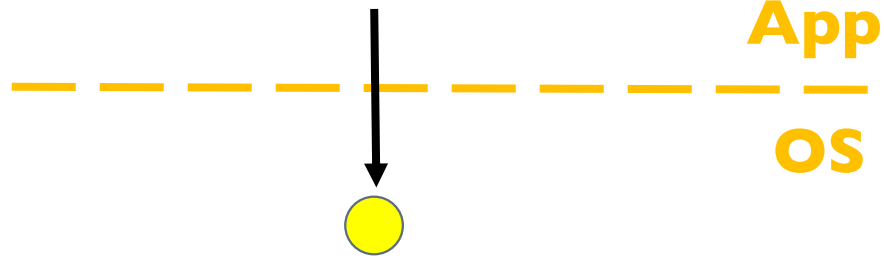
Prediction

ret = read(.., < 20ms)



Prediction

ret = read(.., < 20ms)

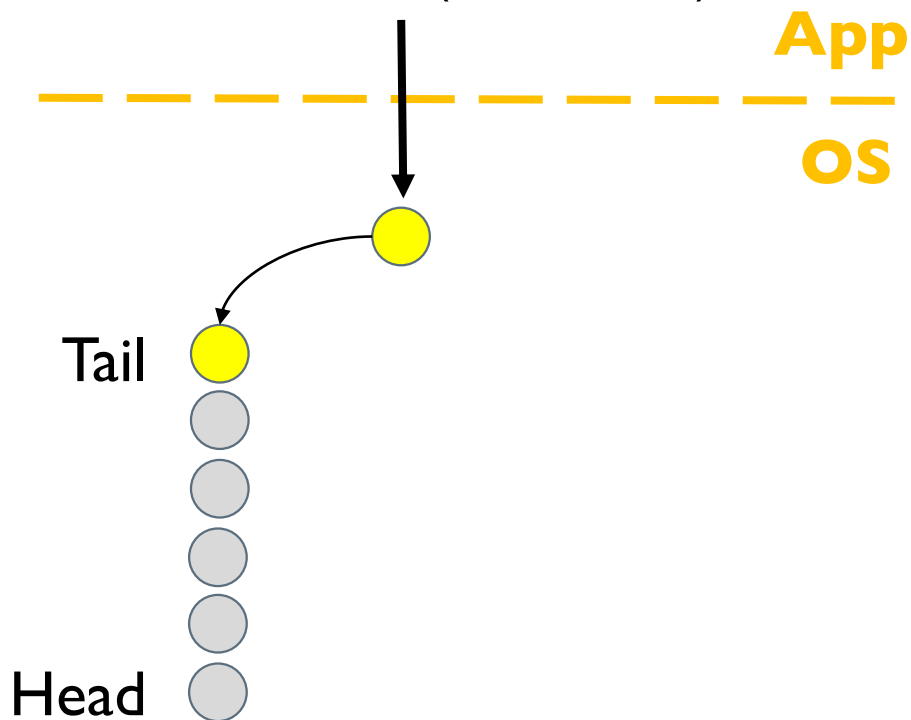


How to **predict** latency *before* submitting to the device?



Prediction

ret = read(.., < 20ms)



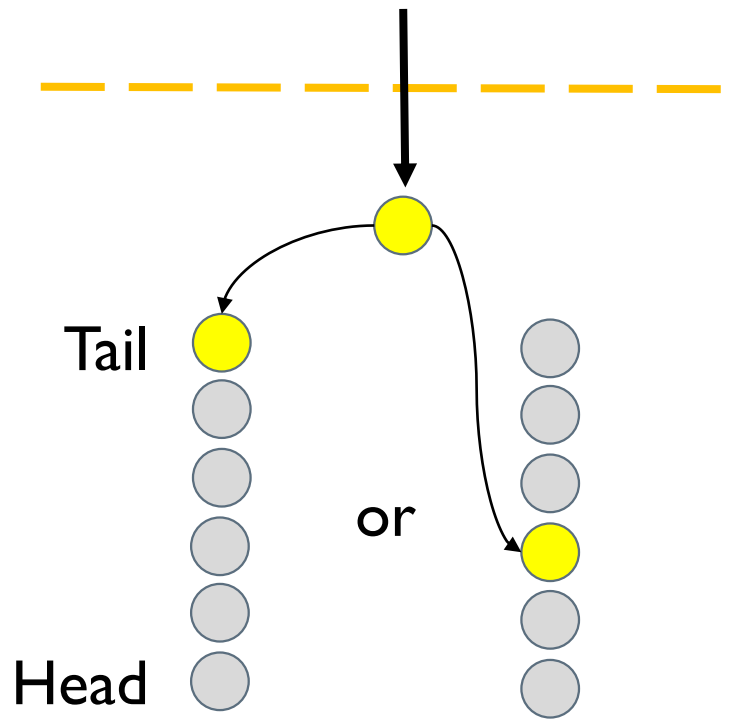
How to **predict** latency *before* submitting to the device?

Prediction

ret = read(.., < 20ms)

App
OS

How to **predict** latency *before* submitting to the device?

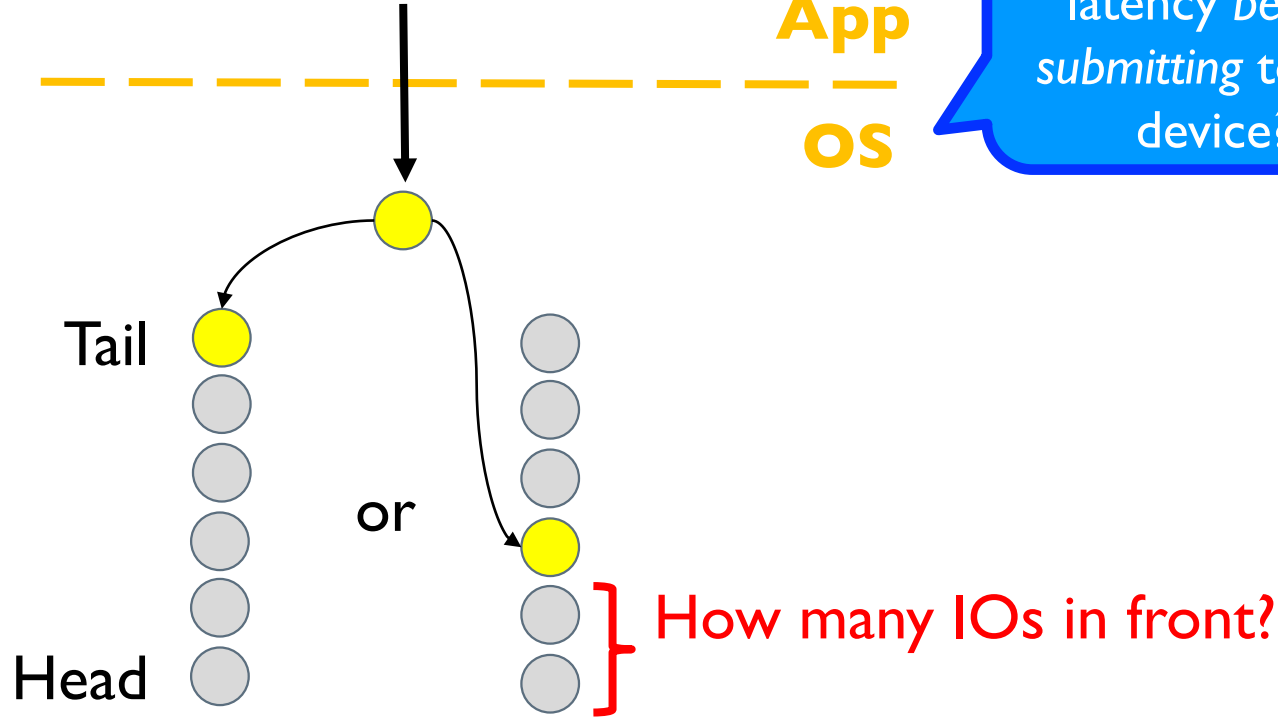


Prediction

ret = read(.., < 20ms)

App
OS

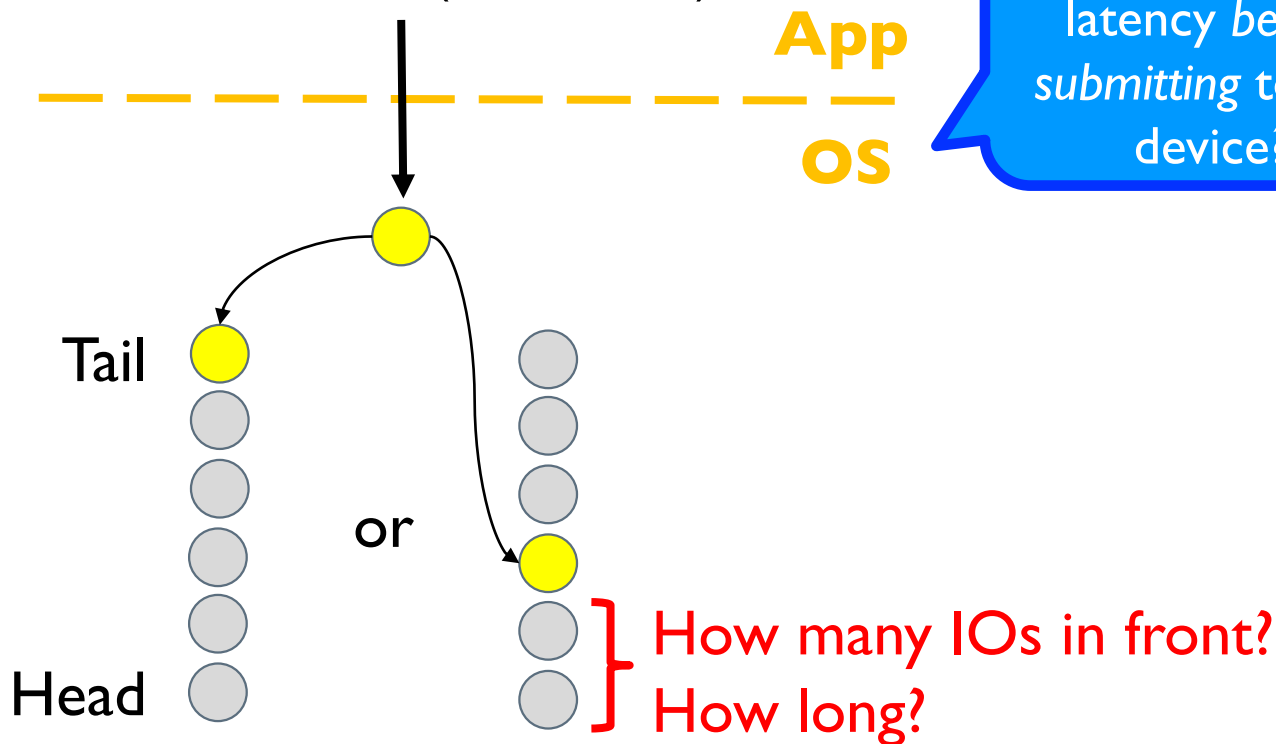
How to **predict** latency *before* submitting to the device?





Prediction

ret = read(.., < 20ms)



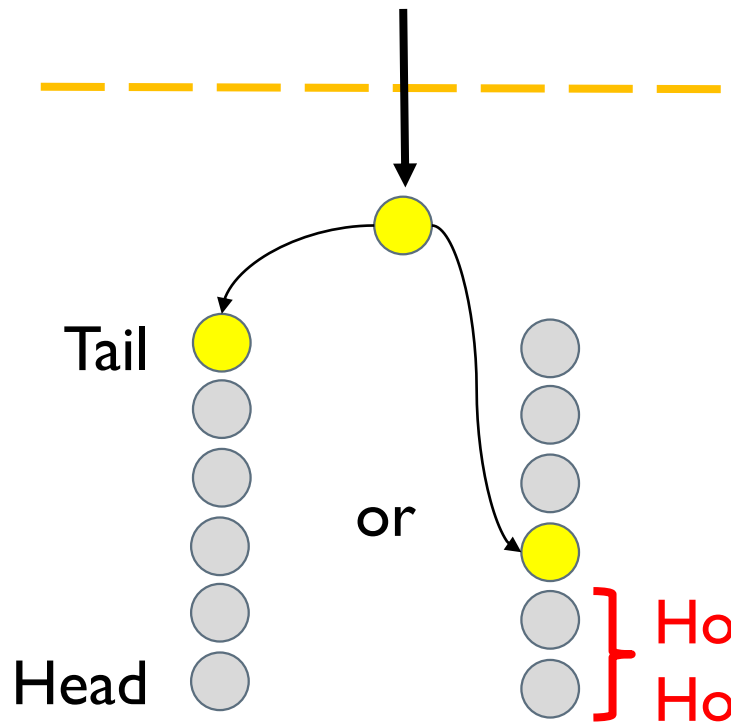
How to **predict** latency *before* submitting to the device?

Prediction

ret = read(.., < 20ms)

App
OS

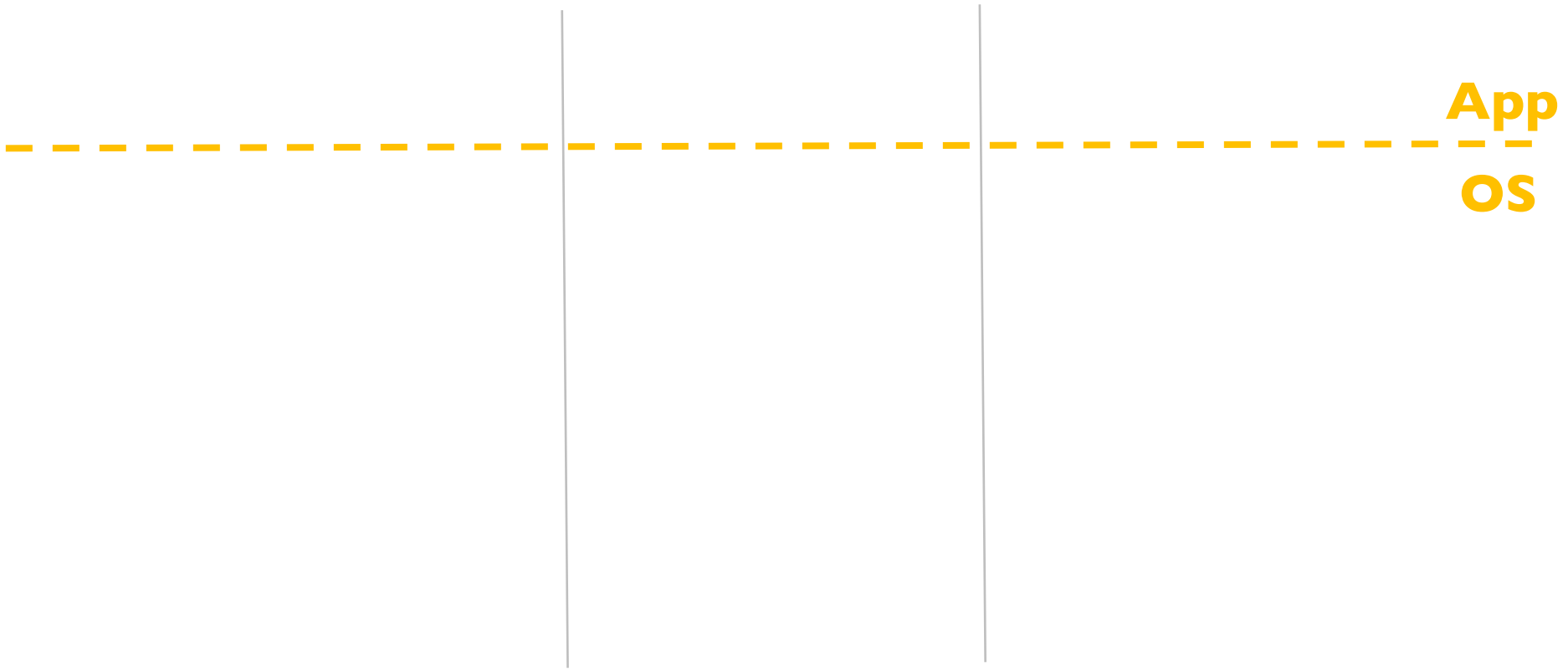
How to **predict** latency *before* submitting to the device?



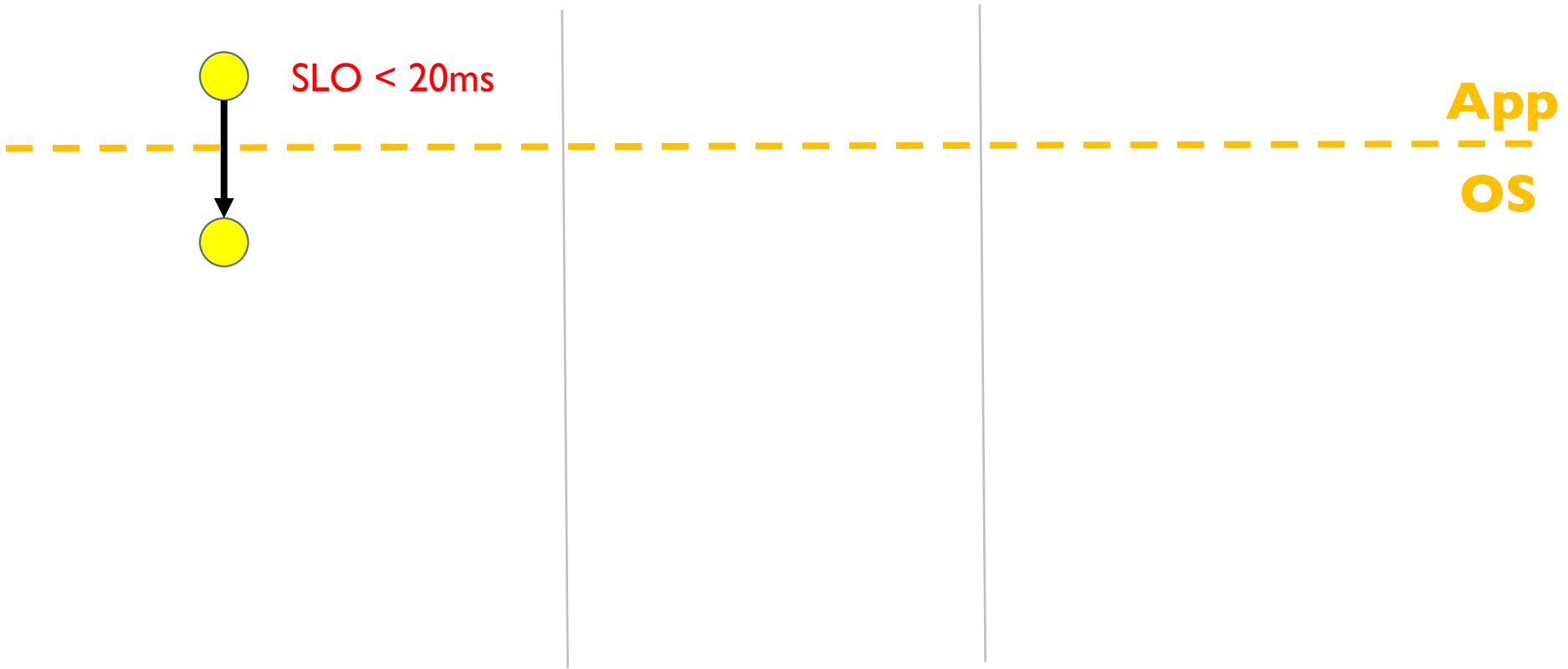
Latency < SLO → **Accept**
Latency > SLO → **Reject**

How many IOs in front?
How long?

Challenge #1: Modeling Queue Policy



Challenge #1: Modeling Queue Policy



Challenge #1: Modeling Queue Policy



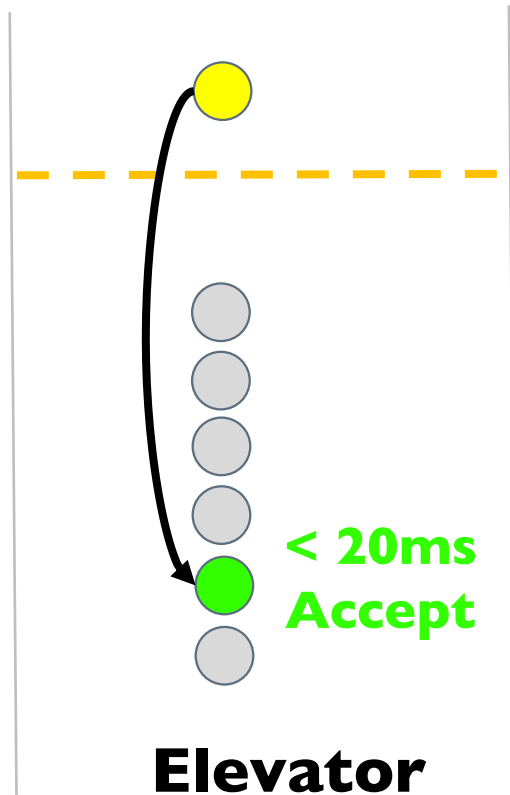
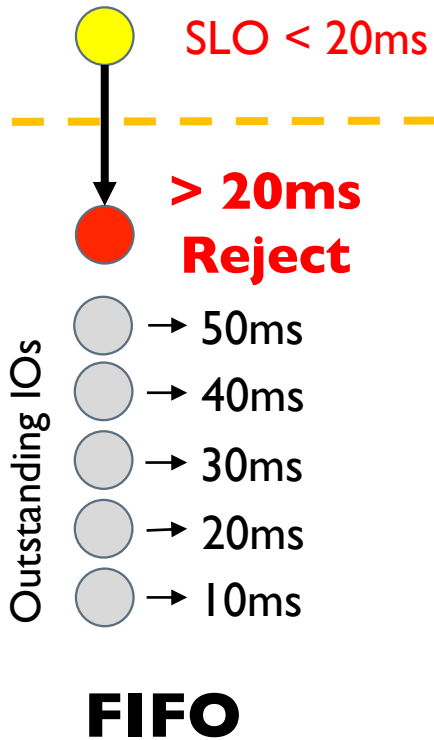
Challenge #1: Modeling Queue Policy



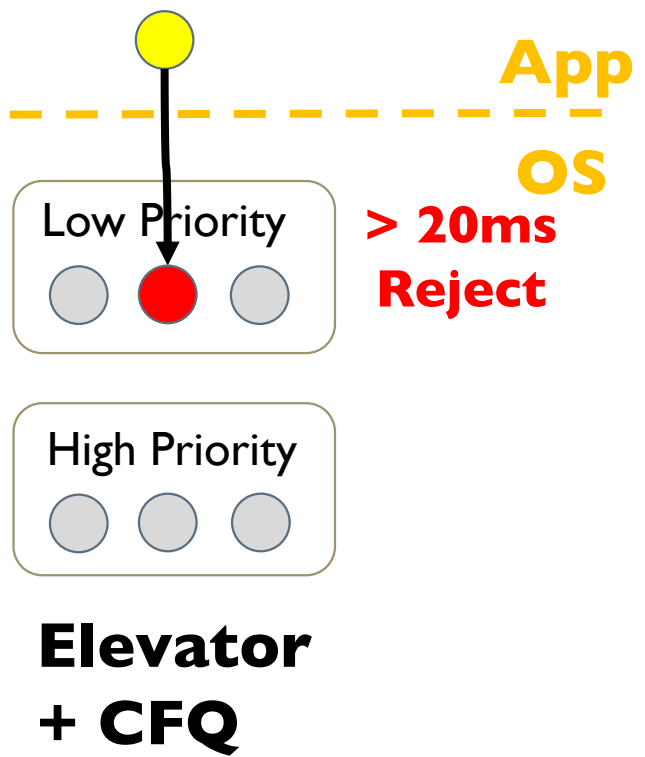
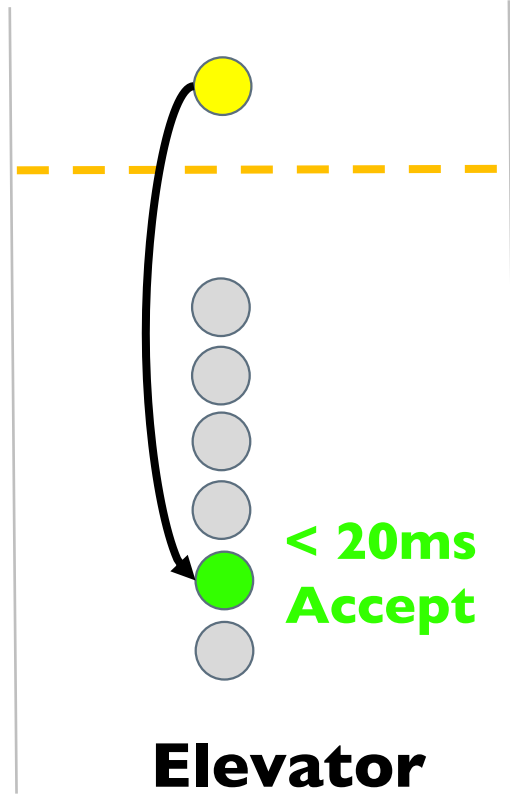
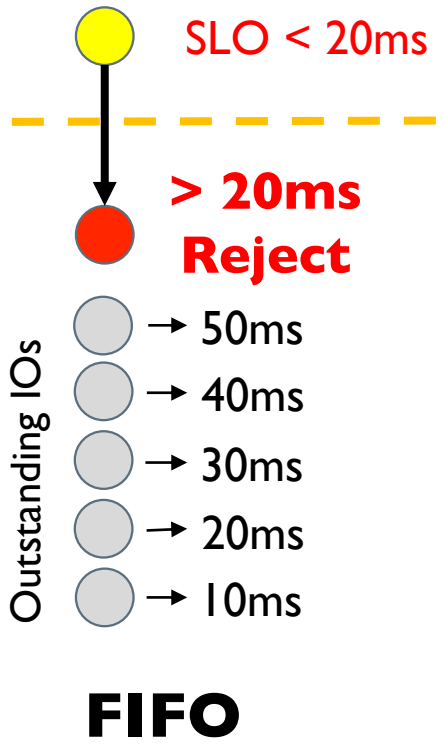


Challenge #1: Modeling Queue Policy

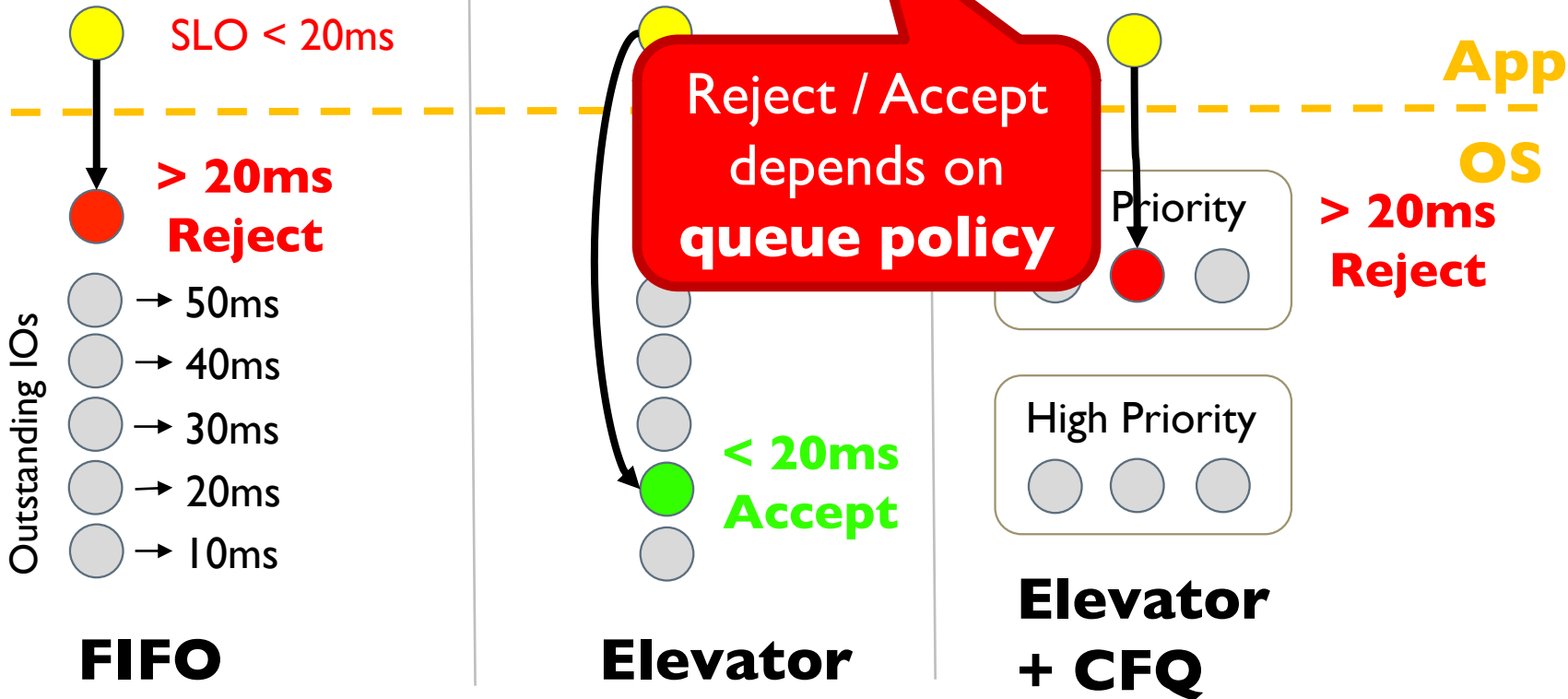
App
OS



Challenge #1: Modeling Queue Policy

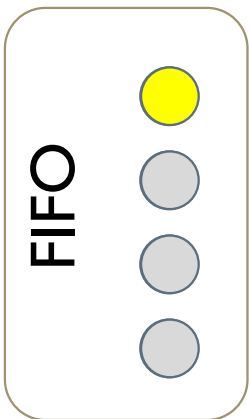


Challenge #1: Modeling Queue Policy

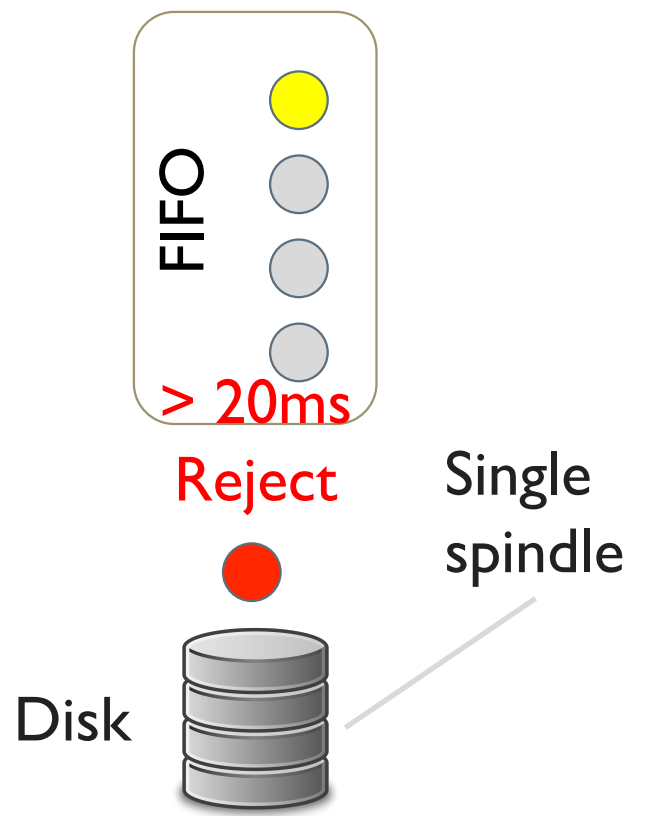




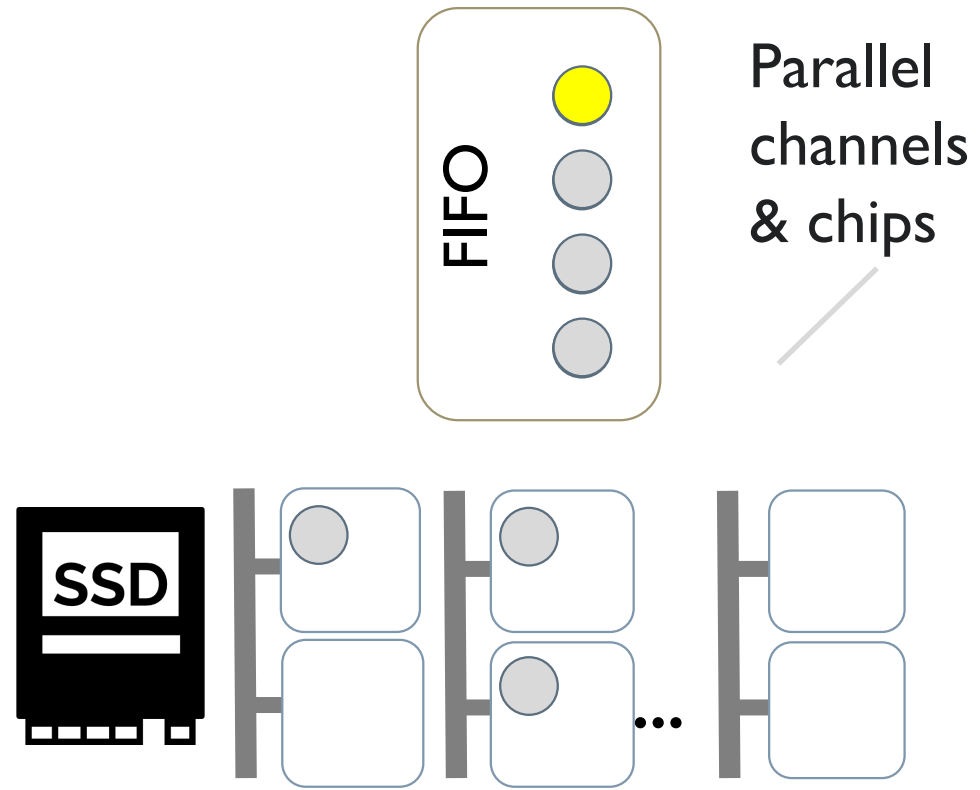
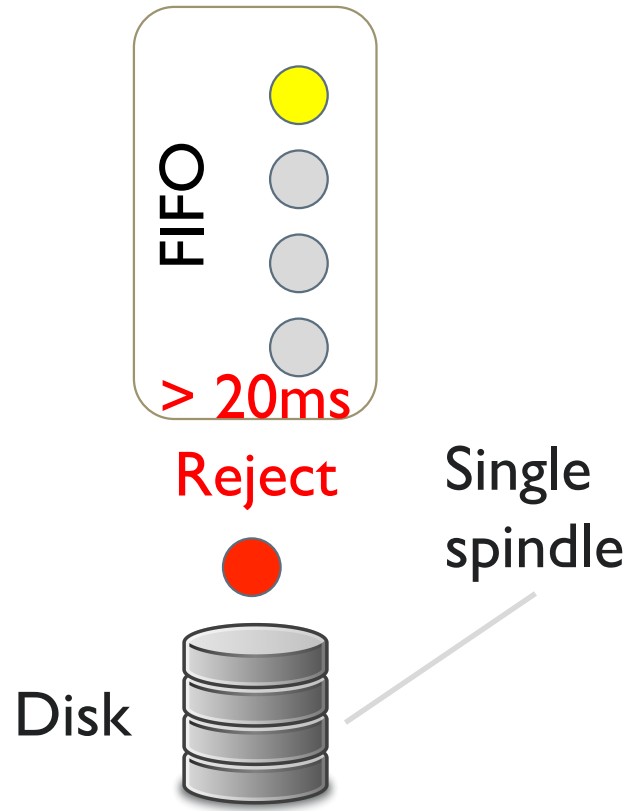
Challenge #2: Device Type



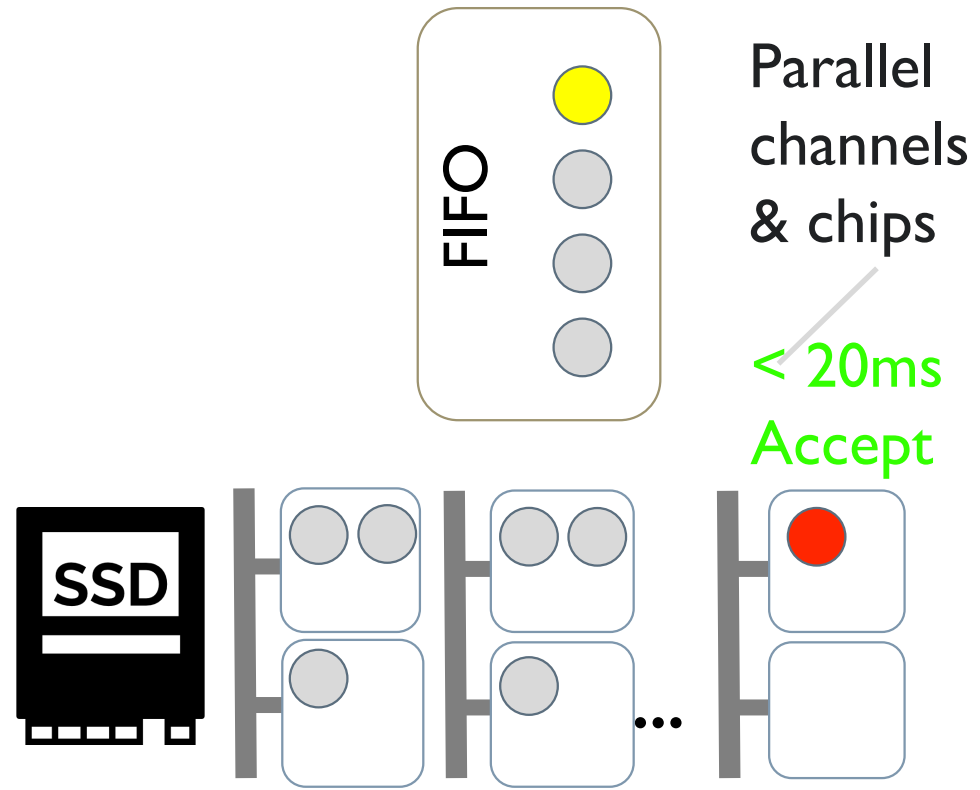
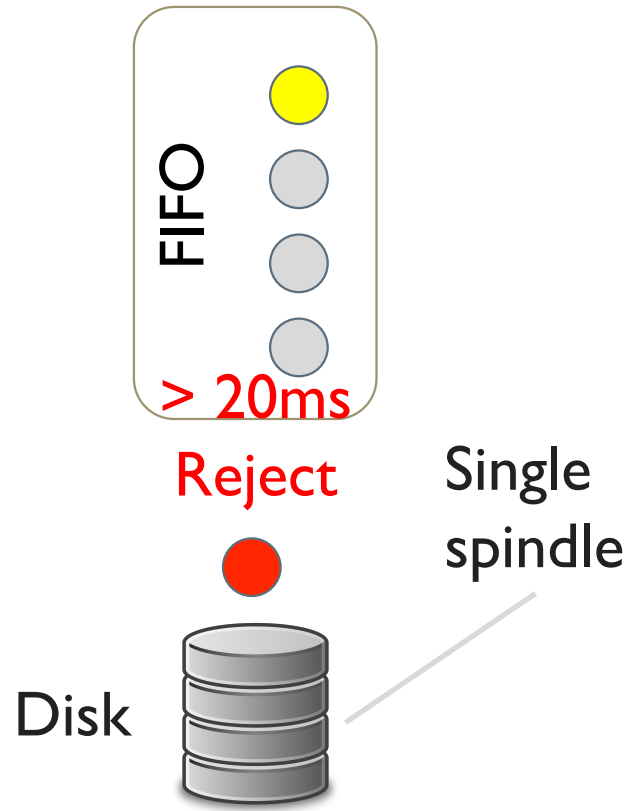
Challenge #2: Device Type



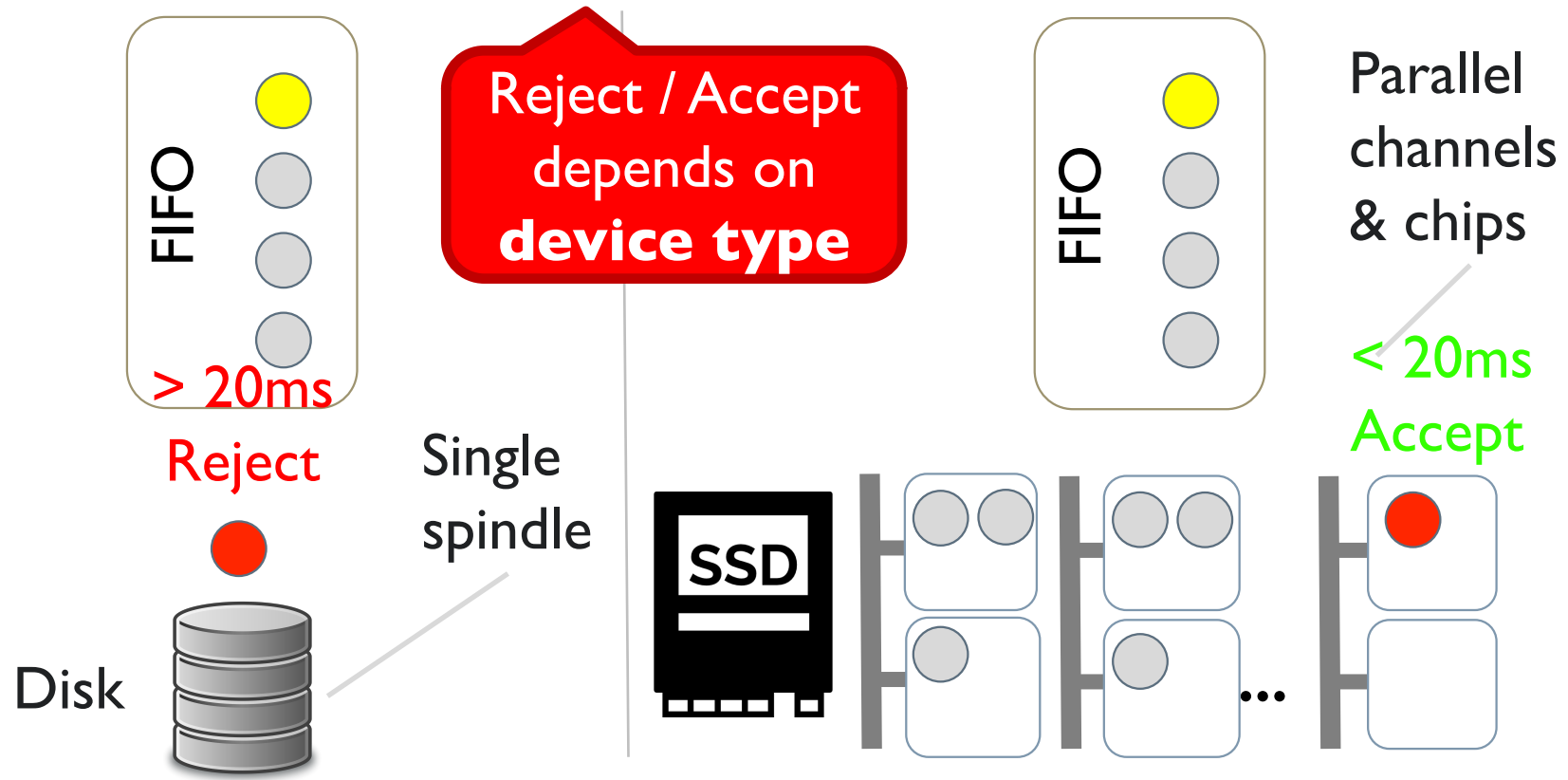
Challenge #2: Device Type



Challenge #2: Device Type

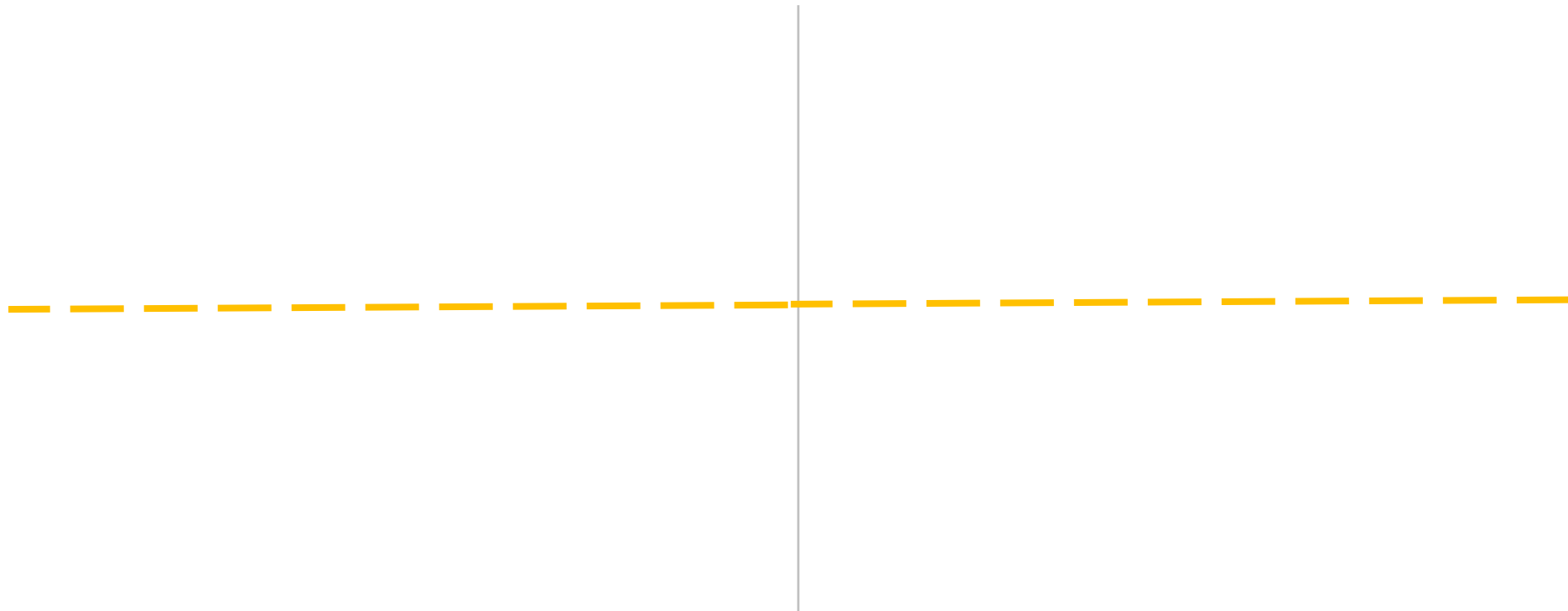


Challenge #2: Device Type



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

IO Offset

Elevator

OS

Too many!
Reject!

● 200

● 700

● 600

● 250



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

IO Offset

Elevator

OS

Too many!
Reject!

● 200

● 700

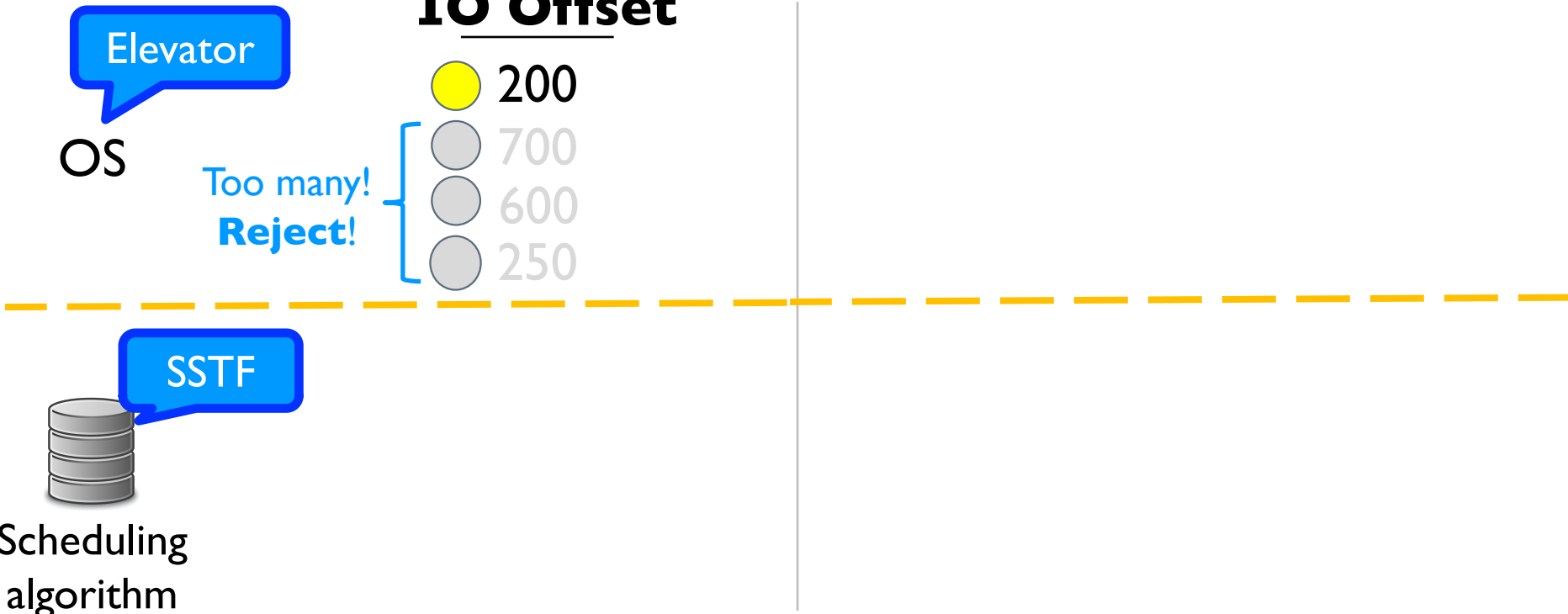
● 600

● 250

SSTF



Scheduling algorithm



Challenge #2: Device Type

Idiosyncrasies of devices are mostly **unrevealed**

IO Offset

Elevator

OS

Too many!
Reject!

- 200
- 700
- 600
- 250

SSTF



Scheduling algorithm

Re-sort,
thus fast,
Accept!

- 700
- 600
- 200
- 250



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

IO Offset

Elevator


OS

Too many!
Reject!

- 200
- 700
- 600
- 250

OS prediction incorrect!

SSTF

 Scheduling algorithm

Re-sort, thus fast,
Accept!

- 700
- 600
- 200
- 250



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

End of queue!
Reject?

IO Offset

Elevator

OS

Too many!
Reject!

- 200
- 700
- 600
- 250

OS prediction incorrect!

-
-
-
-

OS

SSTF



Scheduling algorithm

Re-sort, thus fast, Accept!

- 700
- 600
- 200
- 250



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

IO Offset

Elevator

OS

Too many!
Reject!

- 200
- 700
- 600
- 250

OS prediction incorrect!

SSTF

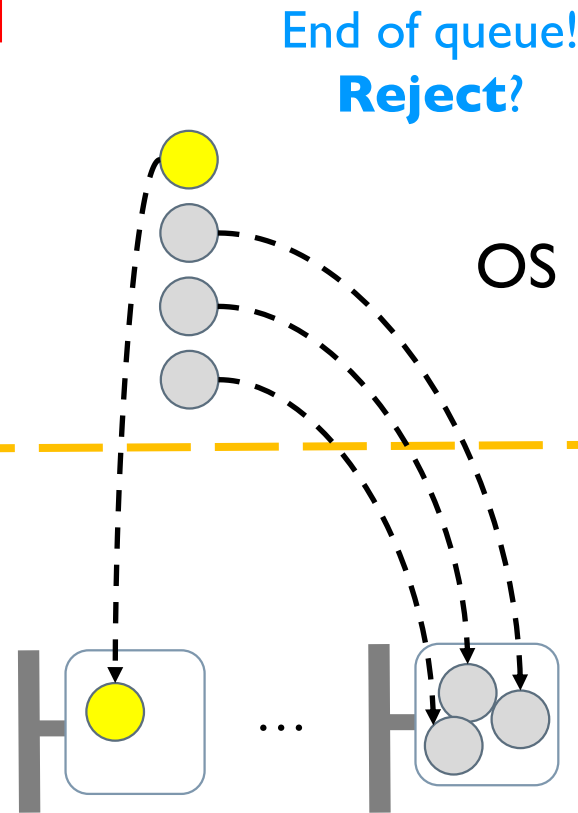
Scheduling algorithm

Re-sort, thus fast,
Accept!

- 700
- 600
- 200
- 250

SSD

Remap to fast chip,
Accept!



Challenge #2: Device Type

Idiosyncrasies of devices are mostly unrevealed

IO Offset

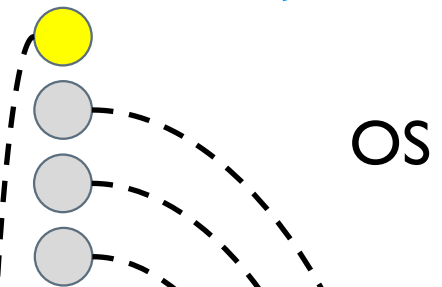
Elevator
OS

Too many!
Reject!

- 200
- 700
- 600
- 250

OS prediction incorrect!

End of queue!
Reject?

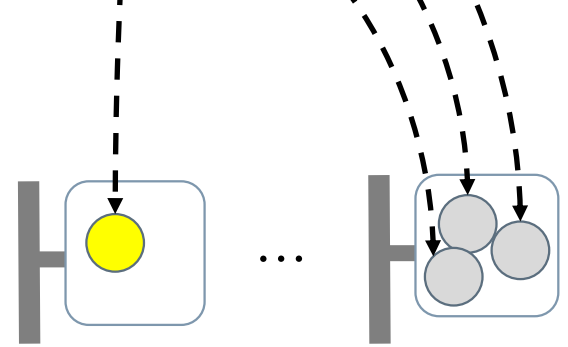


SSTF
Scheduling algorithm

Re-sort, thus fast, Accept!

- 700
- 600
- 200
- 250

SSD
Remap to fast chip, Accept!





Outline

□ Introduction

□ **Design**

- Challenges
- Solutions

□ Evaluation

□ Conclusion

Reject/Latency Prediction

Reject? = $f(\quad)$

Reject/Latency Prediction

Reject? = $f(\text{SLO}, \dots)$

Reject/Latency Prediction

Reject? = $f(\text{SLO}, \text{queue policy}, \dots)$

Reject/Latency Prediction

$$\text{Reject?} = f(\text{SLO}, \text{queue policy}, \text{device type})$$

Reject/Latency Prediction

Reject? = f (SLO, queue policy, device type)

Get from
source-code.
e.g. CFQ, noop



Reject/Latency Prediction

$$\text{Reject?} = f(\text{SLO}, \text{queue policy}, \text{device type})$$

Get from source-code. e.g. CFQ, noop



Simple type

Profiling is enough



Reject/Latency Prediction

$$\text{Reject?} = f(\text{SLO}, \text{queue policy}, \text{device type})$$

Get from source-code.
e.g. CFQ, noop



Simple type

Profiling is enough



Complicated type

White-box knowledge required



Reject/Latency Prediction

$$\text{Reject?} = f(\text{SLO}, \text{queue policy}, \text{device type})$$

Get from source-code.
e.g. CFQ, noop



Simple type

Profiling is enough

MittCFQ



Complicated type

White-box knowledge required

MittSSD

MittCFQ



MittCFQ

Contains user group management

CFQ OS





MittCFQ

Contains user group management
Contains 3 service trees

CFQ

OS



MittCFQ

Contains user group management
Contains 3 service trees
Contains 7 different IO priorities

CFQ OS



MittCFQ

Contains user group management
Contains 3 service trees
Contains 7 different IO priorities
Contains ~4500 LOC

CFQ

OS



MittCFQ

Open-sourced

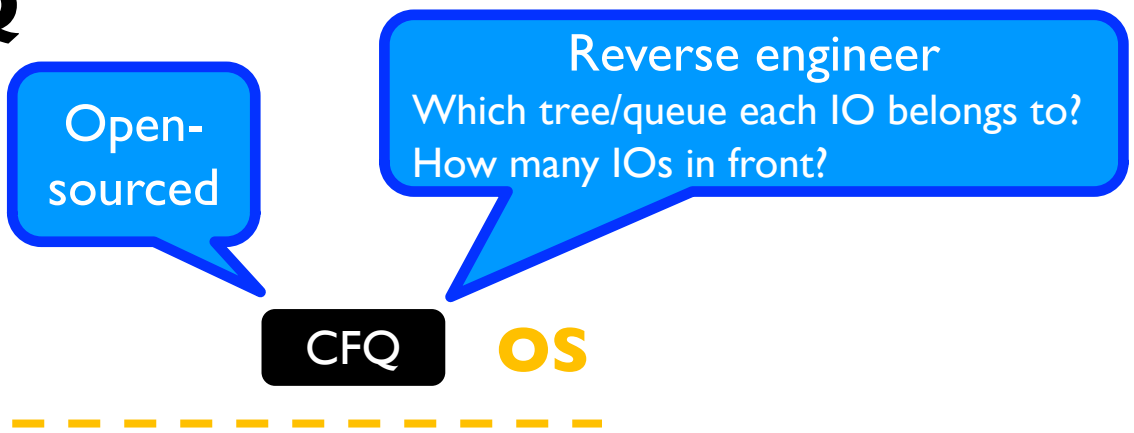
Contains user group management
Contains 3 service trees
Contains 7 different IO priorities
Contains ~4500 LOC

CFQ

OS

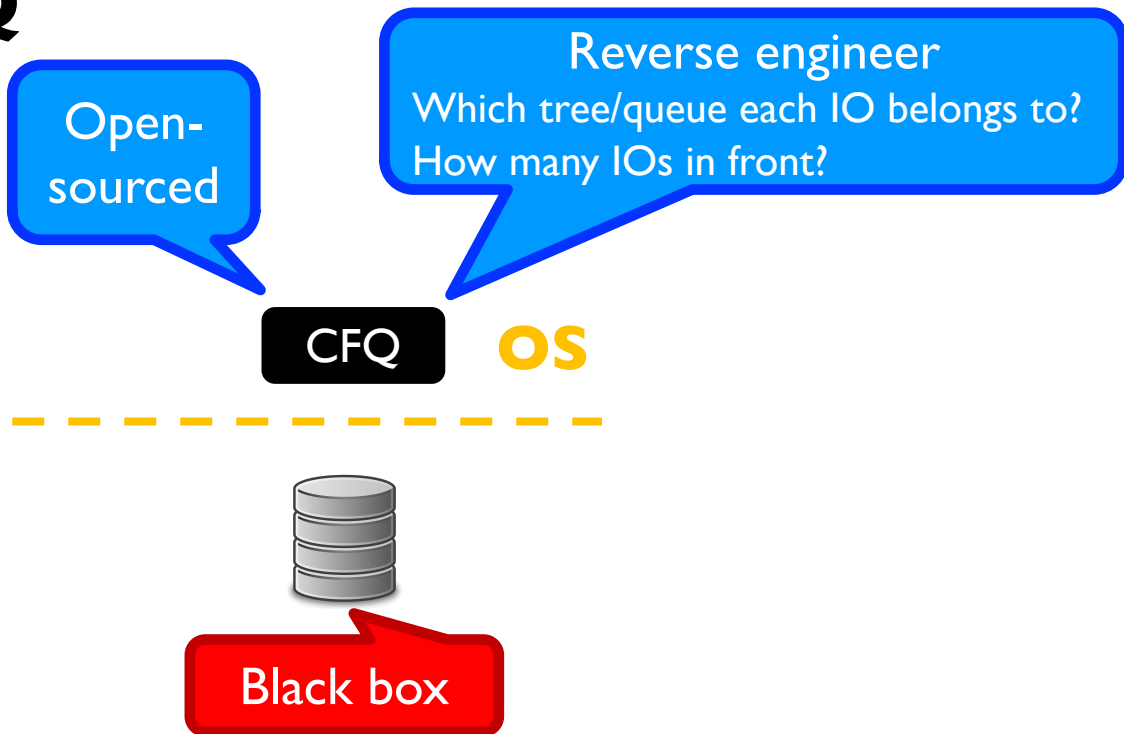


MittCFQ



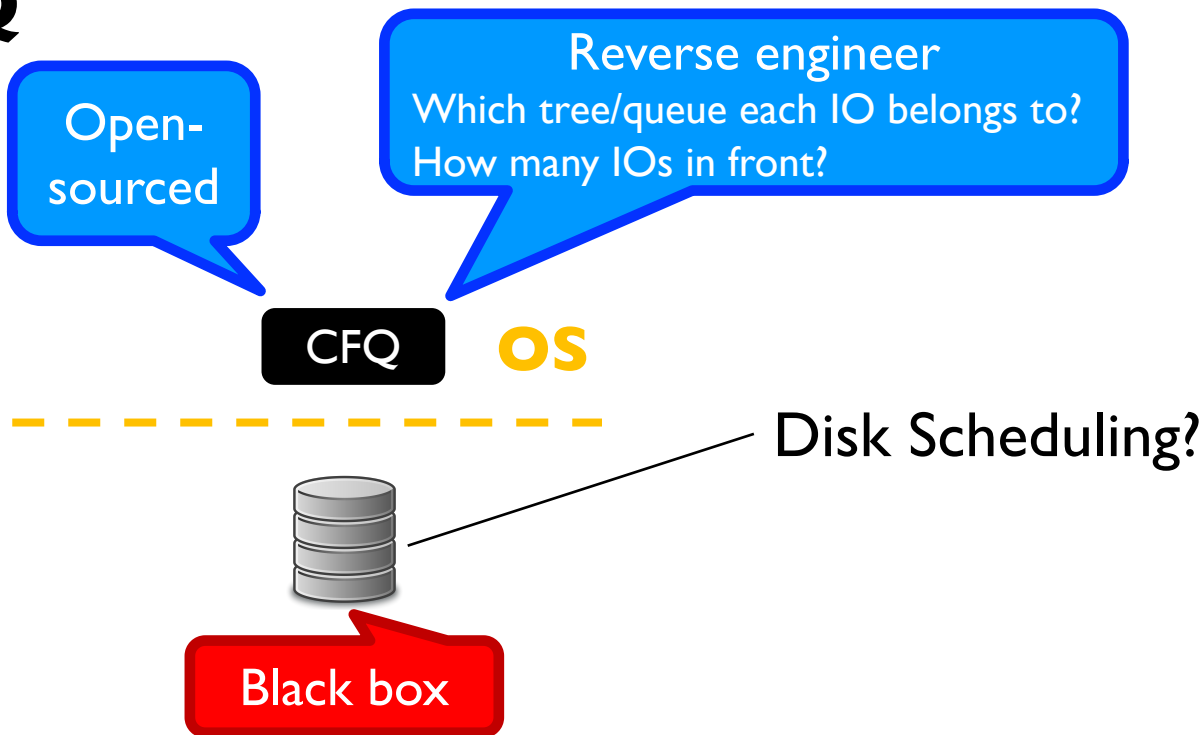


MittCFQ

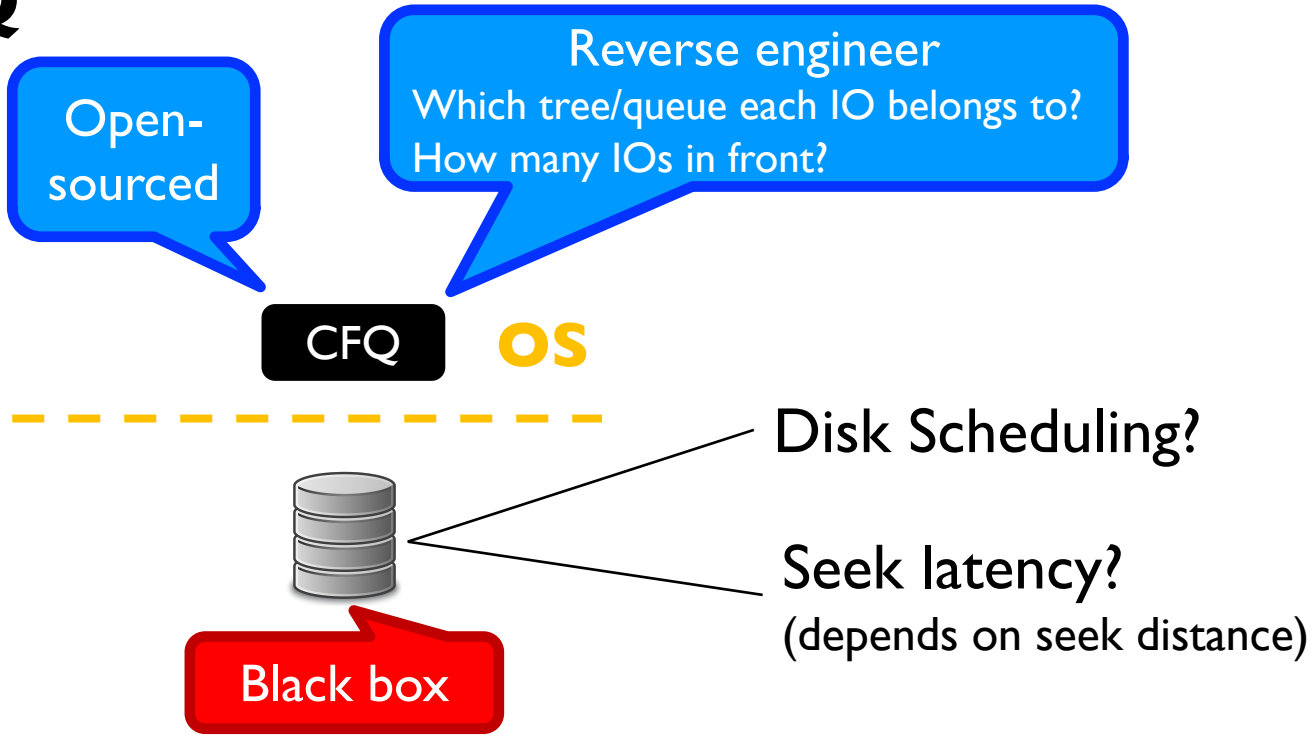




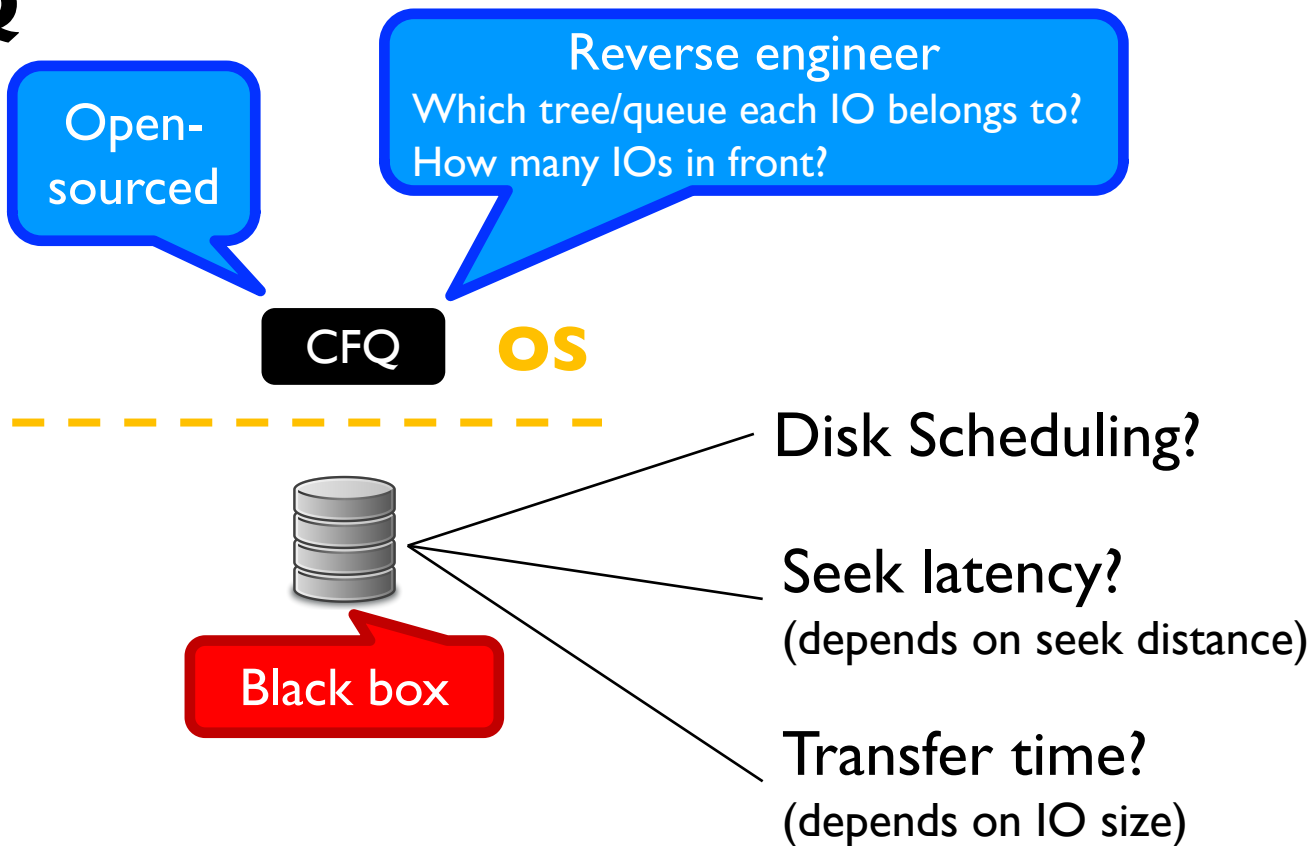
MittCFQ



MittCFQ



MittCFQ





For each interval in [100MB, 200MB, ..., 1GB] do:

```
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {  
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {  
    for (size = 0; size < maxSize; size += sizeInterval){  
      start_ts = gettimeofday();  
      seek (startOffset);  
      read (endOffset, size);  
      end_ts = gettimeofday();  
      latency = start_ts - end_ts;  
      print (endOffset - startOffset, size, latency);  
    }  
  }  
}
```

MittCFQ Profiling



```
For each interval in [ 100MB, 200MB, ..., 1GB ] do:
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek (startOffset) Random seek
      read (endOffset, size);
      end_ts = gettimeofday();
      latency = start_ts - end_ts;
      print (endOffset - endOffset, size, latency);
    }
  }
}
```

MittCFQ Profiling



For each interval in [100MB, 200MB, ..., 1GB] do:

```
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {  
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {  
    for (size = 0; size < maxSize; size += sizeInterval){  
      start_ts = gettimeofday();  
      seek (startOffset) Random seek  
      read (endOffset, size); Random read  
      end_ts = gettimeofday();  
      latency = start_ts - end_ts;  
      print (endOffset - endOffset, size, latency);  
    }  
  }  
}
```

MittCFQ Profiling



For each interval in [100MB, 200MB, ..., 1GB] do:

```
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {  
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {  
    for (size = 0; size < maxSize; size += sizeInterval){  
      start_ts = gettimeofday();  
      seek (startOffset) Random seek  
      read (endOffset, size); Random read  
      end_ts = gettimeofday();  
      latency = start_ts - end_ts; Collect latency  
      print (endOffset - endOffset, size, latency);  
    }  
  }  
}
```

MittCFQ Profiling



MittCFQ Profiling

For each interval in [100MB, 200MB, ..., 1GB] do:

```
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {  
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {  
    for (size = 0; size < maxSize; size += sizeInterval){  
      start_ts = gettimeofday();  
      seek(startOffset) Random seek  
      read(endOffset, size) Random read  
      end_ts = gettimeofday();  
      latency = start_ts - end_ts; Collect latency  
      print (endOffset - endOffset, size, latency);  
    }  
  }  
}
```

2 disk models
11-hour profiling



MittCFQ Profiling

For each interval in [100MB, 200MB, ..., 1GB] do:

```

for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek(startOffset)
      read(endOffset, size);
      end_ts = gettimeofday();
      latency = start_ts - end_ts;
      print (endOffset, startOffset, size, latency);
    }
  }
}

```

seek(startOffset)

Random seek

read(endOffset, size);

Random read

latency = start_ts - end_ts;

Collect latency

scikit-learn

Linear Regression

2 disk models
11-hour profiling

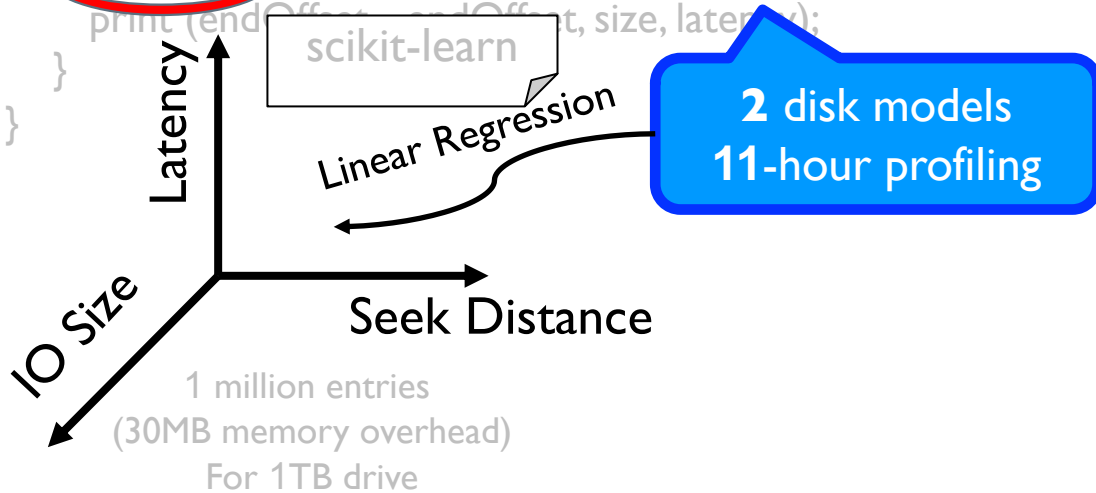
MittCFQ Profiling

For each interval in [100MB, 200MB, ..., 1GB] do:

```

for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek(startOffset) Random seek
      read(endOffset, size); Random read
      end_ts = gettimeofday();
      latency = start_ts - end_ts; Collect latency
      print(endOffset, startOffset, size, latency);
    }
  }
}

```

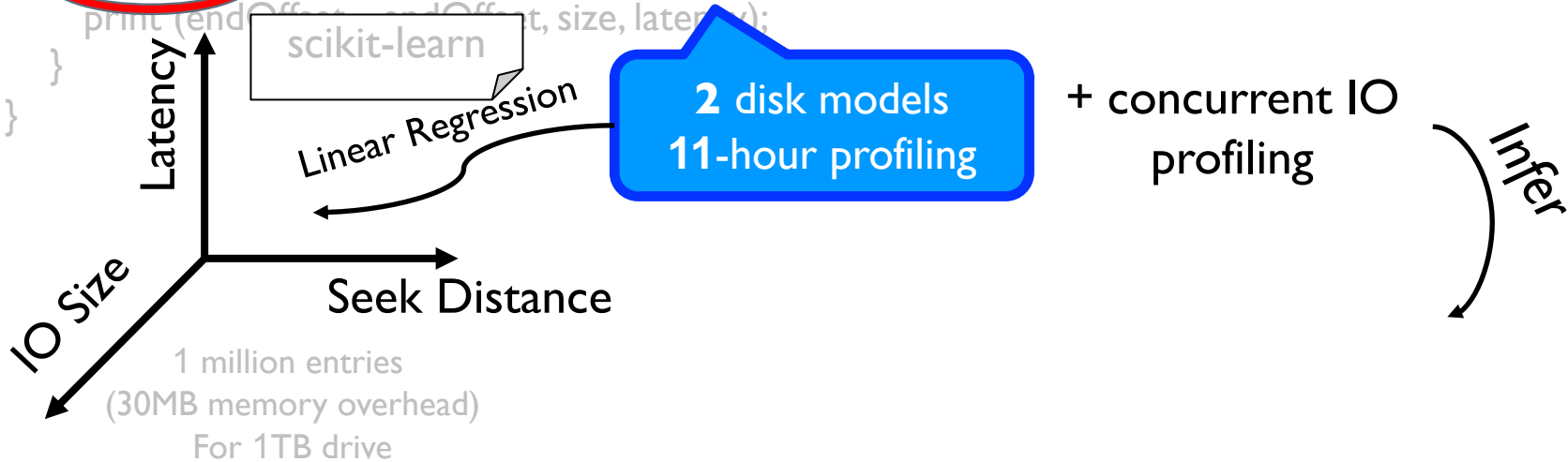


```

For each interval in [ 100MB, 200MB, ..., 1GB ] do:
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek(startOffset)
      read(endOffset, size);
      end_ts = gettimeofday();
      latency = start_ts - end_ts;
      print (endOffset, startOffset, size, latency);
    }
  }
}
  
```

Random seek
Random read
Collect latency

MittCFQ Profiling



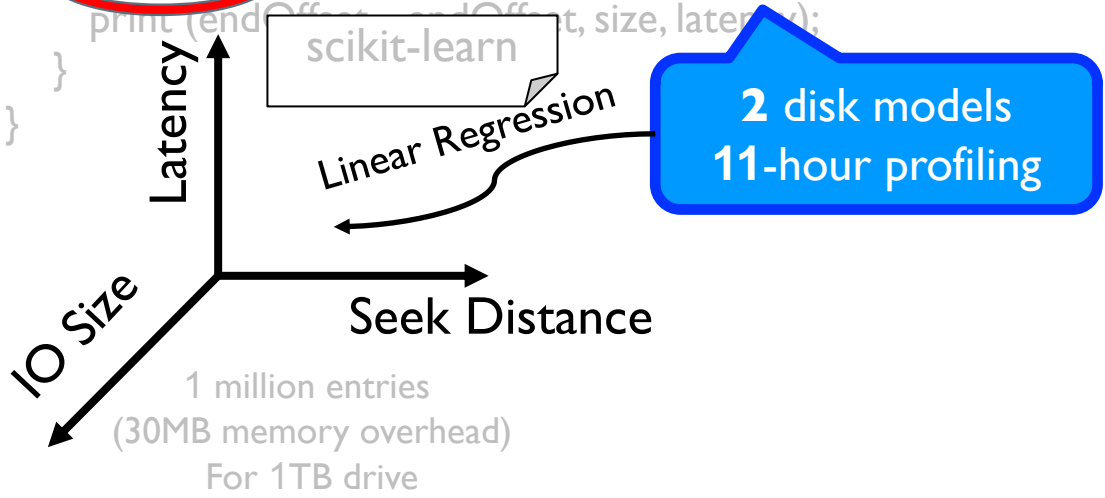
```

For each interval in [ 100MB, 200MB, ..., 1GB ] do:
for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek(startOffset);
      read(endOffset, size);
      end_ts = gettimeofday();
      latency = start_ts - end_ts;
      print(endOffset, startOffset, size, latency);
    }
  }
}

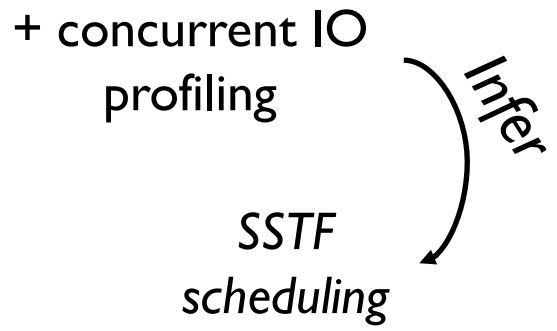
```

Random seek
Random read

Collect latency



MittCFQ Profiling





For each interval in [100MB, 200MB, ..., 1GB] do:

```

for (startOffset = 0; startOffset < maxOffset; startOffset += interval) {
  for (endOffset = 0; endOffset < maxOffset; endOffset += interval) {
    for (size = 0; size < maxSize; size += sizeInterval){
      start_ts = gettimeofday();
      seek(startOffset)
      read(endOffset, size);
      end_ts = gettimeofday();
      latency = start_ts - end_ts;
      print(endOffset, startOffset, size, latency);
    }
  }
}

```

seek(startOffset)

Random seek

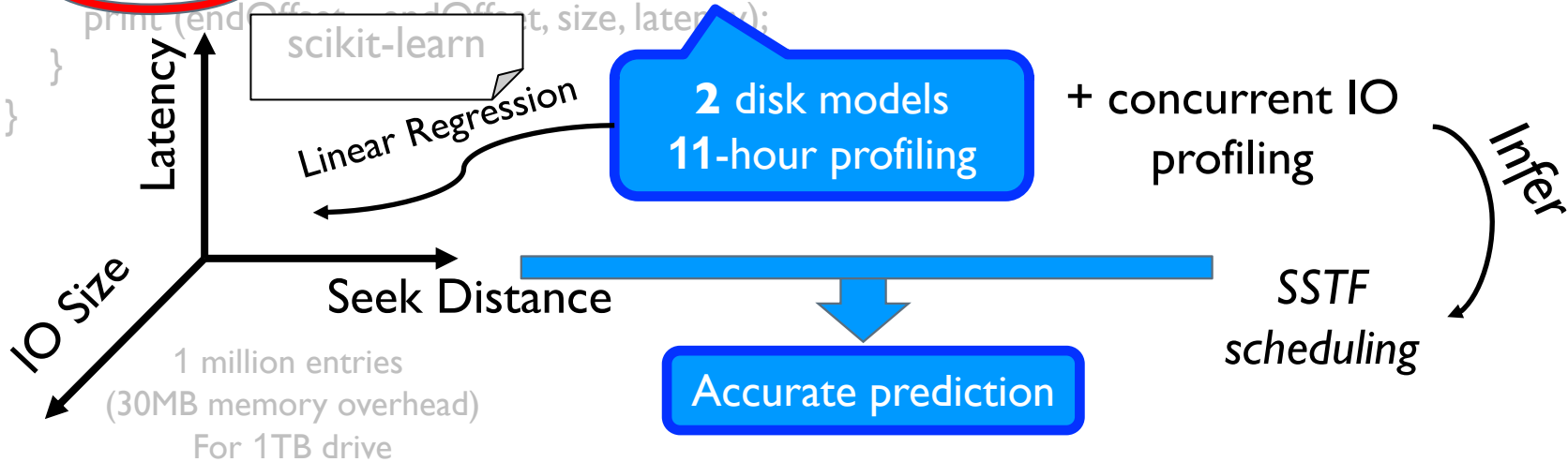
read(endOffset, size);

Random read

latency = start_ts - end_ts;

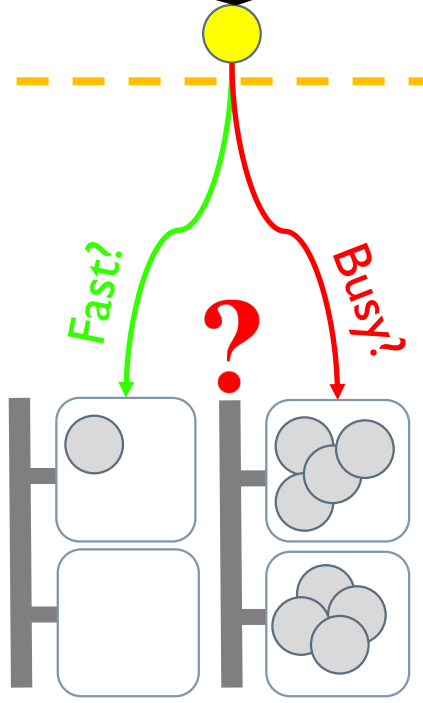
Collect latency

MittCFQ Profiling



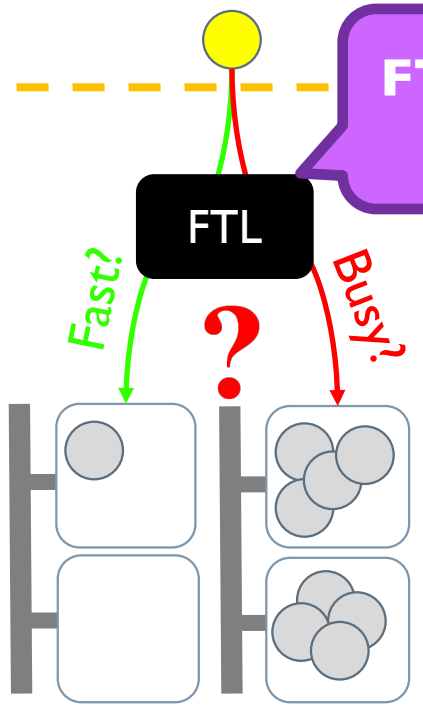
MittSSD OS

Which channel/chip?
Fast? Busy?



MittSSD

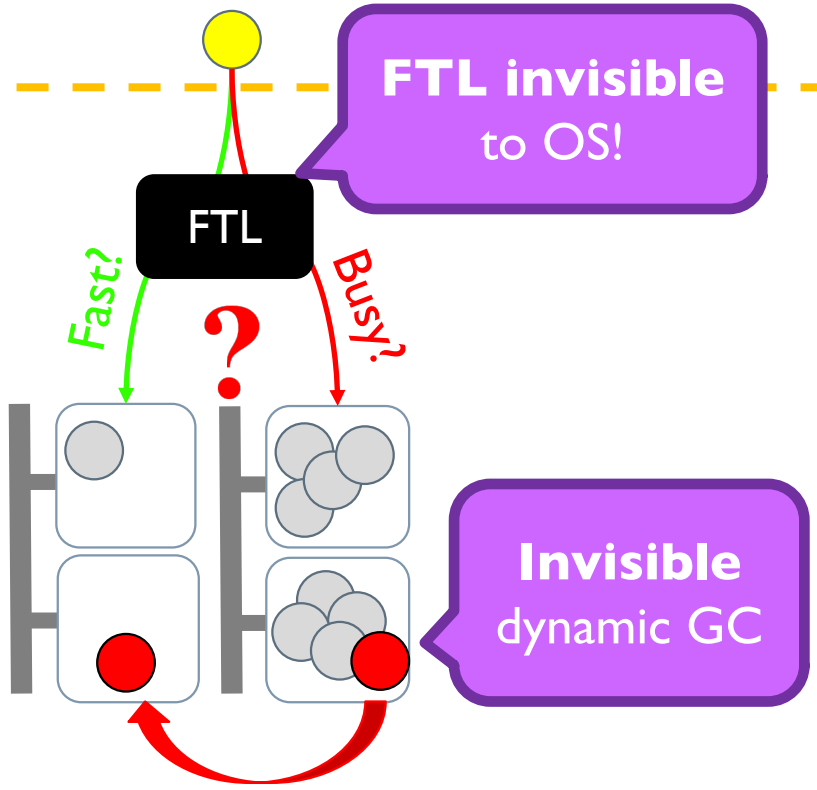
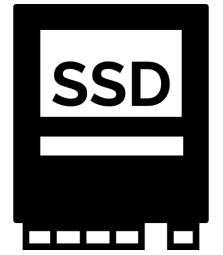
OS



FTL invisible to OS!

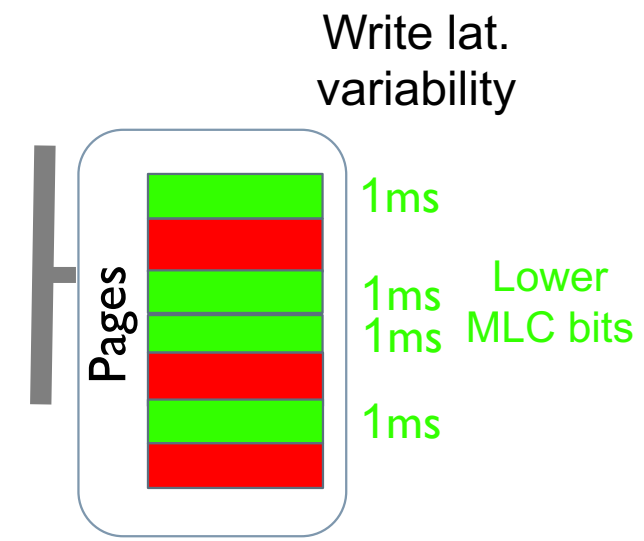
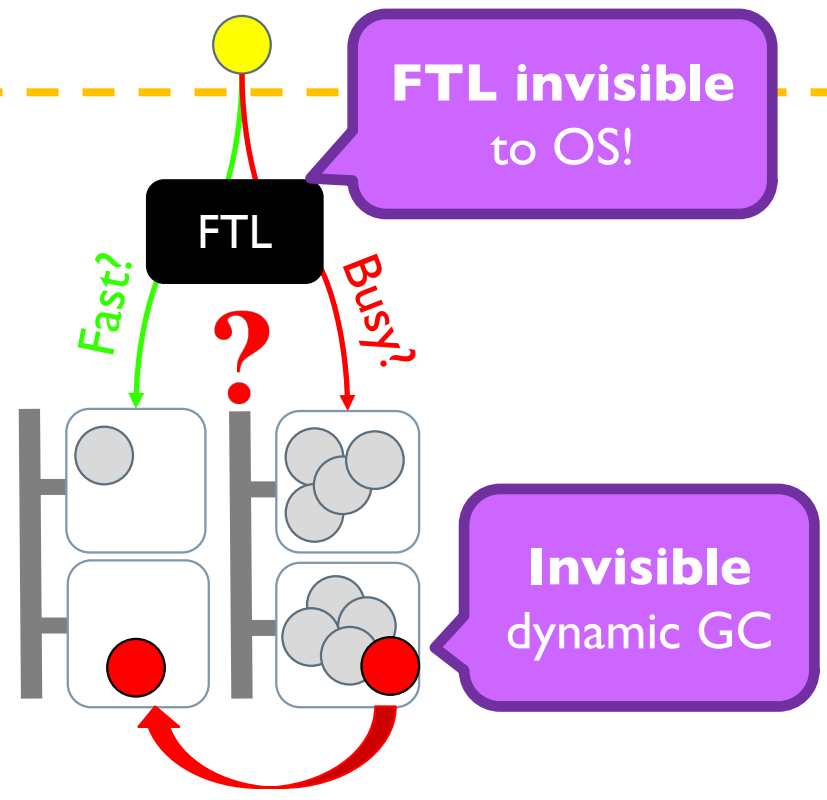
MittSSD

OS



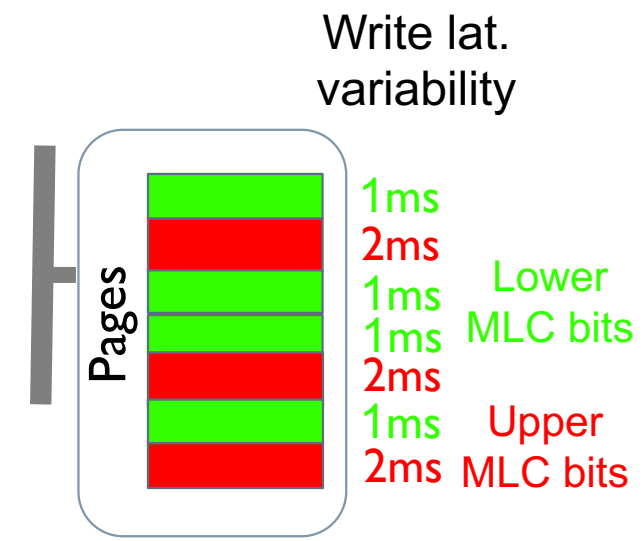
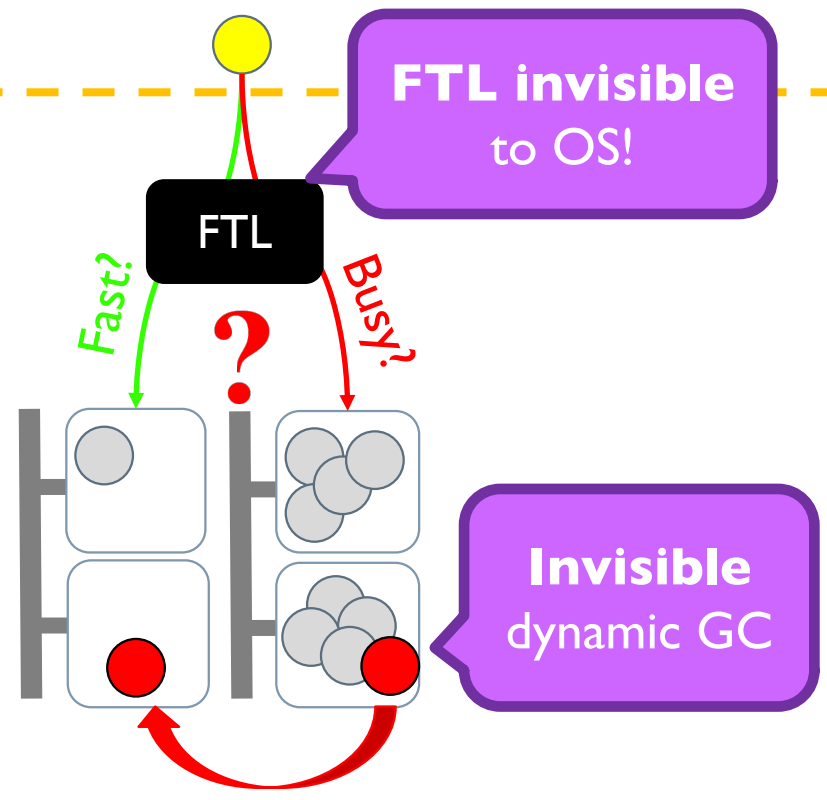
MittSSD

OS



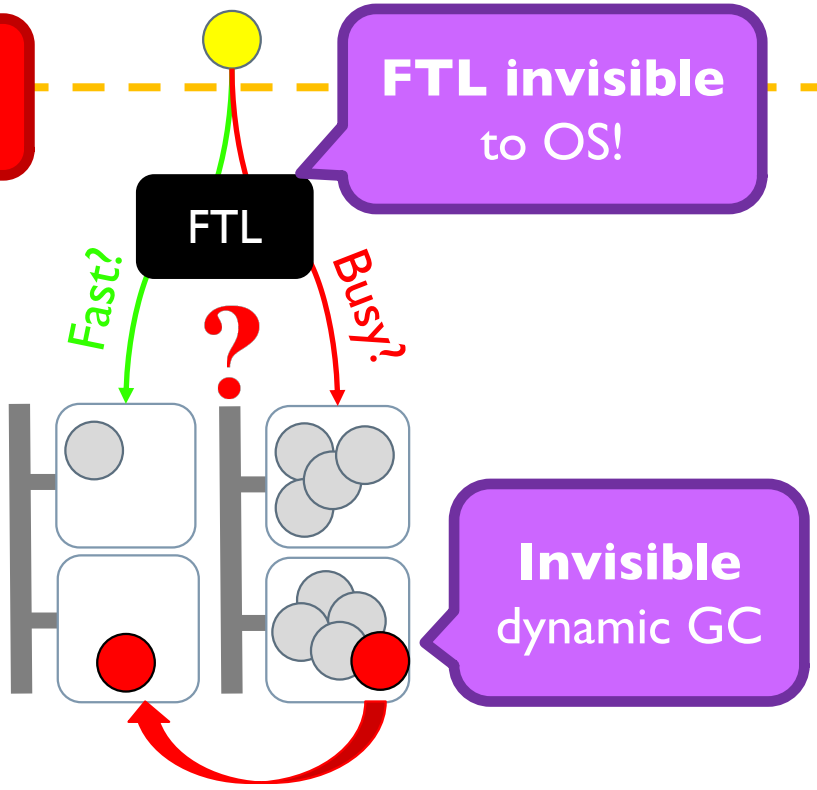
MittSSD

OS

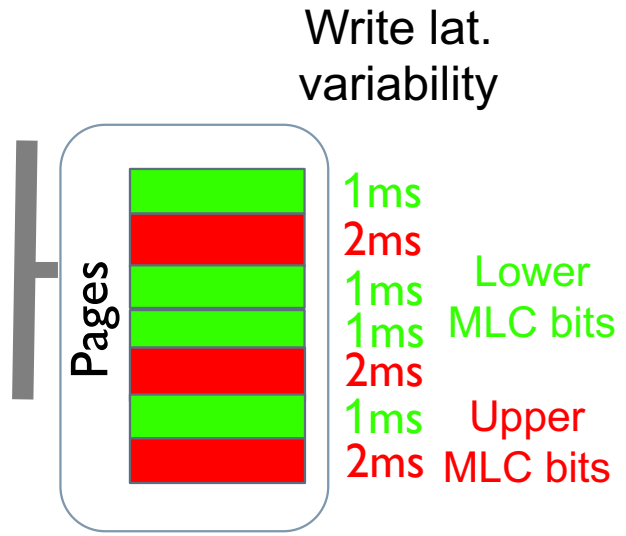


MittSSD

Too complex to model!



FTL invisible to OS!



MittSSD

OS



MittSSD

Software-defined flash

OS



MittSSD

Software-defined flash

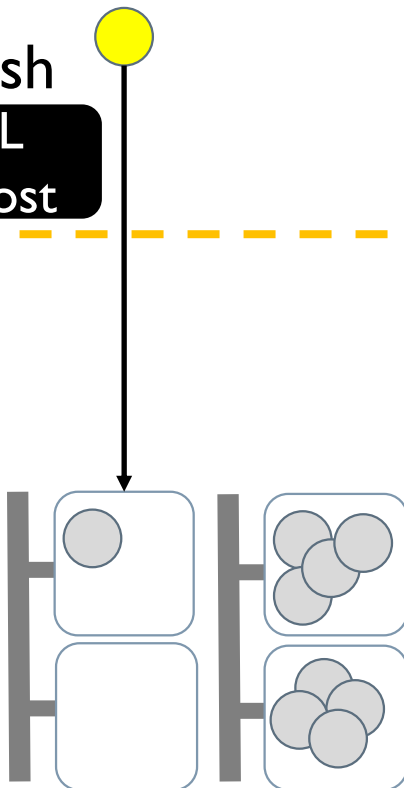
LightNVM

FTL
at host

OS



Open-Channel
SSD



MittSSD

Software-defined flash

LightNVM

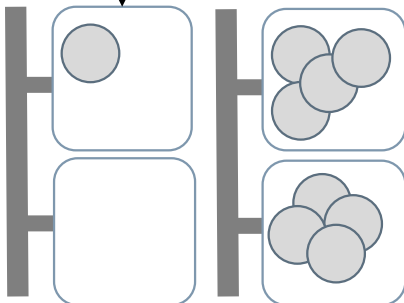
FTL
at host

OS knows where
IOs are mapped

OS



Open-Channel
SSD



MittSSD

Software-defined flash

LightNVM

FTL
at host

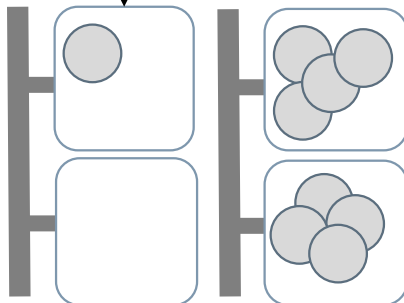


OS knows where
IOs are mapped

OS



Open-Channel
SSD



OS can track
every single IO

MittSSD

Software-defined flash

LightNVM

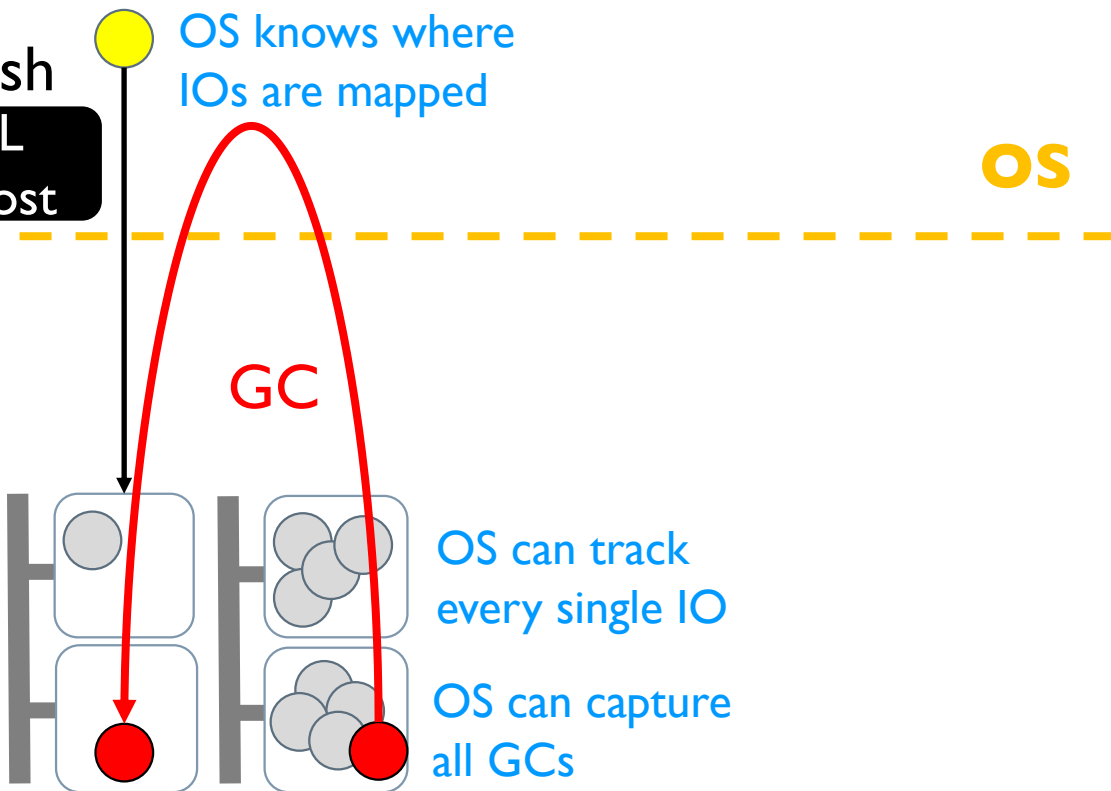
FTL
at host

OS knows where
IOs are mapped

OS



Open-Channel
SSD



MittSSD

Software-defined flash

LightNVM

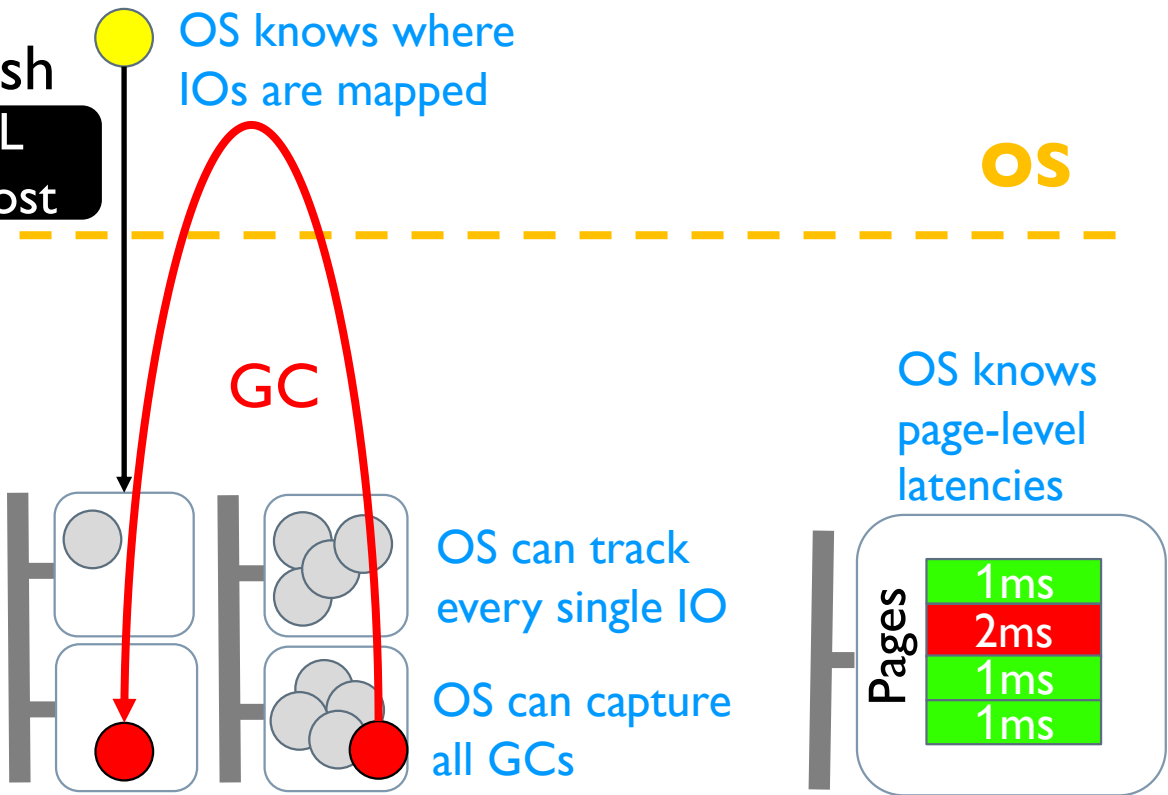
FTL at host

OS knows where IOs are mapped

OS



Open-Channel SSD



MittSSD

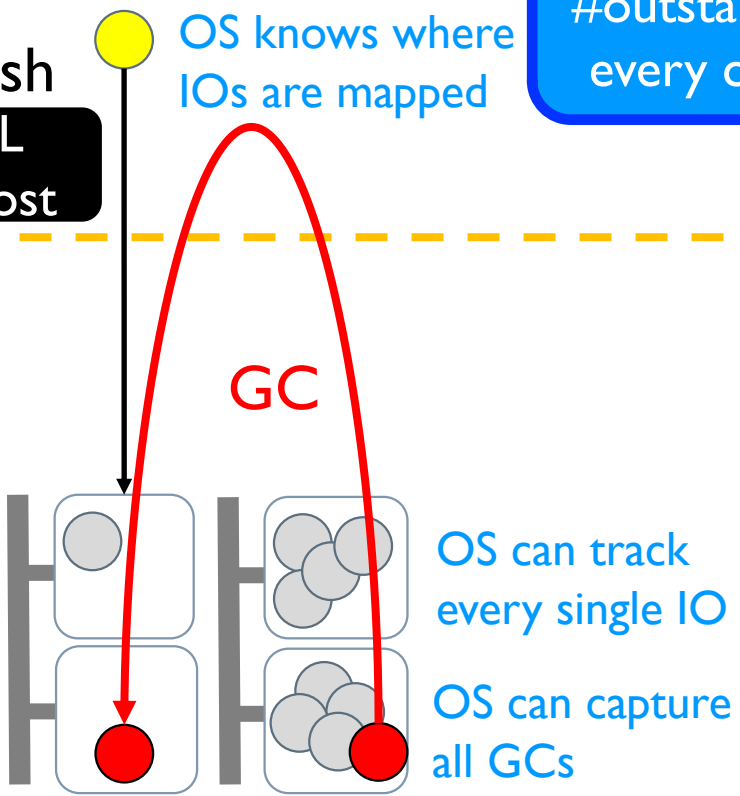
Software-defined flash

LightNVM

FTL at host



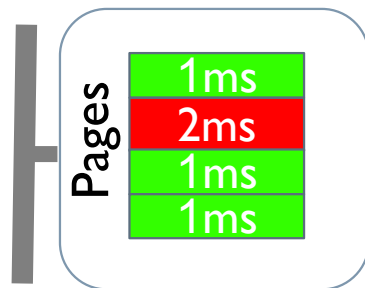
Open-Channel SSD



OS can see #outstanding IOs to every chip/channel

OS

OS knows page-level latencies



MittSSD

Software-defined flash
LightNVM FTL at host



Open-Channel SSD

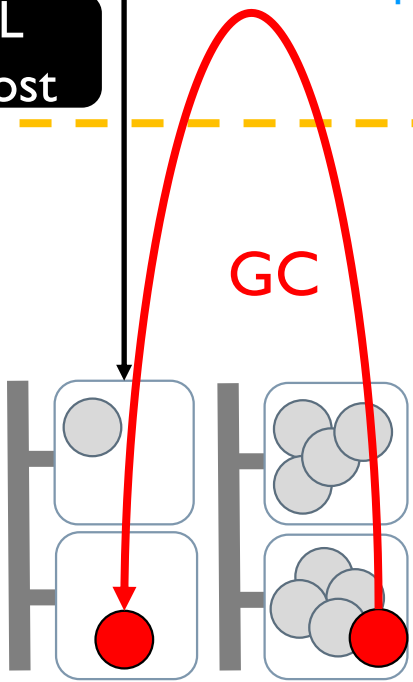


OS knows where IOs are mapped

OS can see #outstanding IOs to every chip/channel

Accurate prediction

OS

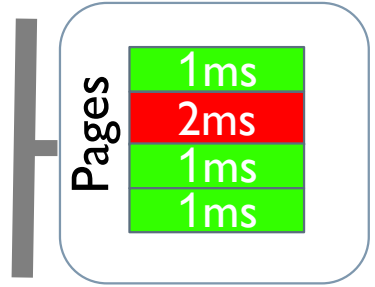


GC

OS can track every single IO

OS can capture all GCs

OS knows page-level latencies



Pages

1ms
2ms
1ms
1ms

Reject/Latency Prediction

$$\text{Reject?} = f(\text{SLO}, \text{queue policy}, \text{device type})$$

Reverse engineering based on source code



Simple type

Profiling is enough

device type

Complicated type

White-box knowledge required



Other Solved Challenges

Other Solved Challenges

- Prediction overhead optimizations

Other Solved Challenges

- Prediction overhead optimizations
 - Avoids going through every IO in the queue

Other Solved Challenges

- Prediction overhead optimizations
 - Avoids going through every IO in the queue
 - Reduces overhead from $O(n)$ to roughly $O(1)$

Other Solved Challenges

- Prediction overhead optimizations
 - Avoids going through every IO in the queue
 - Reduces overhead from $O(n)$ to roughly $O(1)$
 - Shows $< 5\mu\text{s}$ overhead for MittCFQ prediction
 - $< 300\text{ns}$ for MittSSD prediction

Other Solved Challenges

- Prediction overhead optimizations
 - Avoids going through every IO in the queue
 - Reduces overhead from $O(n)$ to roughly $O(1)$
 - Shows $< 5\mu\text{s}$ overhead for MittCFQ prediction
 - $< 300\text{ns}$ for MittSSD prediction
- MittCache
 - Prediction for OS Cache

Other Solved Challenges

- Prediction overhead optimizations
 - Avoids going through every IO in the queue
 - Reduces overhead from $O(n)$ to roughly $O(1)$
 - Shows $< 5\mu s$ overhead for MittCFQ prediction
 - $< 300ns$ for MittSSD prediction
- MittCache
 - Prediction for OS Cache

MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface

Mingzhe Hao, Huaicheng Li, Michael Hao Tong, Chrisma Pakha, Riza O. Suminto, Cesar A. Stuardo, Andrew A. Chien, and Haryadi S. Gunawi

University of Chicago

ABSTRACT
 In this paper, we address the challenge of supporting millisecond-level tail latencies for data-parallel applications. In MittOS, we advocate a new principle that operating systems should not only reject requests that cannot be promptly served, but also fast reject requests that are not expected to finish within the SLO. If an application can predict the SLO (e.g., IO deadlines), if MittOS predicts that the IO SLOs cannot be met, MittOS will promptly return EBUSY signal, allowing the application to reallocate resources to other requests. This approach helps reduce IO completion time up to 35% compared to wait-then-speculate approaches.

CCS CONCEPTS

• Computer systems organization → Real-time operating systems; Distributed architectures;

KEYWORDS

Data-parallel frameworks, low latency, operating system, performance, real-time, SLO, tail tolerance.

ACM Reference Format:

Mingzhe Hao, Huaicheng Li, Michael Hao Tong, Chrisma Pakha, Riza O. Suminto, Cesar A. Stuardo, Andrew A. Chien, and Haryadi S. Gunawi. 2017. MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface. In *ACM SIGOPS 26th Symposium on Operating Systems Principles*. Shanghai, China, October 28-31, 2017, 16 pages.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
 SOSP '17, October 28, 2017, Shanghai, China
 © 2017 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.
 ACM ISBN 978-1-4503-5085-3/17/10...\$15.00
<https://doi.org/10.1145/3132747.3132774>

1 INTRODUCTION

Low and stable latency is a critical key to the success of many services, but variable load and resource sharing common in cloud environments introduces resource contention that in turn produces “the long tail” of requests that take hundreds of seconds) [20], where there is sufficient time to wait, observe, and launch extra speculative tasks if necessary. Such a “wait-then-speculate” method has proven to be highly effective; many variants of this technique have been proposed and put into wide production use. More challenging are applications that generate millions of small requests, each expected to finish in milliseconds. For these, techniques that “wait-then-speculate” are ineffective, as the time to detect a problem is comparable to the delay caused by it.

One approach to this challenge is *pre-emptive scheduling*, where every request is served multiple times, the first to respond is used to fulfill this request, and the others however *doubles* the IO intensity. To reduce extra load, applications can delay the duplicate request and cancel the clone when a response is received (a “*tied requests*”) [19]. To achieve this, IO queuing and rejection management can be *built* in the application layer. One alternative is “*hedged requests*”, where a duplicate request is sent after the first request is outstanding for more than, for example, the 95th-percentile expected latency; but the slow requests (5% *must wait* before being retried. Finally, “*snitching*” [1, 52] – the application monitoring request latency and picking the fastest replica – can be employed; however, such techniques are *ineffective* if noise is bursty.

All of the techniques discussed above attempt to minimize tail in the *absence* of information about underlying resource business. While the OS layer may have such information, it is *hidden* and *unexposed*. A prime example is the `read()` interface that returns either success or error. However, when resources are busy (disk contention from other tenants, device garbage collection, etc.), a `read()` can be stalled inside the OS for some time. Currently, the OS does not have a direct way to indicate that a request may take a long time, nor is there a way for applications to indicate they would like “to know the OS is busy.”

Please refer to the paper!



Outline

- Introduction
- Design
- **Evaluation**
 - Tail reduction
 - Latency prediction accuracy
- Conclusion

MittCFQ-powered MongoDB



MittCFQ-powered MongoDB



Physical node #1



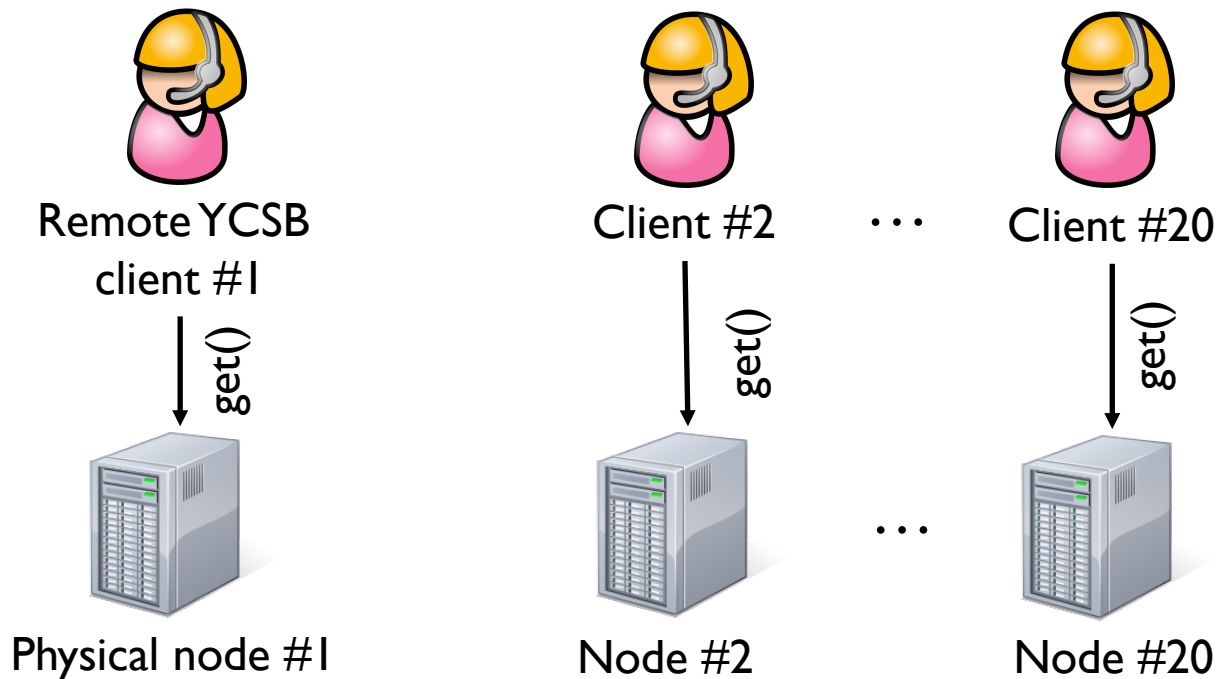
Node #2

...

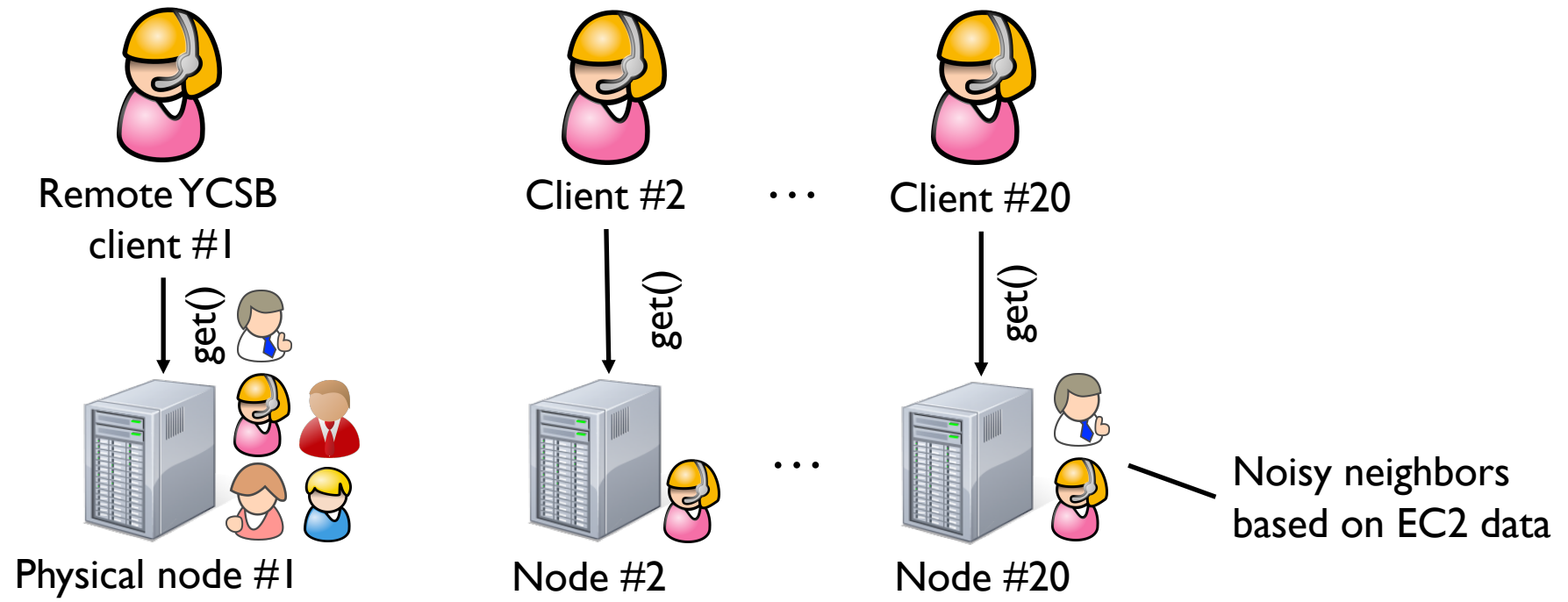


Node #20

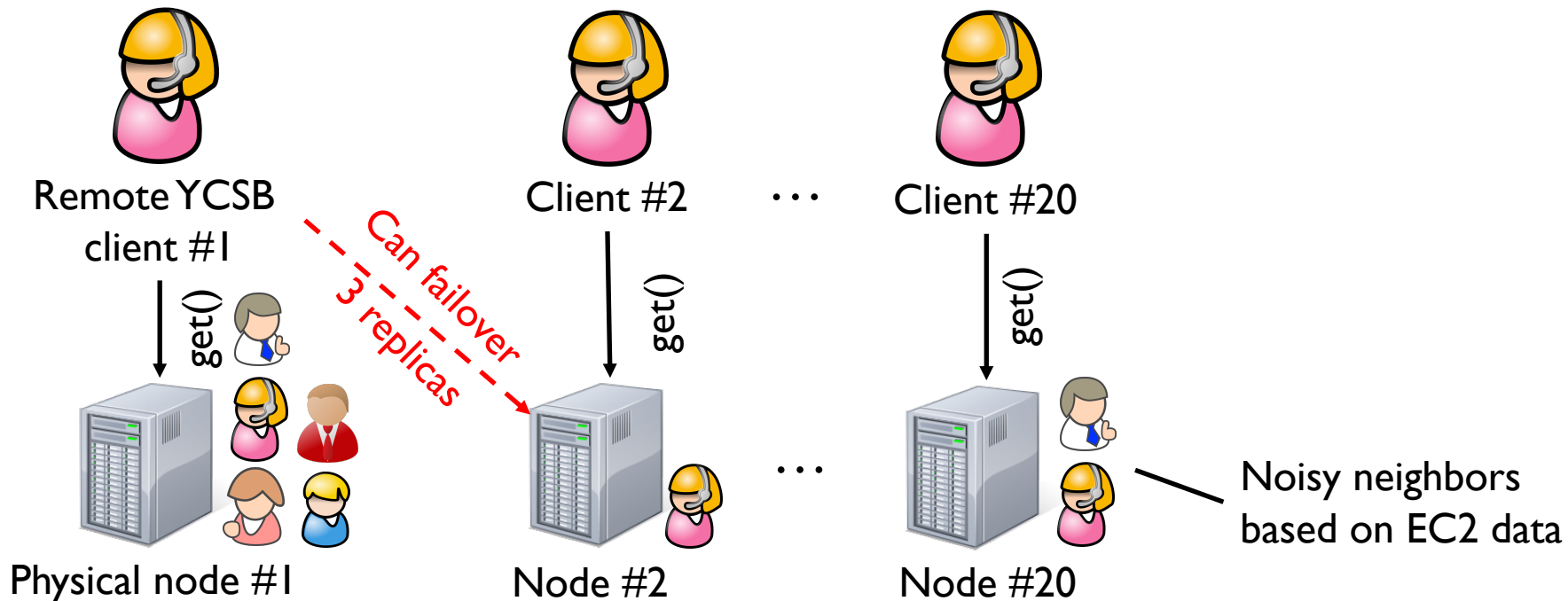
MittCFQ-powered MongoDB



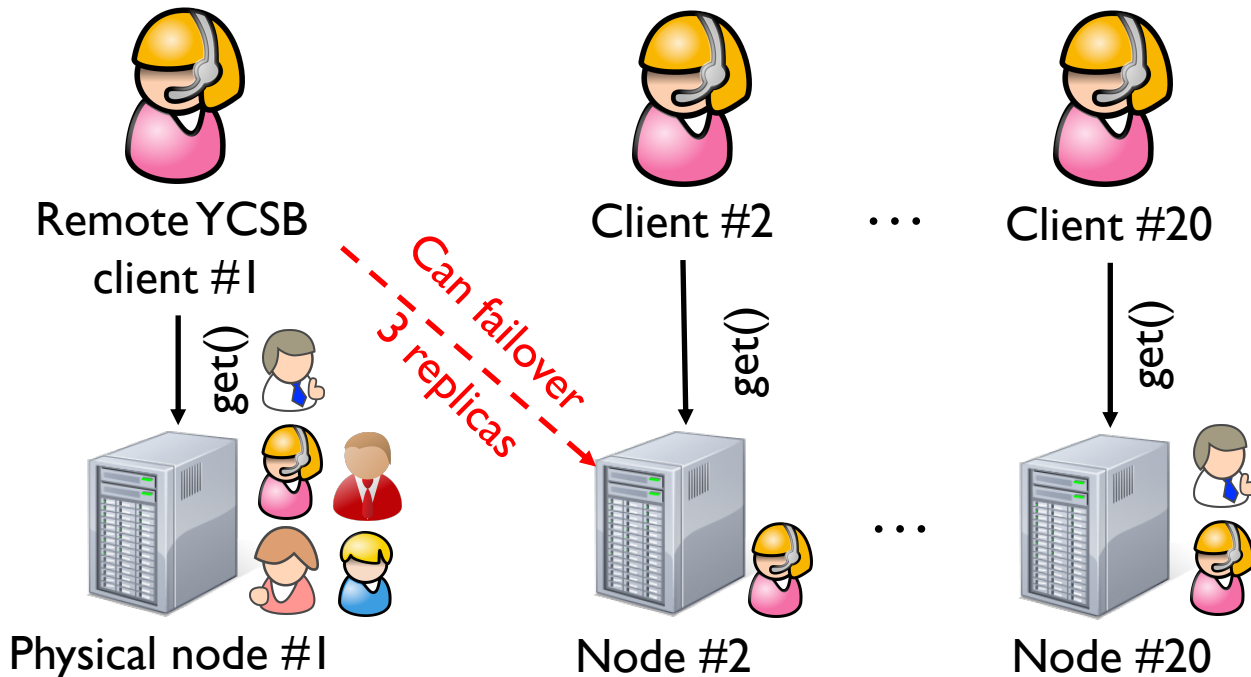
MittCFQ-powered MongoDB



MittCFQ-powered MongoDB



MittCFQ-powered MongoDB



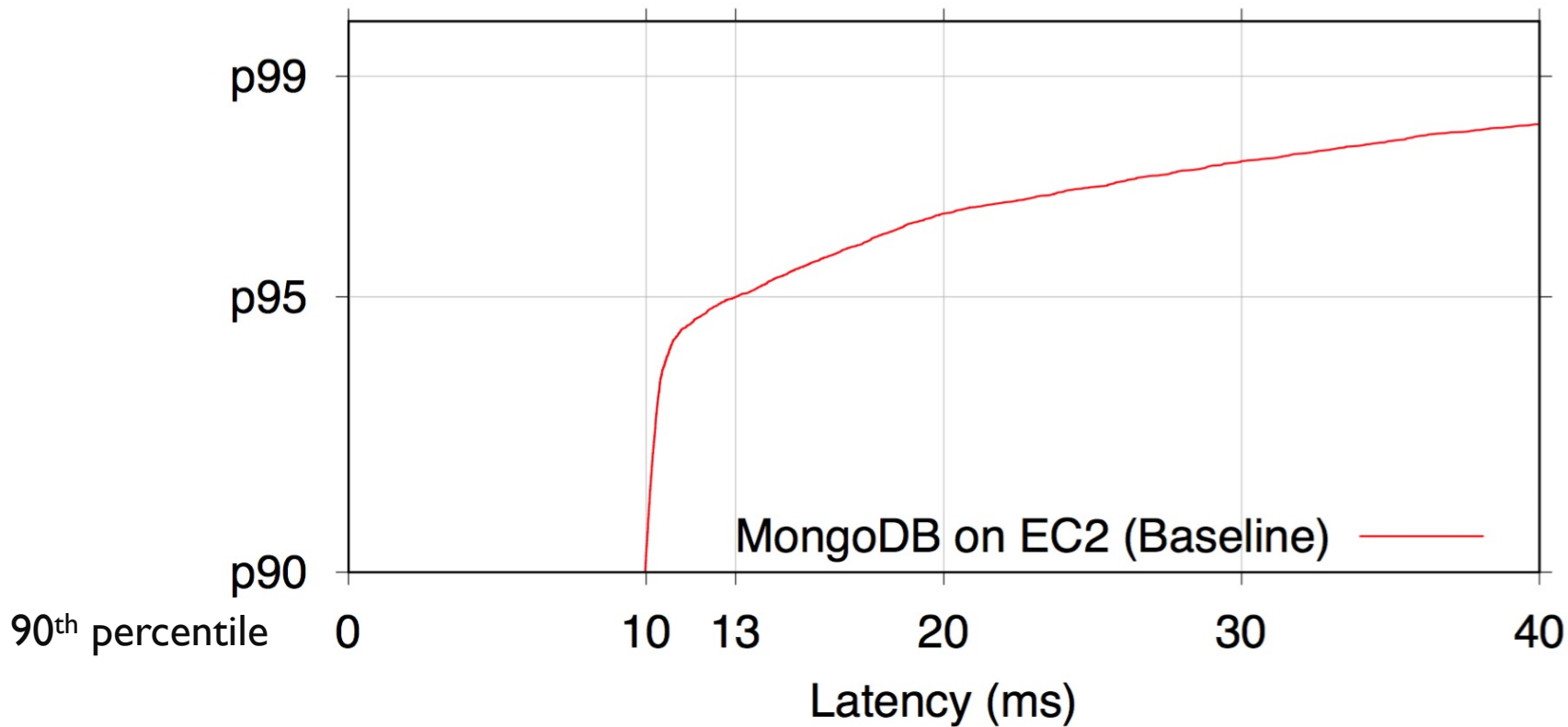
Metric:

CDF of all `get()` requests latencies (total 6 million data points)

Noisy neighbors based on EC2 data

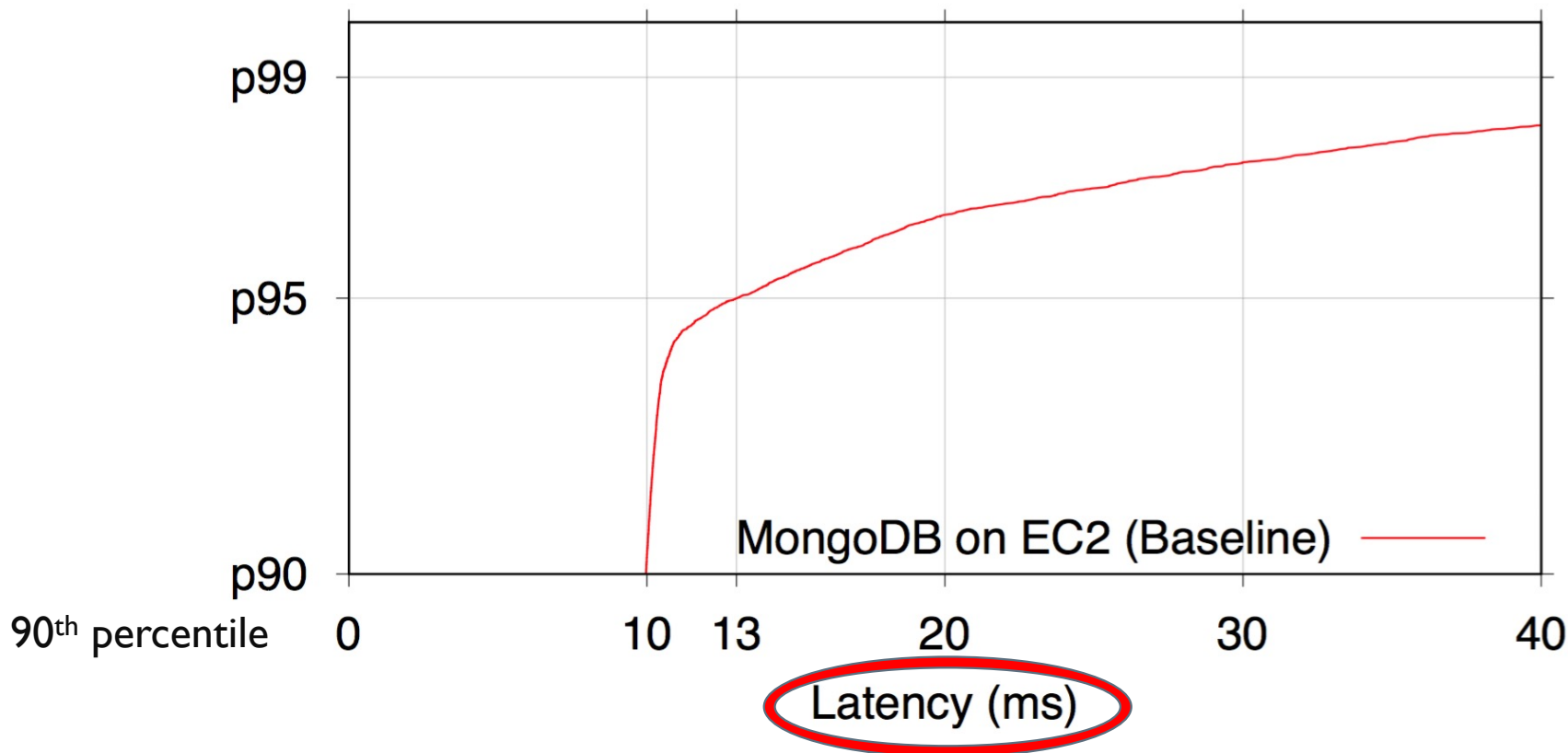
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



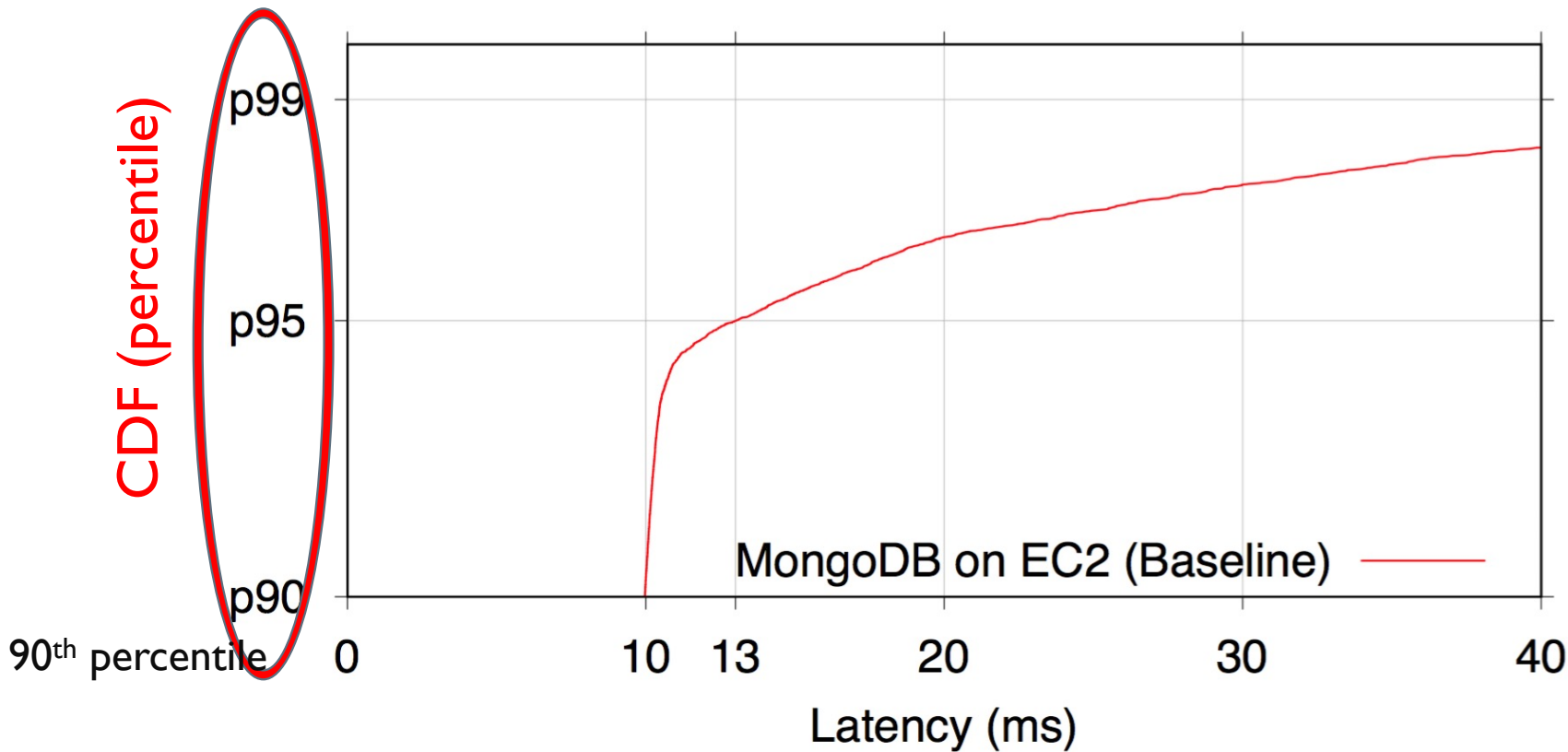
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



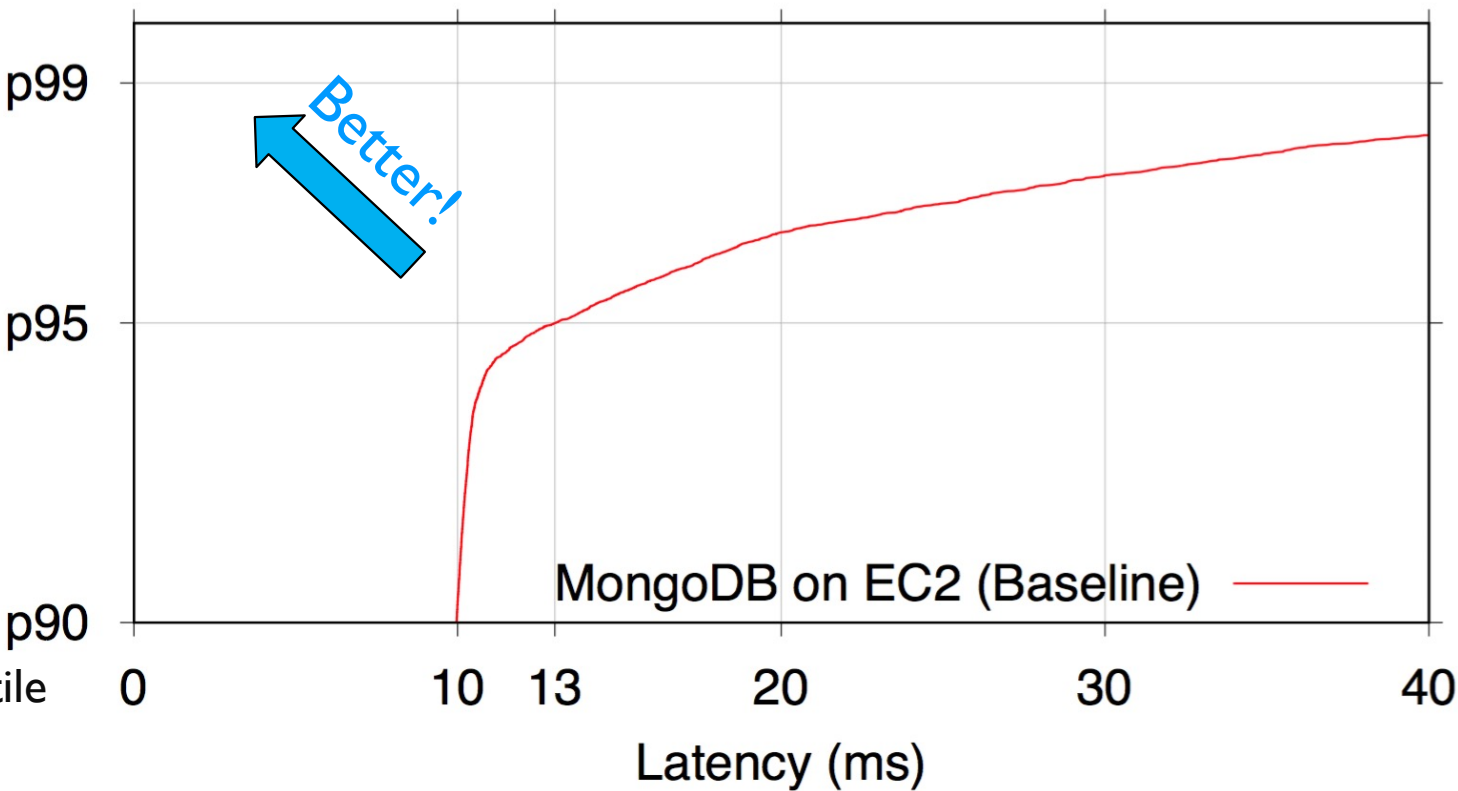
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



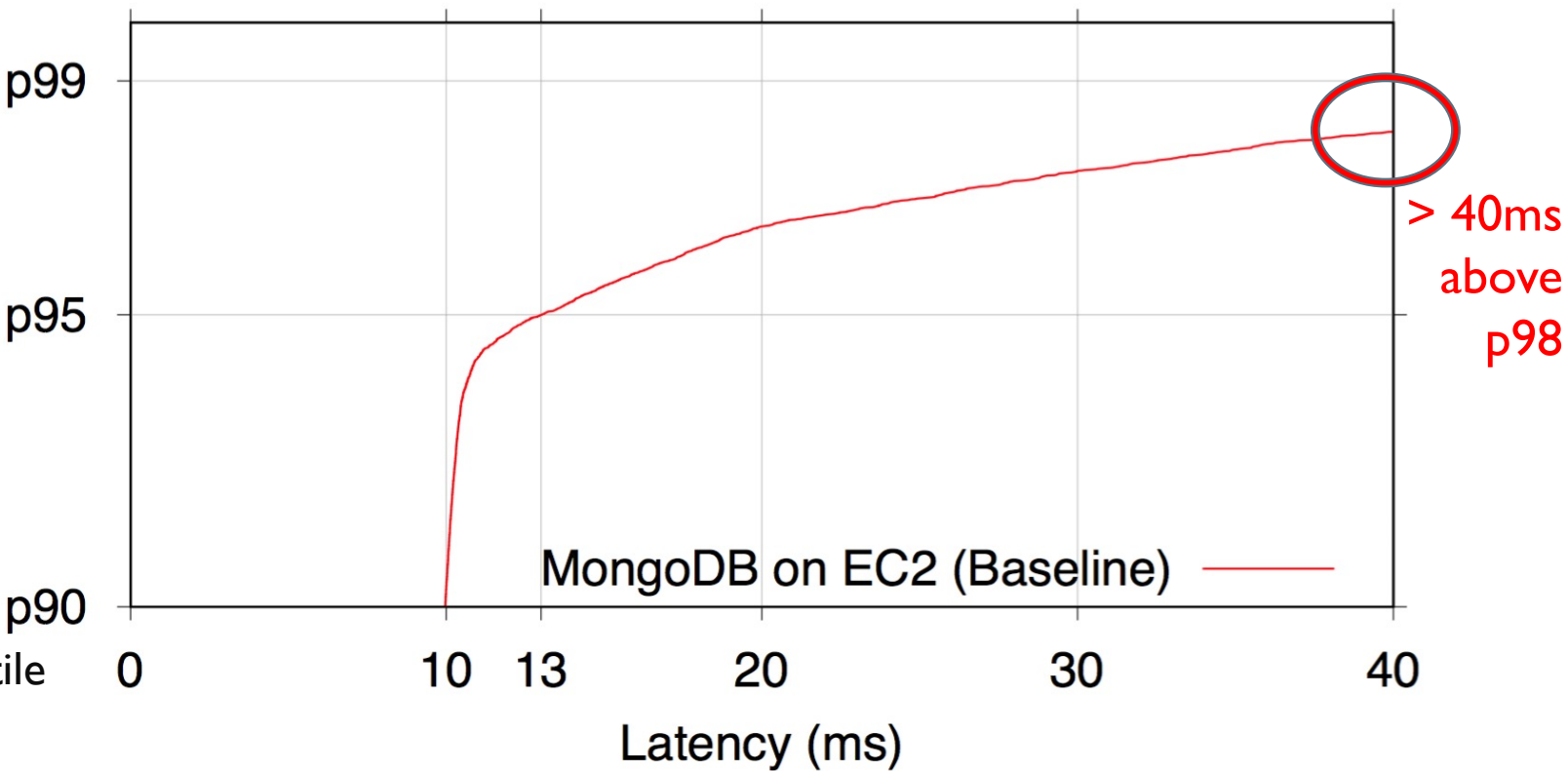
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



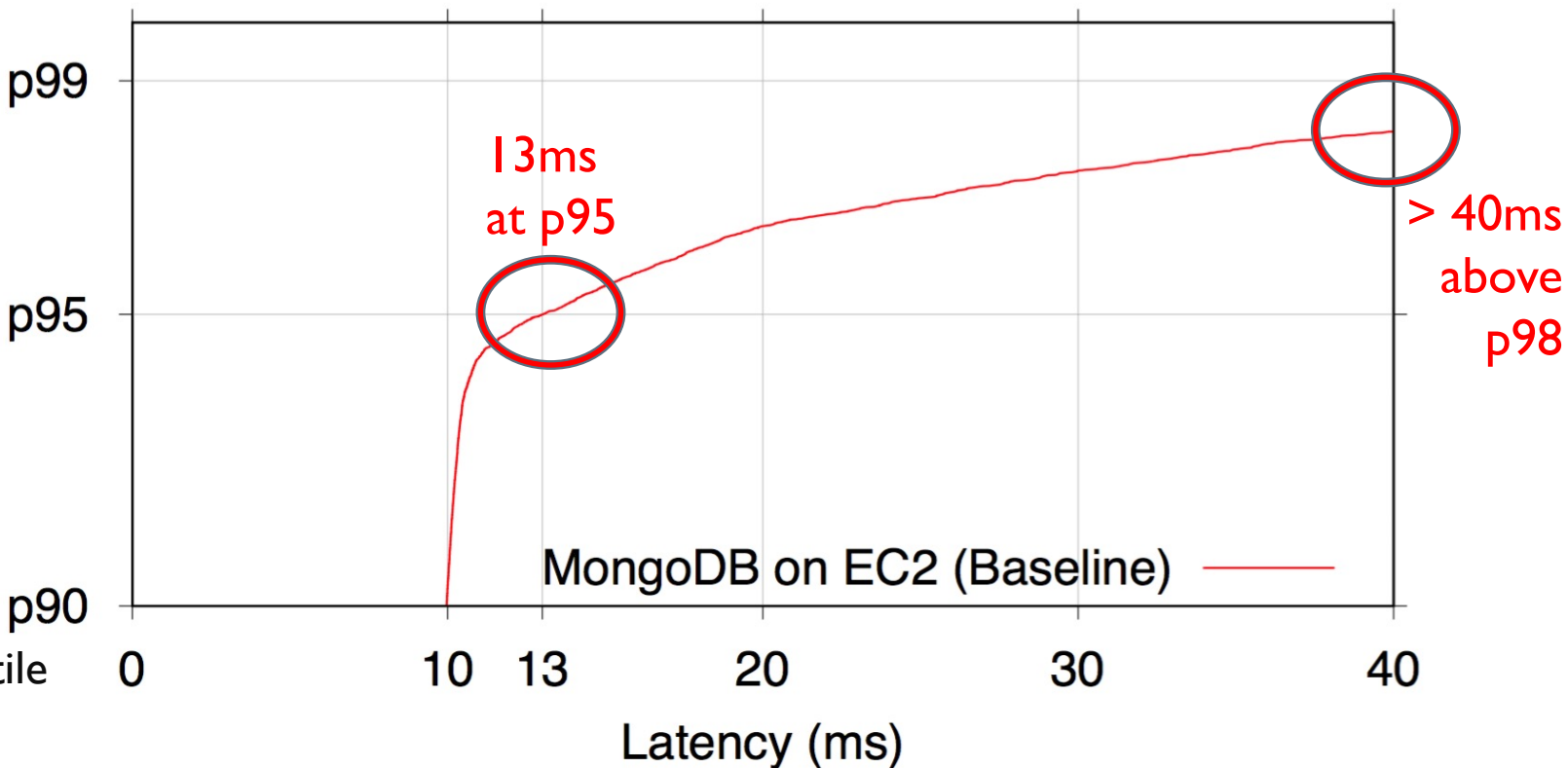
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



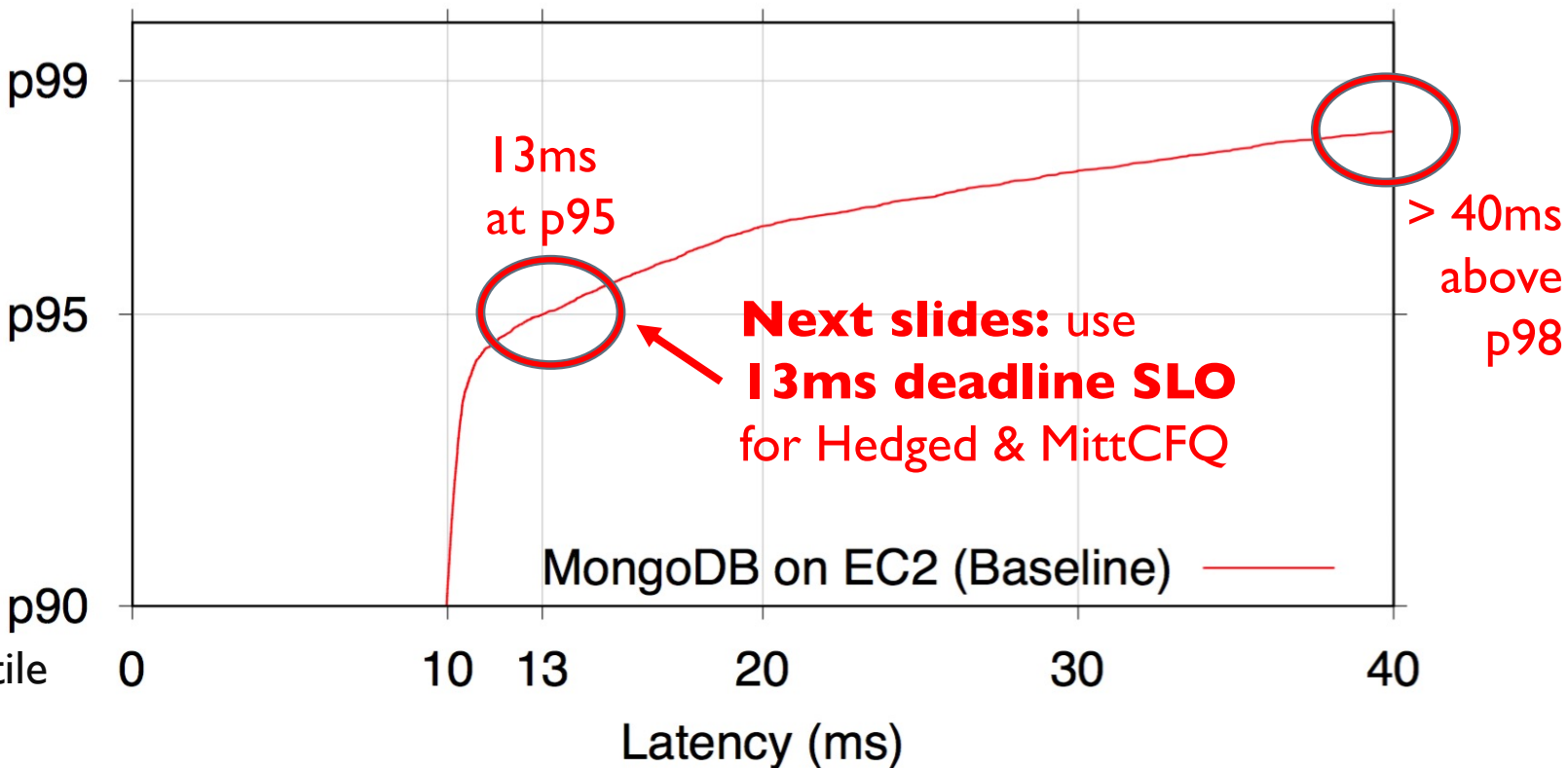
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



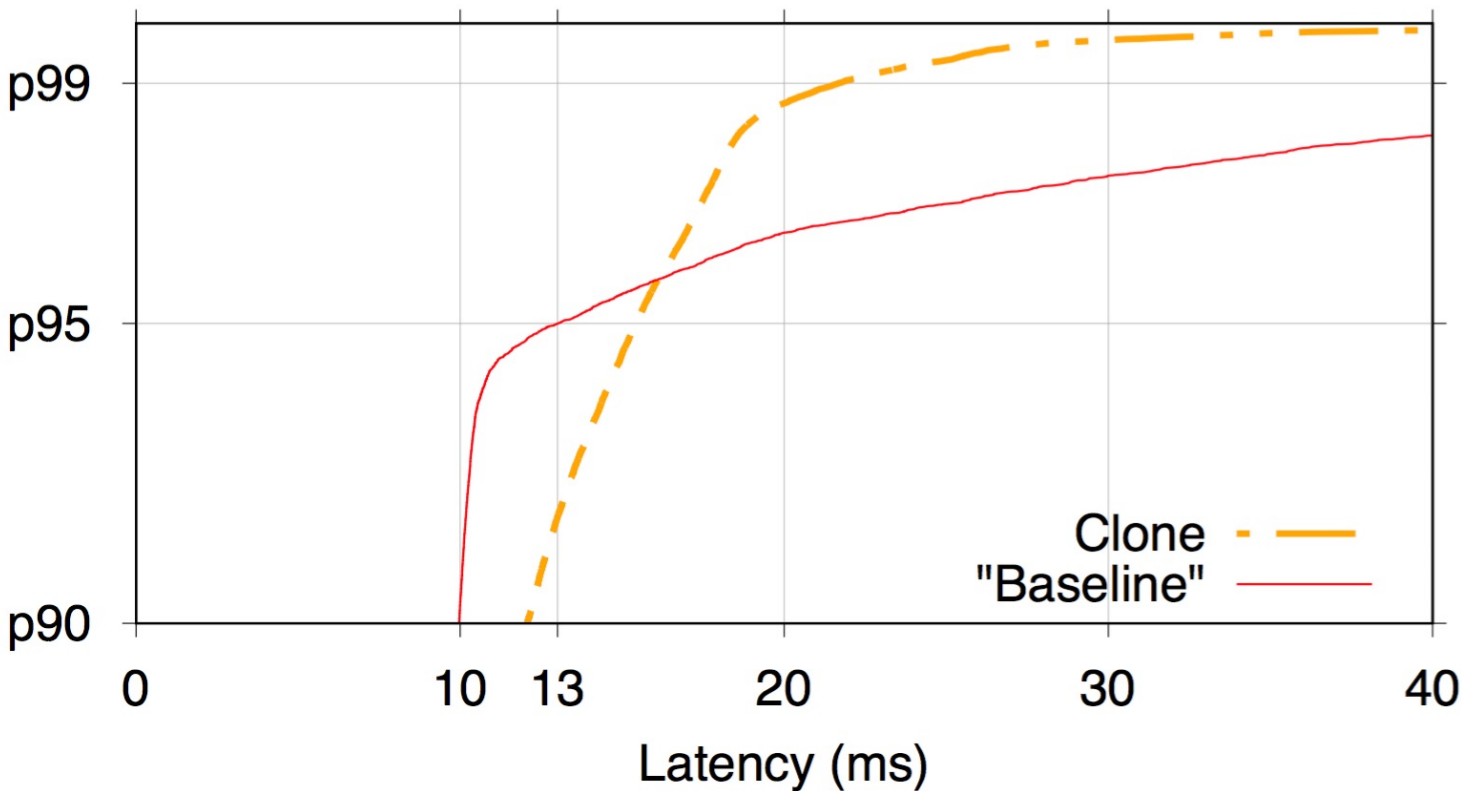
Baseline

CDF of YCSB get() Latencies on 20-node MongoDB



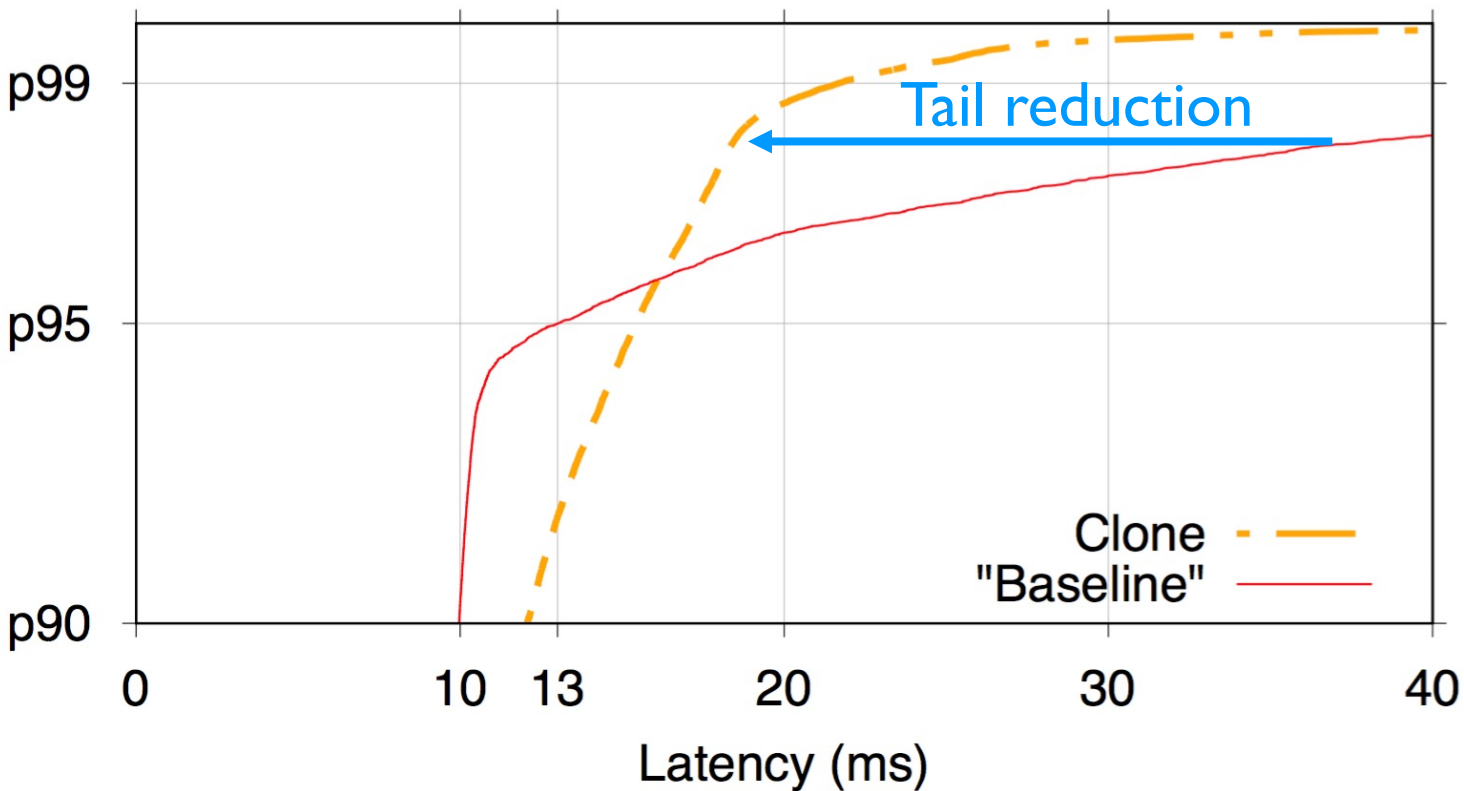
Clone

CDF of YCSB get() Latencies on 20-node MongoDB



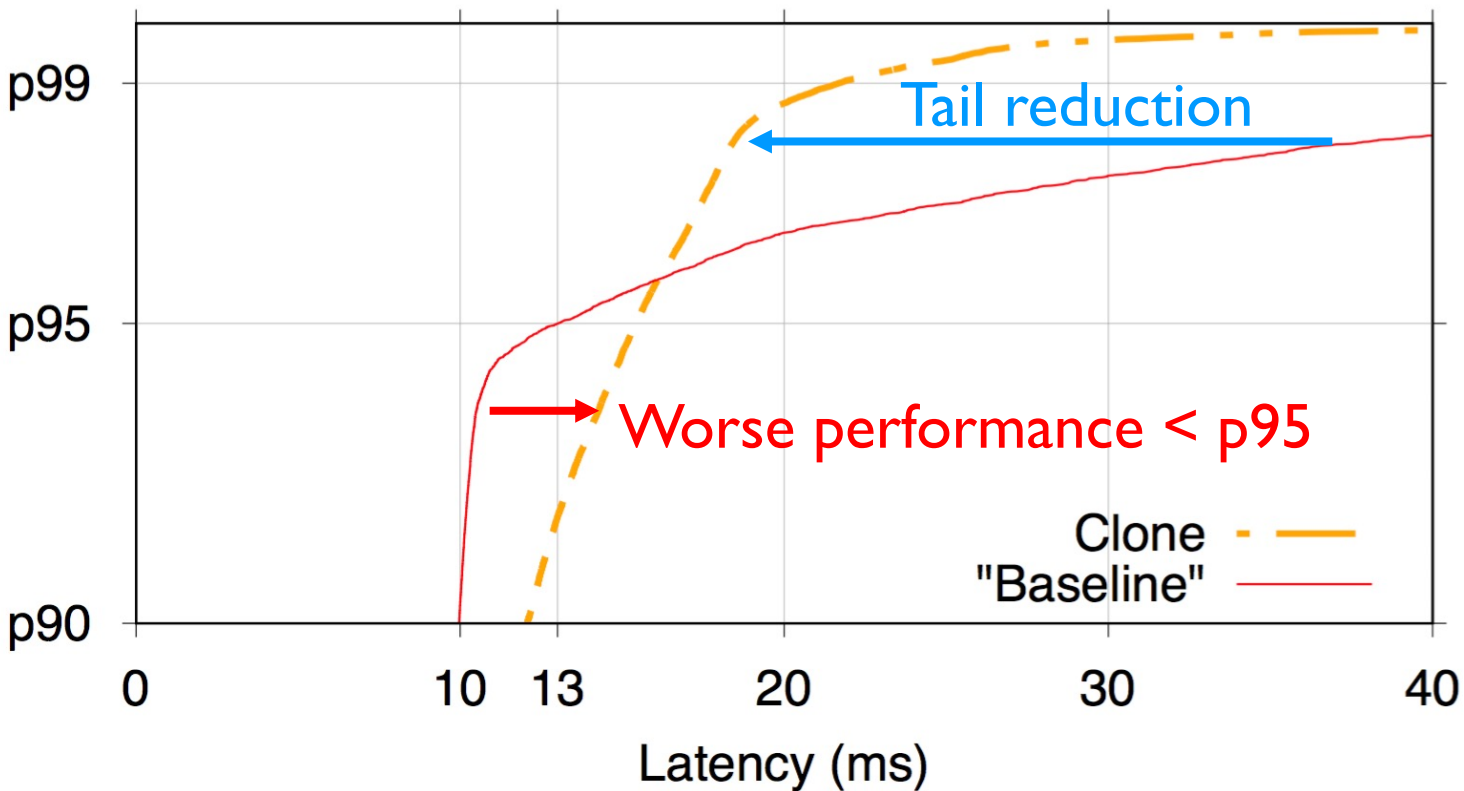
Clone

CDF of YCSB get() Latencies on 20-node MongoDB



Clone

CDF of YCSB get() Latencies on 20-node MongoDB





Hedged Requests

DOI:10.1145/2408776.2408794

Software techniques that tolerate latency variability are vital to building responsive large-scale Web services.

BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

The Tail at Scale

Communications of the ACM,
vol. 56 (2013), pp. 74-80



Hedged Requests

DOI:10.1145/2408776.2408794

Software techniques that tolerate latency variability are vital to building responsive large-scale Web services.

BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

The Tail at Scale

Communications of the ACM,
vol. 56 (2013), pp. 74-80



(a) Sends
first request





Hedged Requests

DOI:10.1145/2408776.2408794

Software techniques that tolerate latency variability are vital to building responsive large-scale Web services.

BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

The Tail at Scale

Communications of the ACM,
vol. 56 (2013), pp. 74-80

(b) Waits for
13ms **timeout**



(a) Sends
first request



Hedged Requests

DOI:10.1145/2408776.2408794

Software techniques that tolerate latency variability are vital to building responsive large-scale Web services.

BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

The Tail at Scale

Communications of the ACM,
vol. 56 (2013), pp. 74-80

(b) Waits for
13ms **timeout**



(a) Sends
first request



(c) Sends
secondary
request



Hedged Requests

DOI:10.1145/2408776.2408794

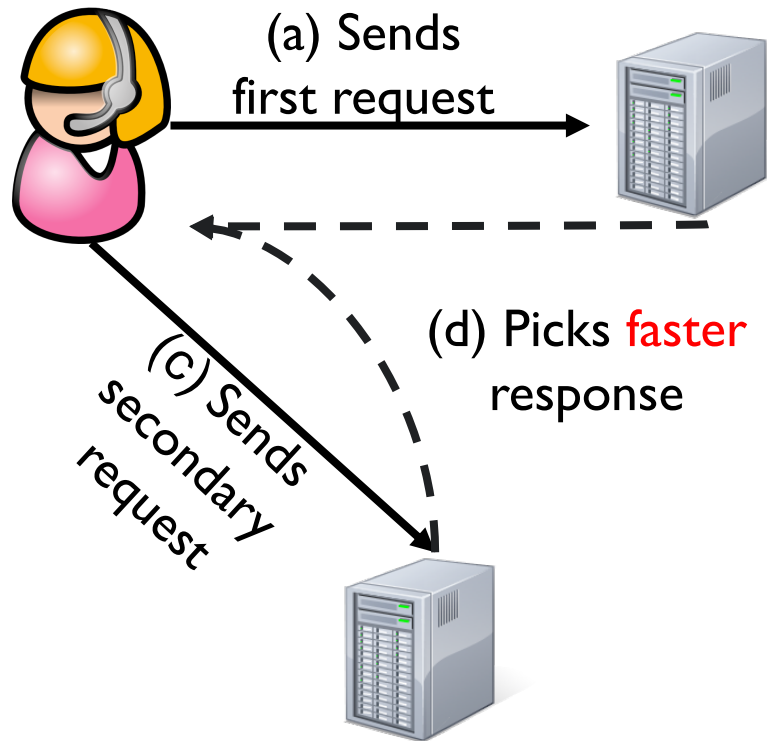
Software techniques that tolerate latency variability are vital to building responsive large-scale Web services.

BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

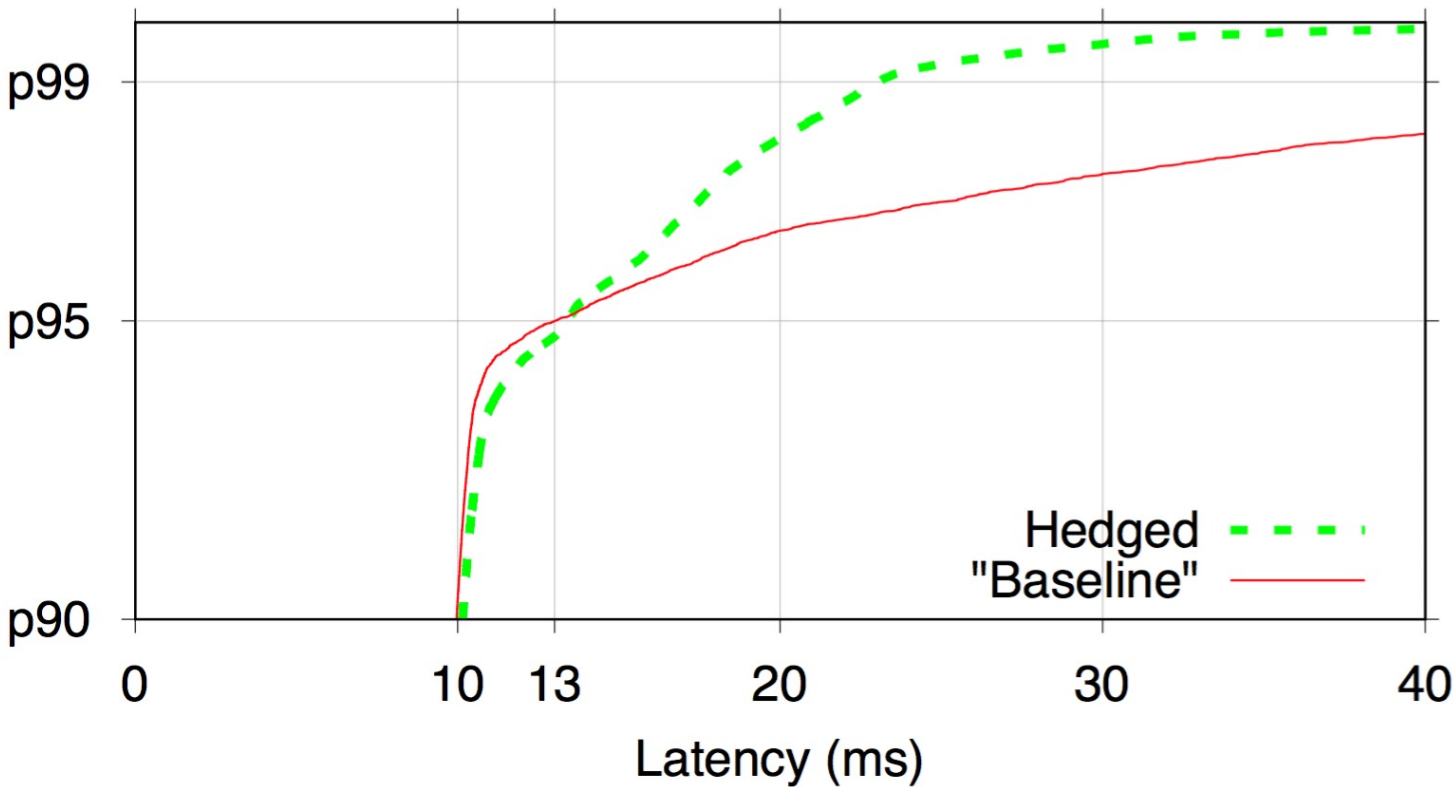
The Tail at Scale

Communications of the ACM,
vol. 56 (2013), pp. 74-80

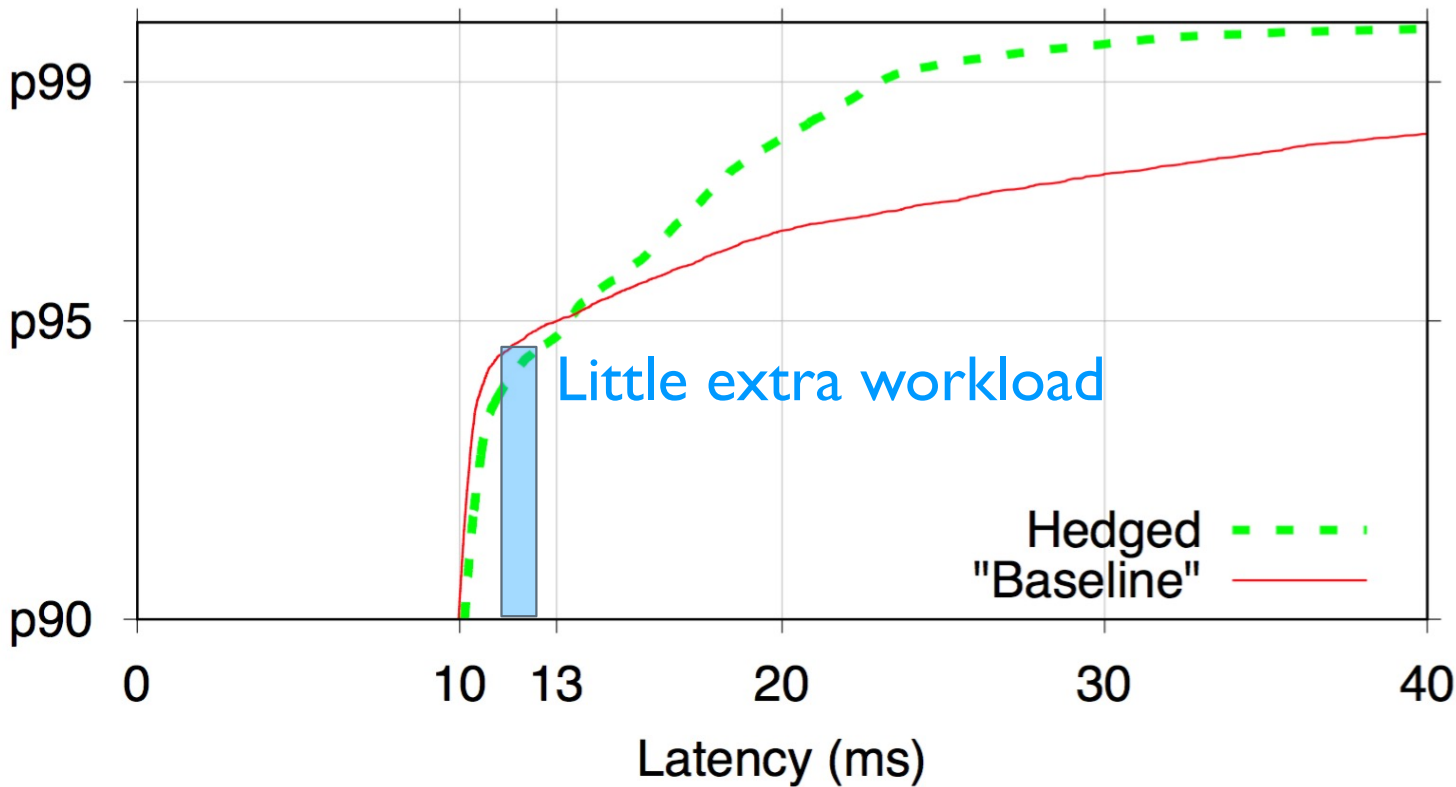
(b) Waits for
13ms **timeout** 



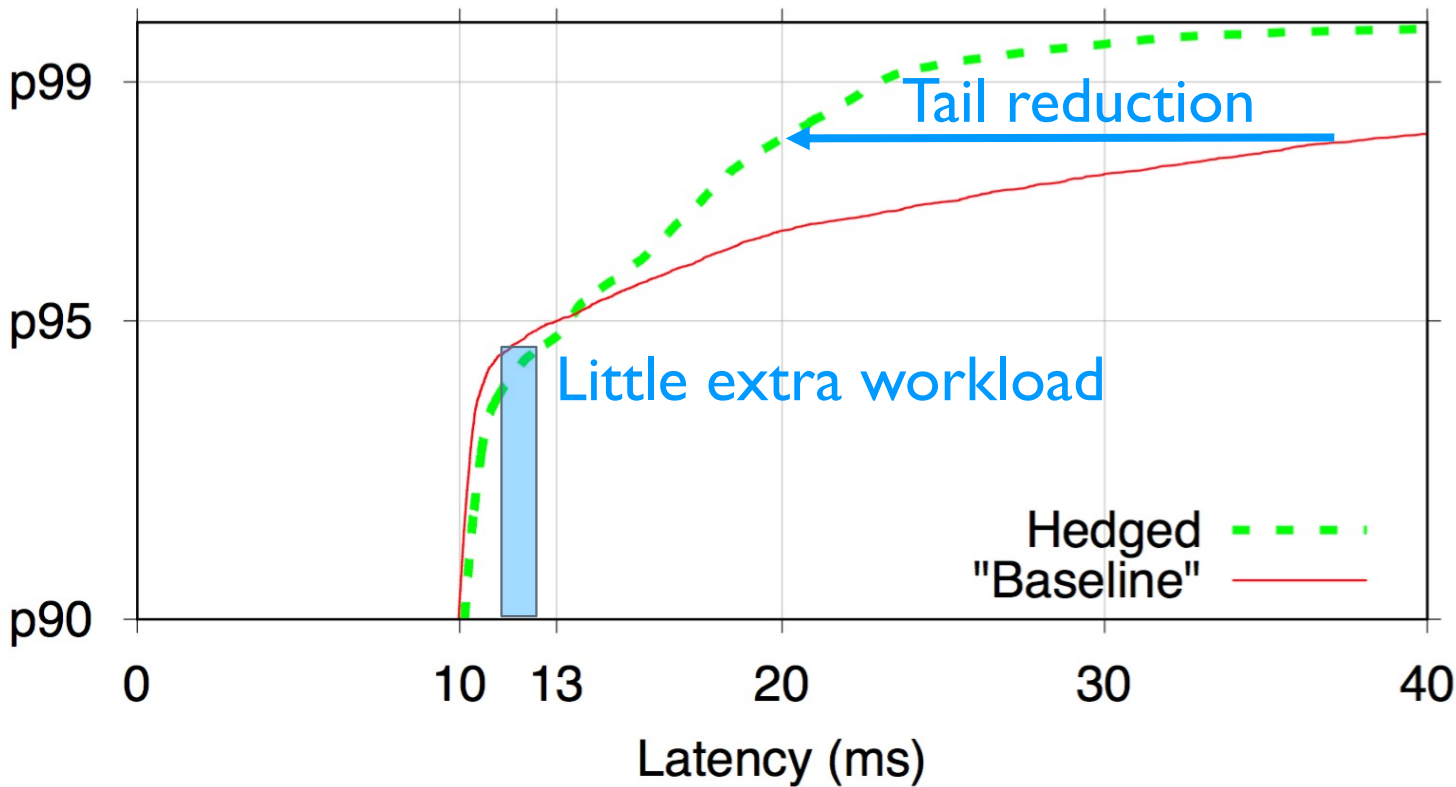
CDF of YCSB get() Latencies on 20-node MongoDB



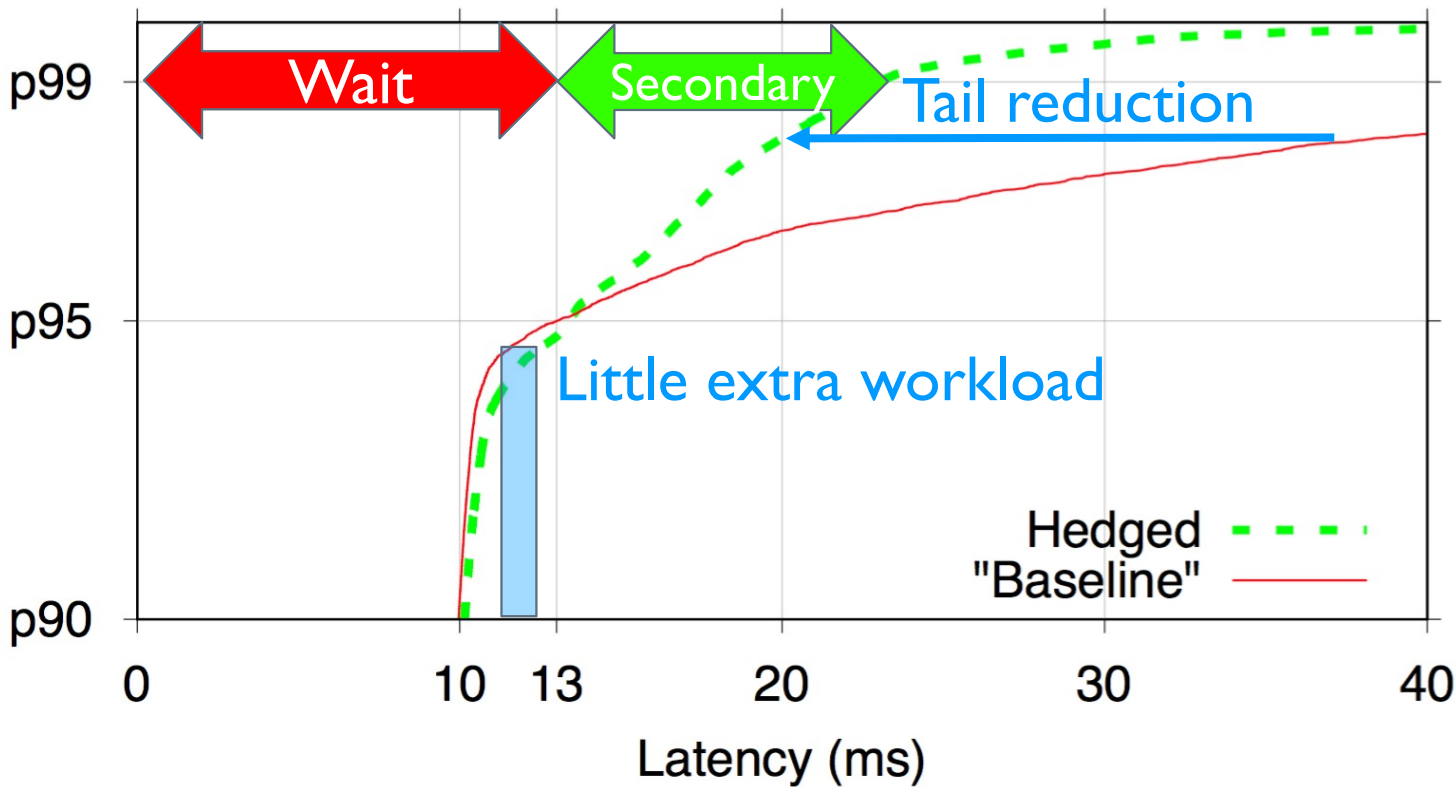
CDF of YCSB get() Latencies on 20-node MongoDB



CDF of YCSB get() Latencies on 20-node MongoDB

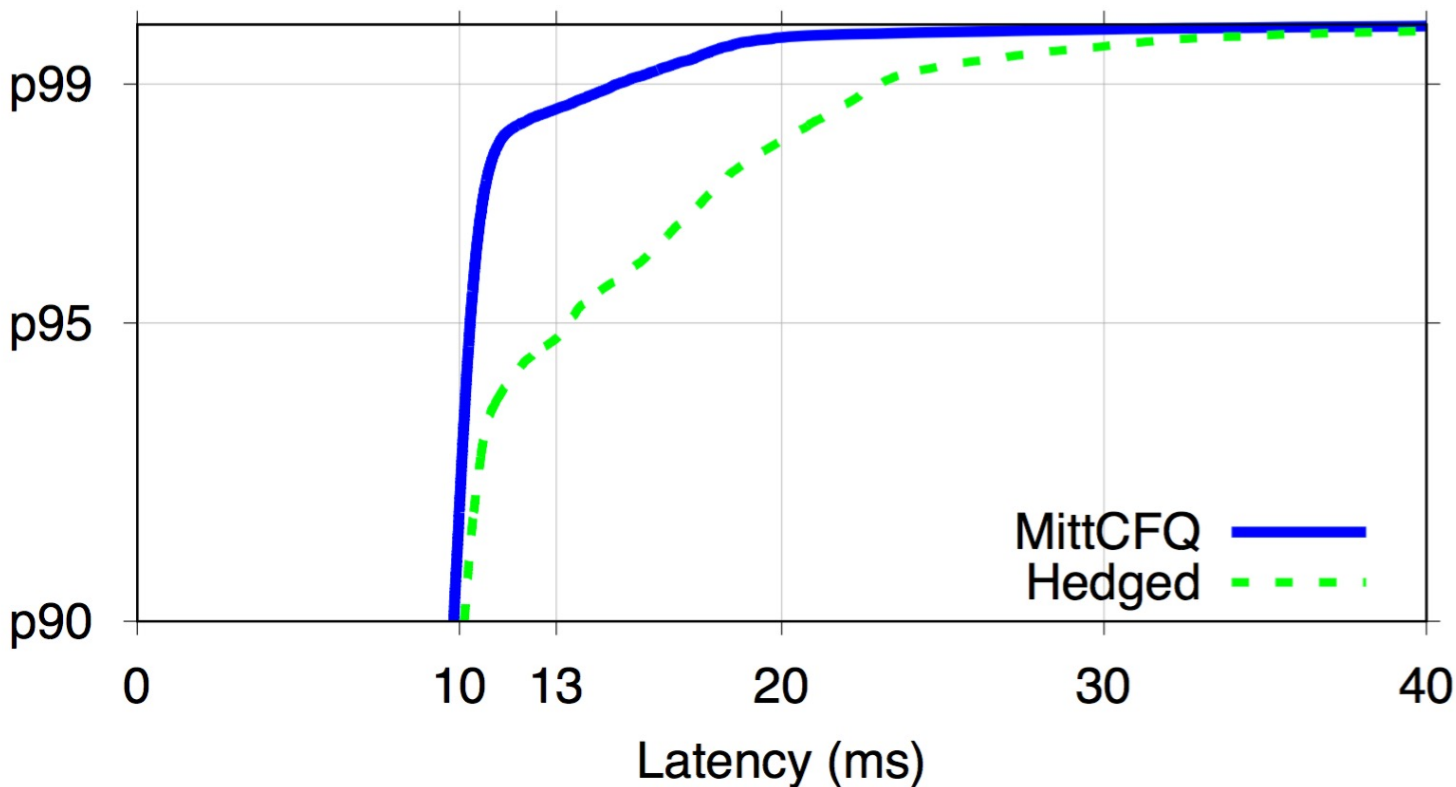


CDF of YCSB get() Latencies on 20-node MongoDB



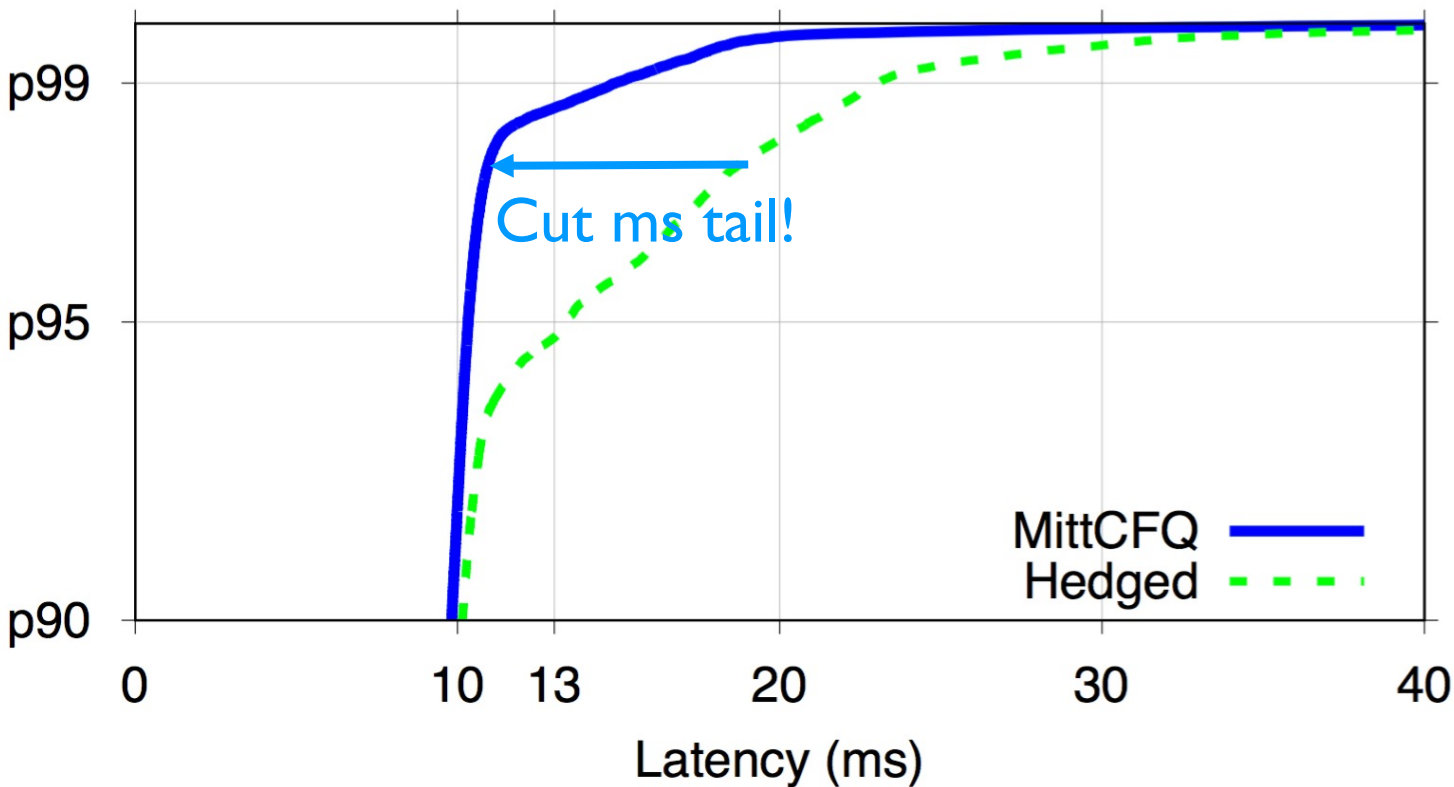
MittCFQ

CDF of YCSB get() Latencies on 20-node MongoDB



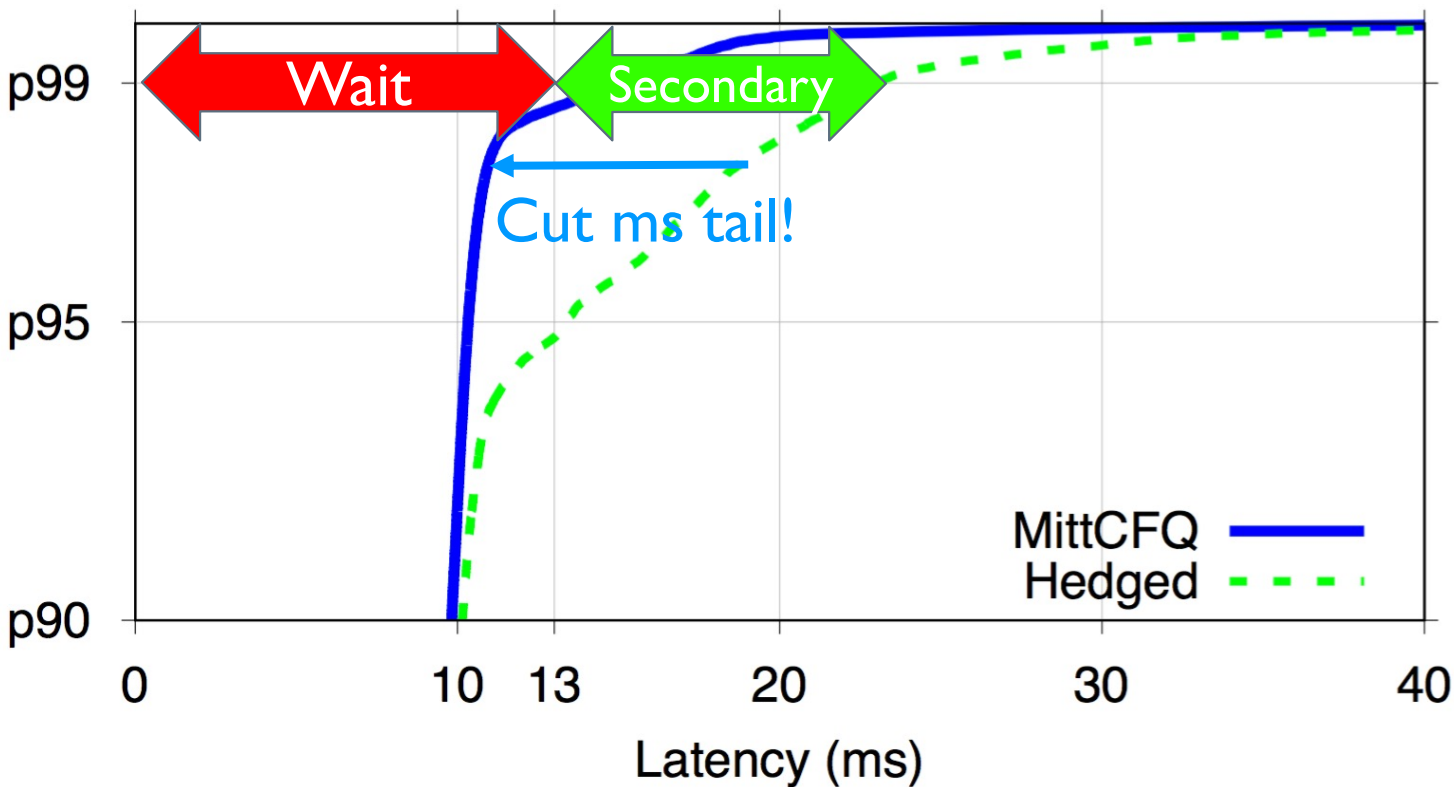
MittCFQ

CDF of YCSB get() Latencies on 20-node MongoDB



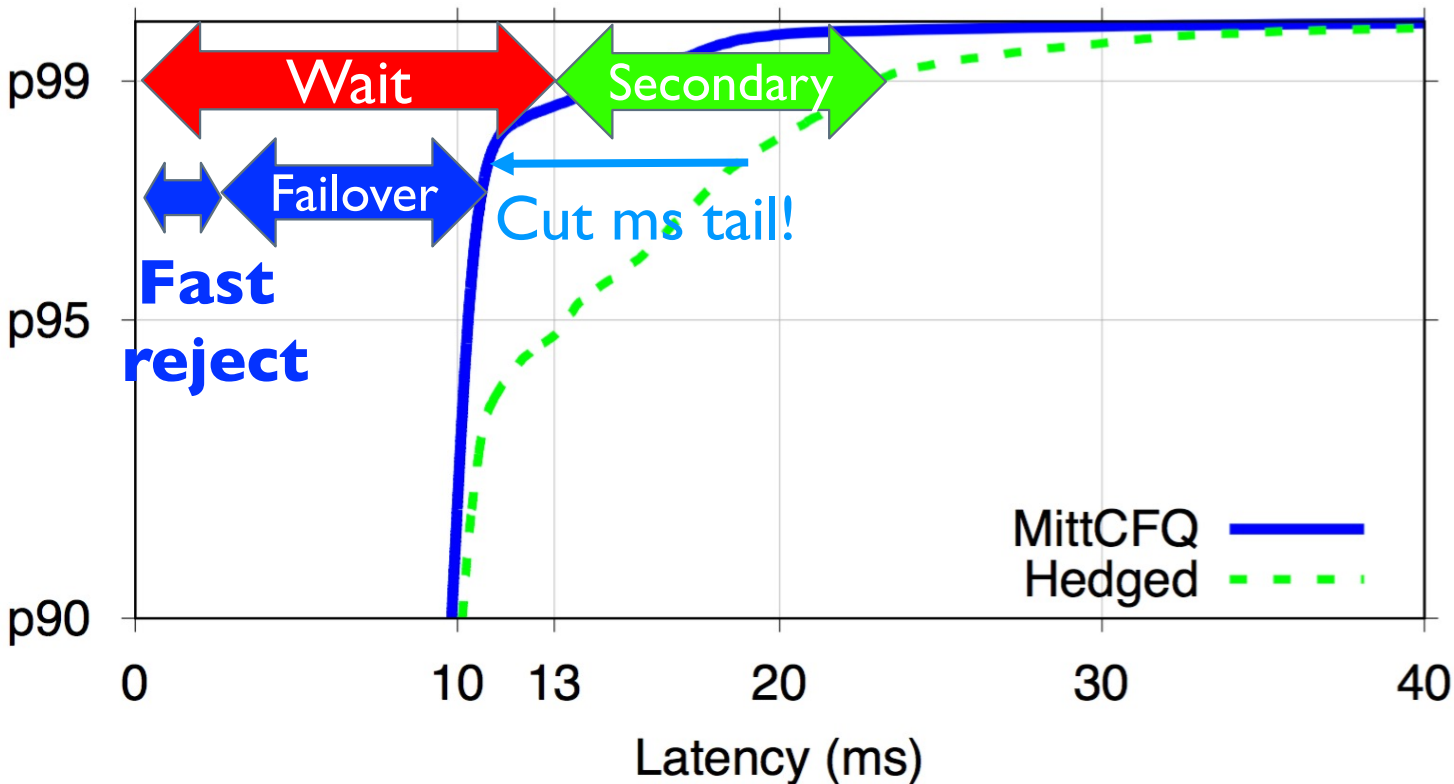
MittCFQ

CDF of YCSB get() Latencies on 20-node MongoDB



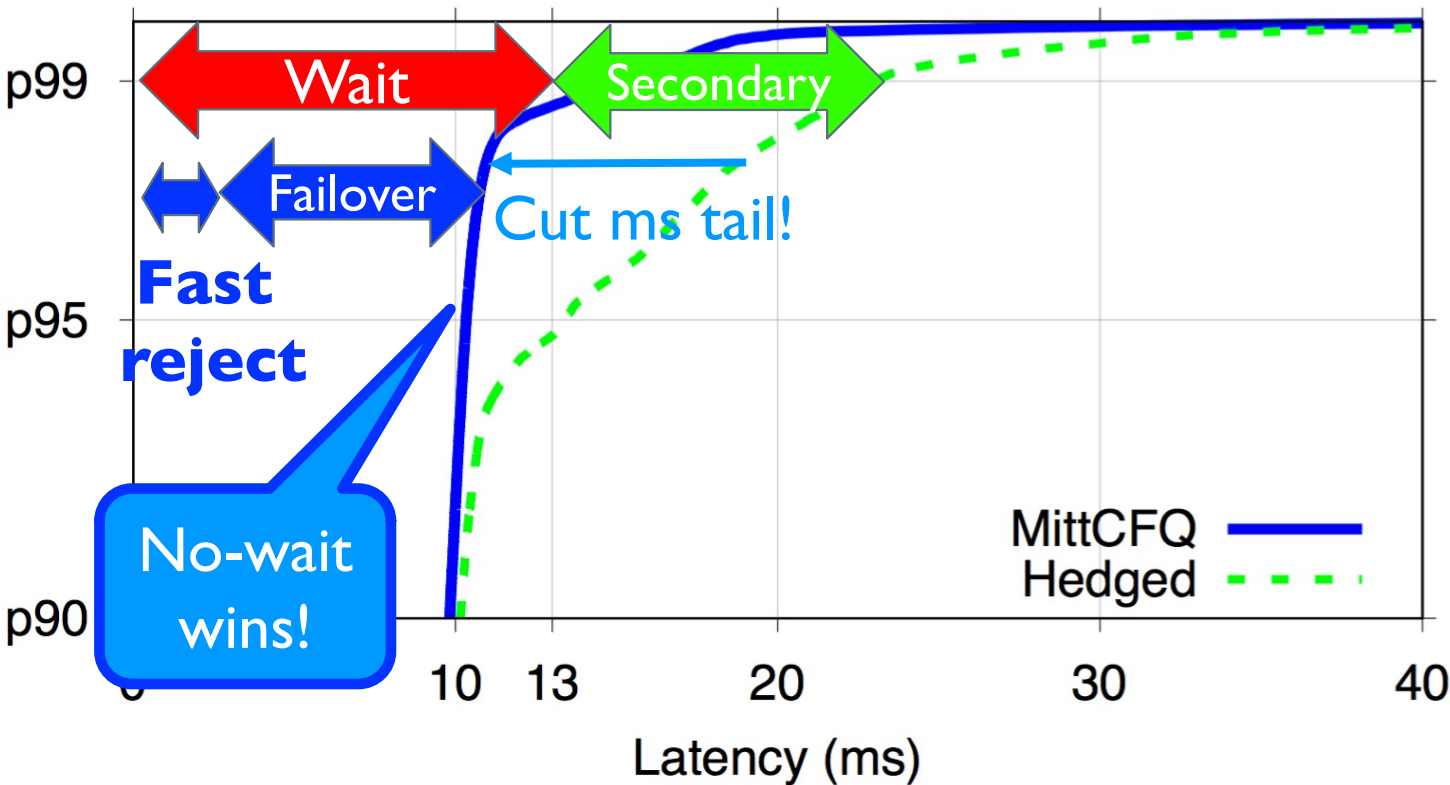
MittCFQ

CDF of YCSB get() Latencies on 20-node MongoDB



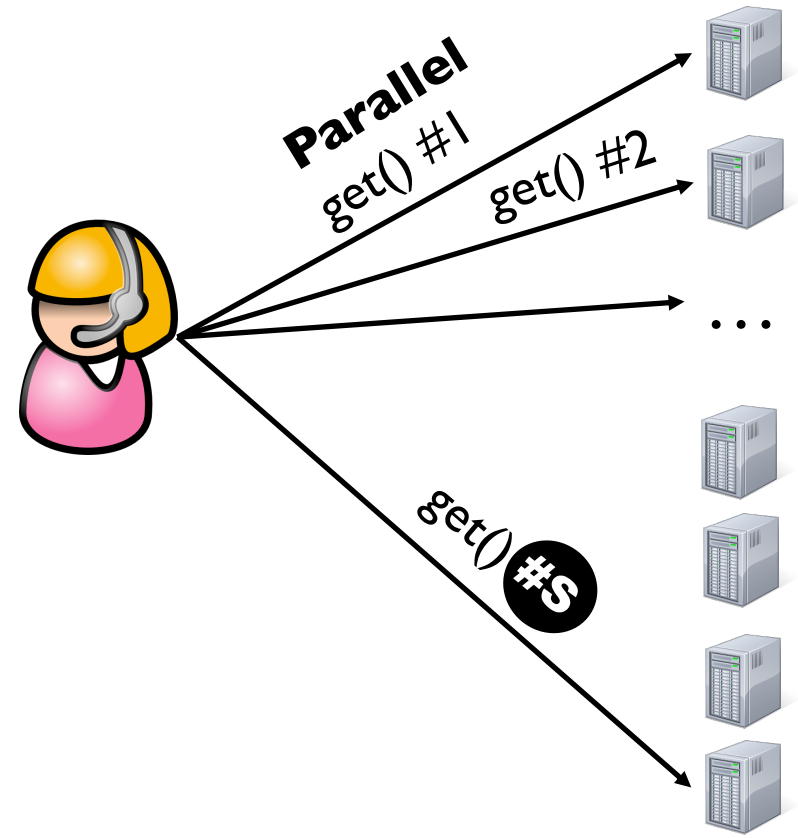
MittCFQ

CDF of YCSB get() Latencies on 20-node MongoDB

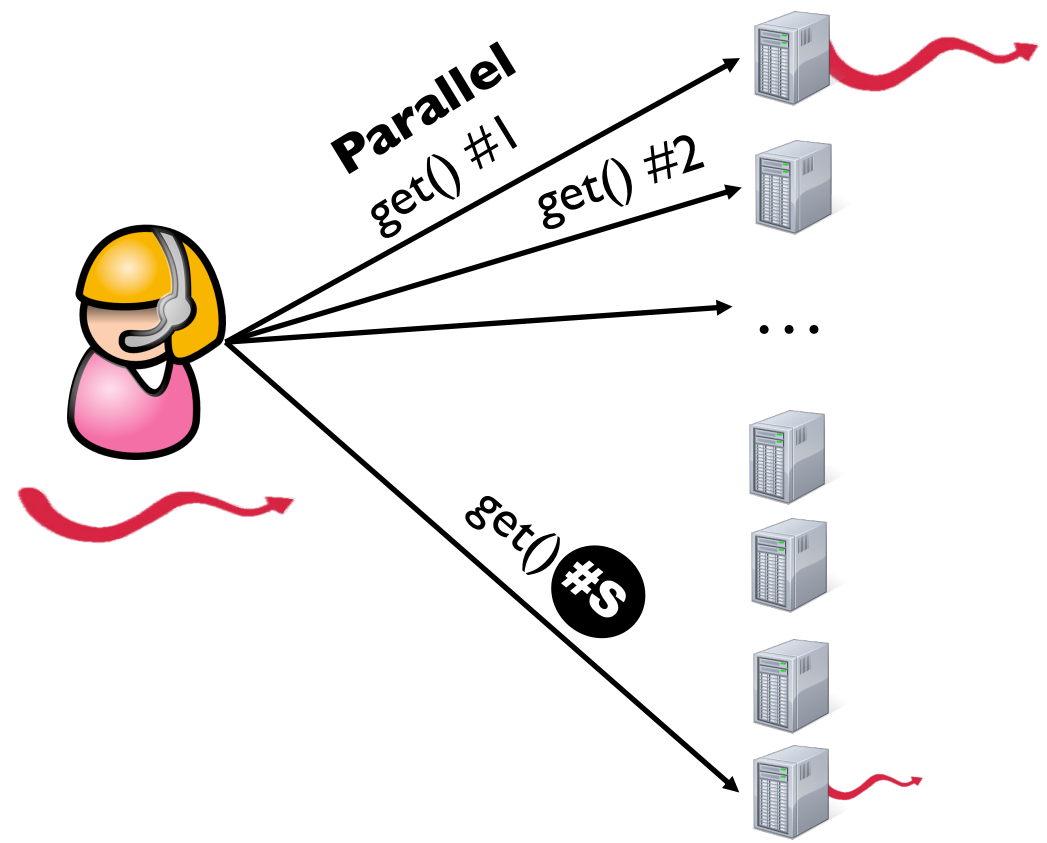


Tail amplified at Scale

Tail amplified at Scale

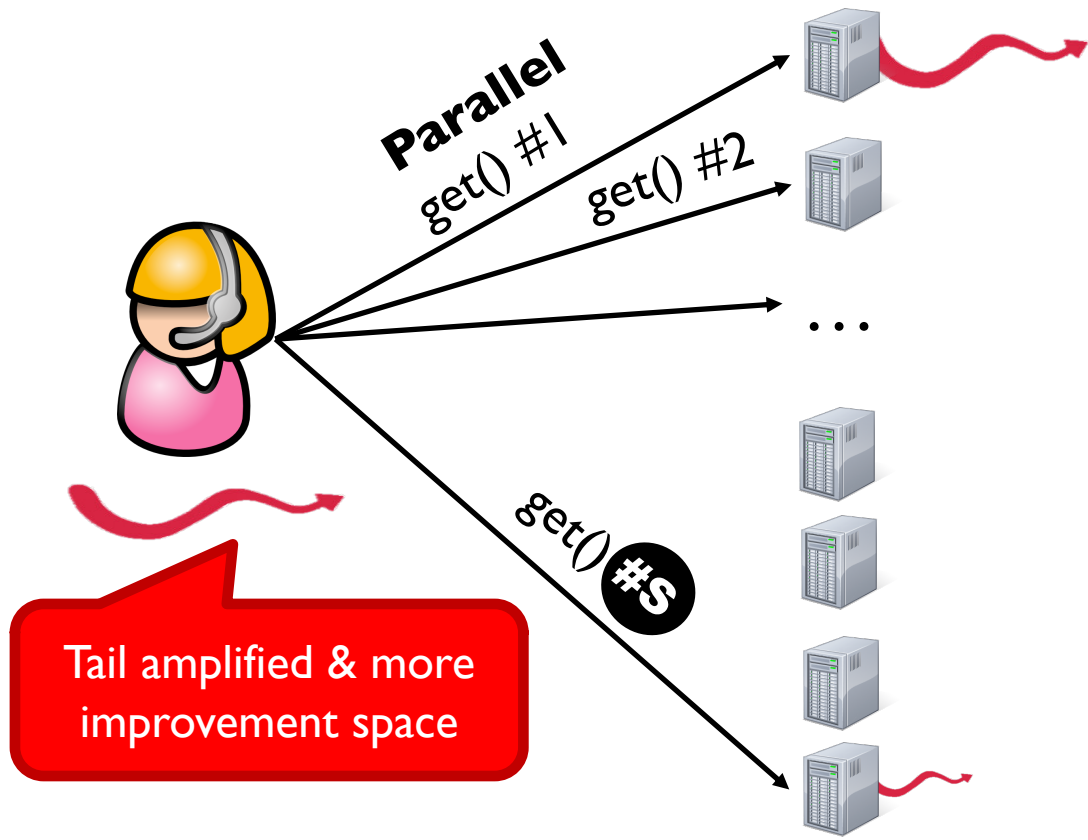


Tail amplified at Scale

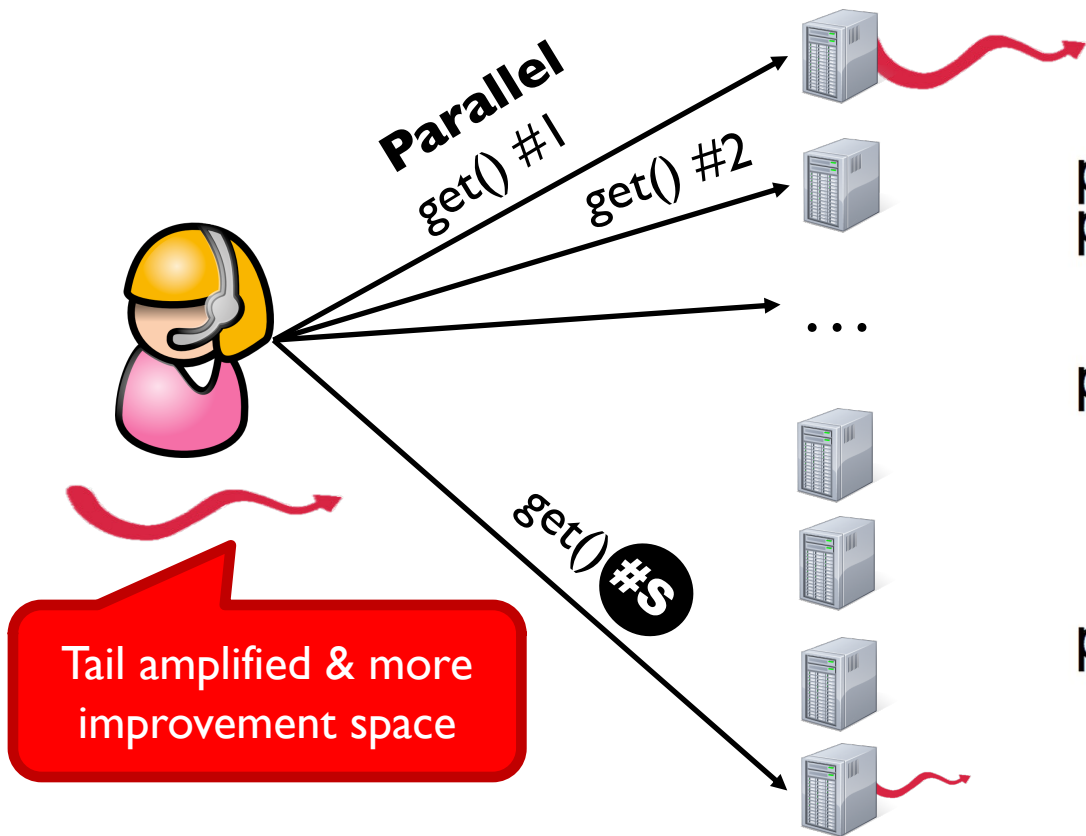




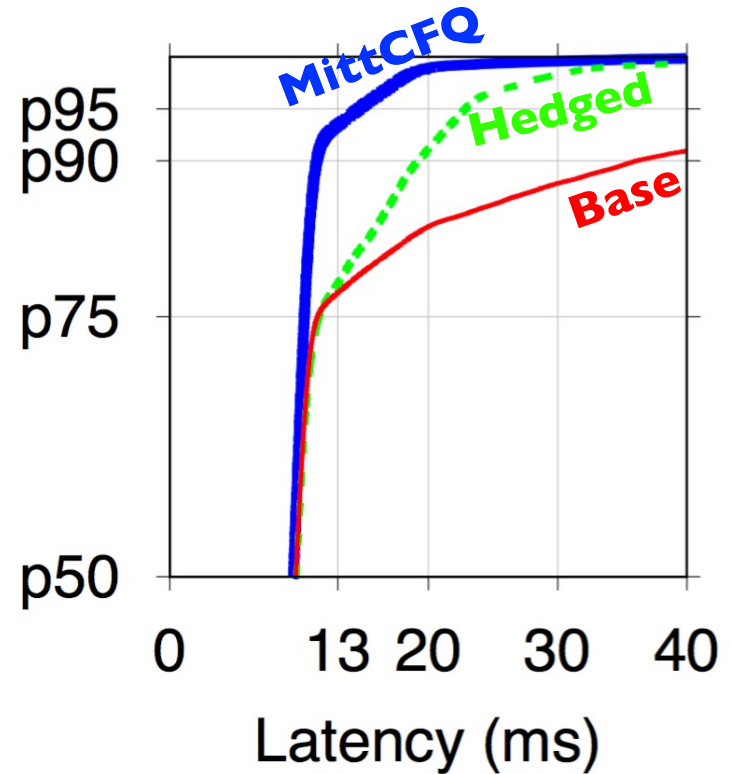
Tail amplified at Scale



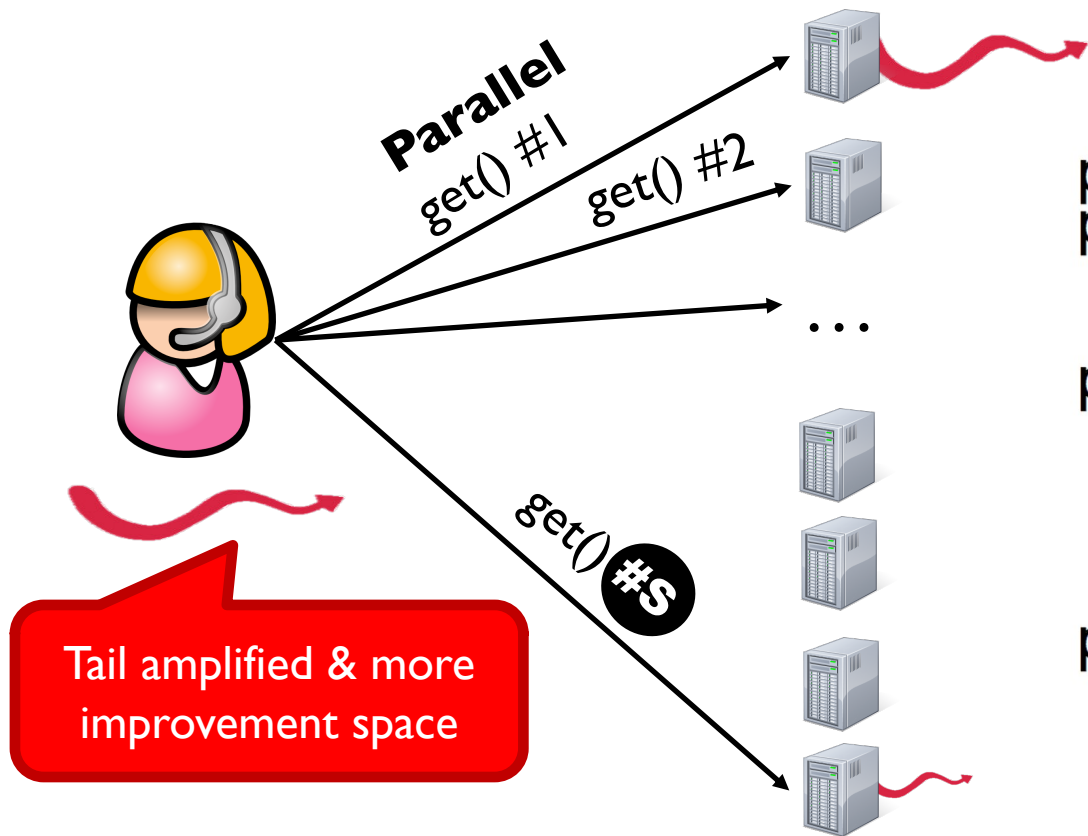
Tail amplified at Scale



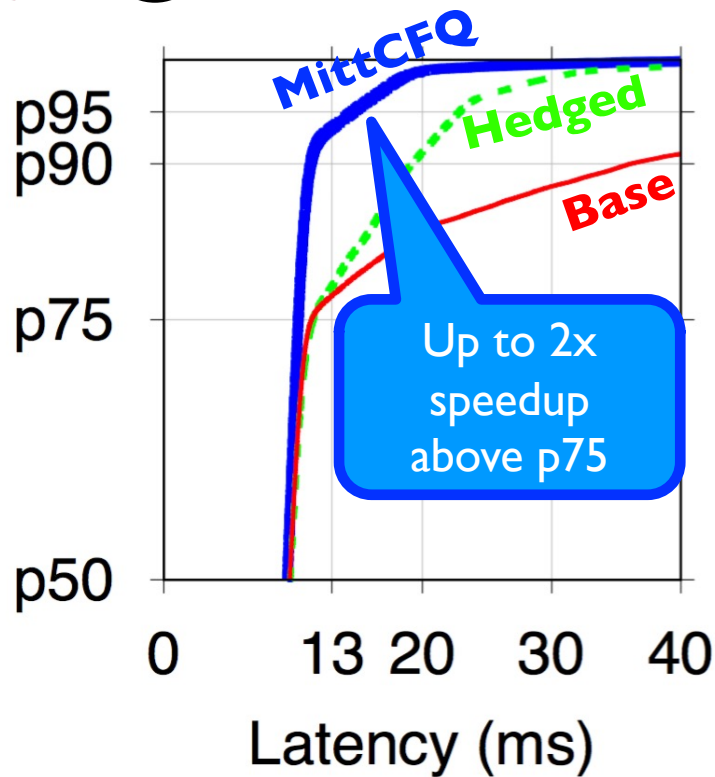
S Scale Factor: 5



Tail amplified at Scale



S Scale Factor: 5



Accuracy Evaluation

Accuracy Evaluation

MittCFQ



Disk

MittSSD



Open-Channel
SSD



Accuracy Evaluation

MittCFQ



Disk

MittSSD



Open-Channel
SSD

5 real-world block-level traces

DAPPS
DTRS TPCC
EXCH LMBE



Accuracy Evaluation

MittCFQ



Disk

MittSSD



Open-Channel
SSD

Metrics:

- **False positive:** IO rejected, but deadline is met

5 real-world block-level traces

DAPPS
DTRS TPCC
EXCH LMBE



Accuracy Evaluation

MittCFQ



Disk

MittSSD



Open-Channel
SSD

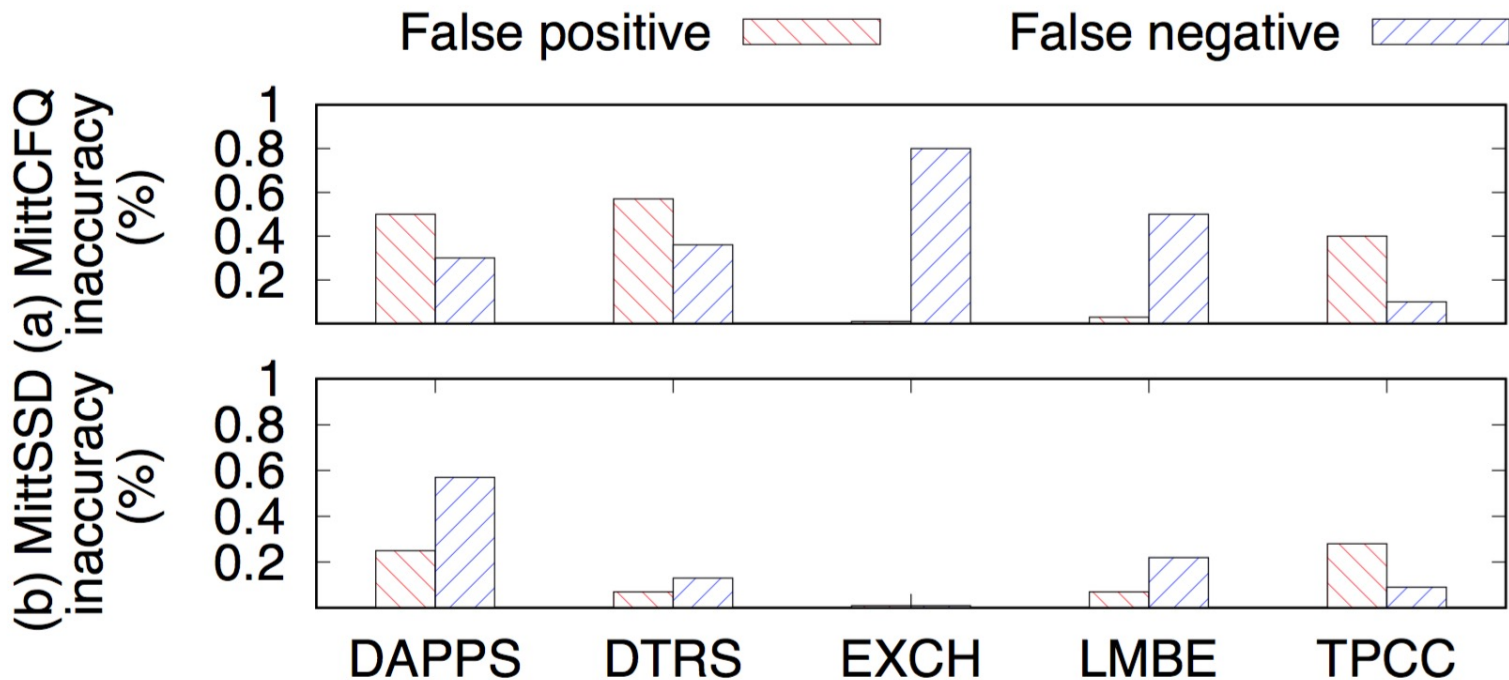
5 real-world block-level traces

DTRS DAPPS
 TPCC
EXCH LMBE

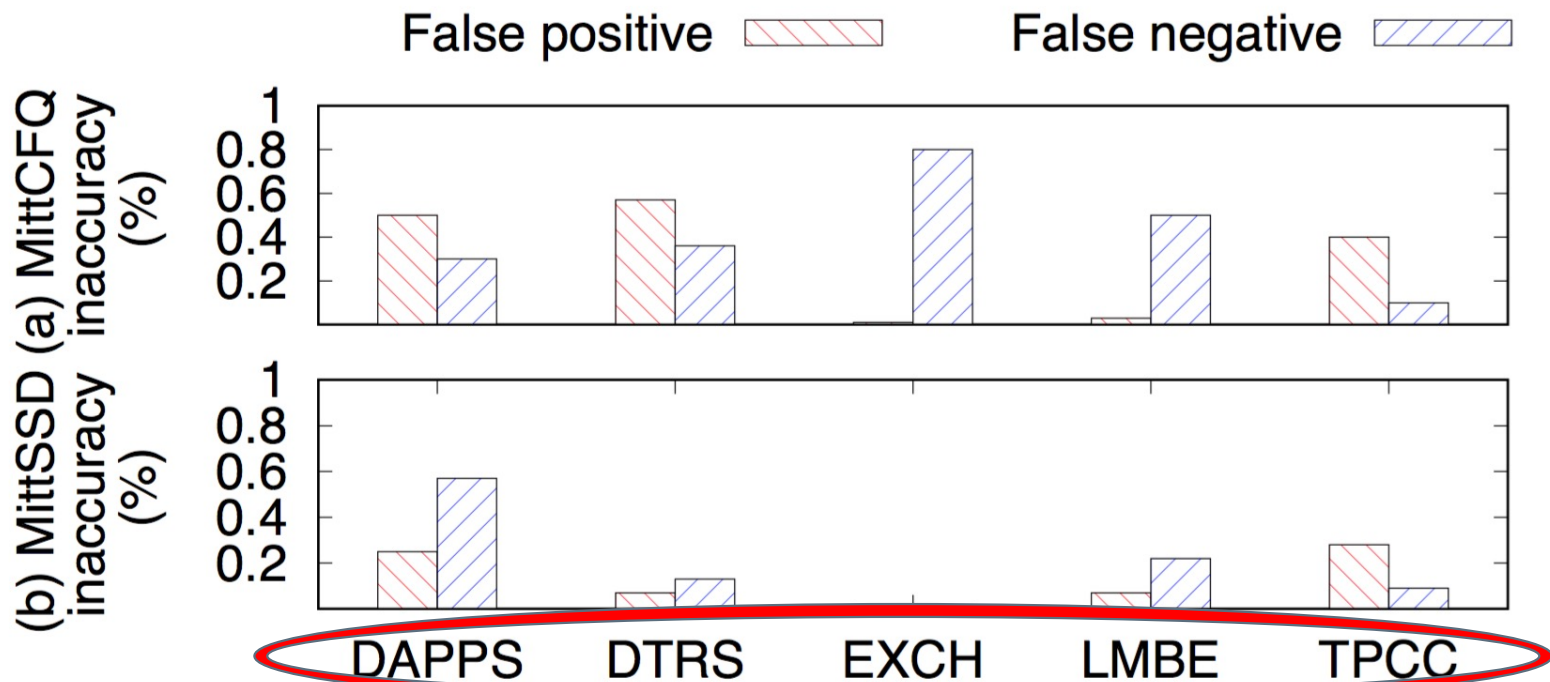
Metrics:

- **False positive:** IO rejected, but deadline is met
- **False negative:** Deadline violated, but IO is not rejected

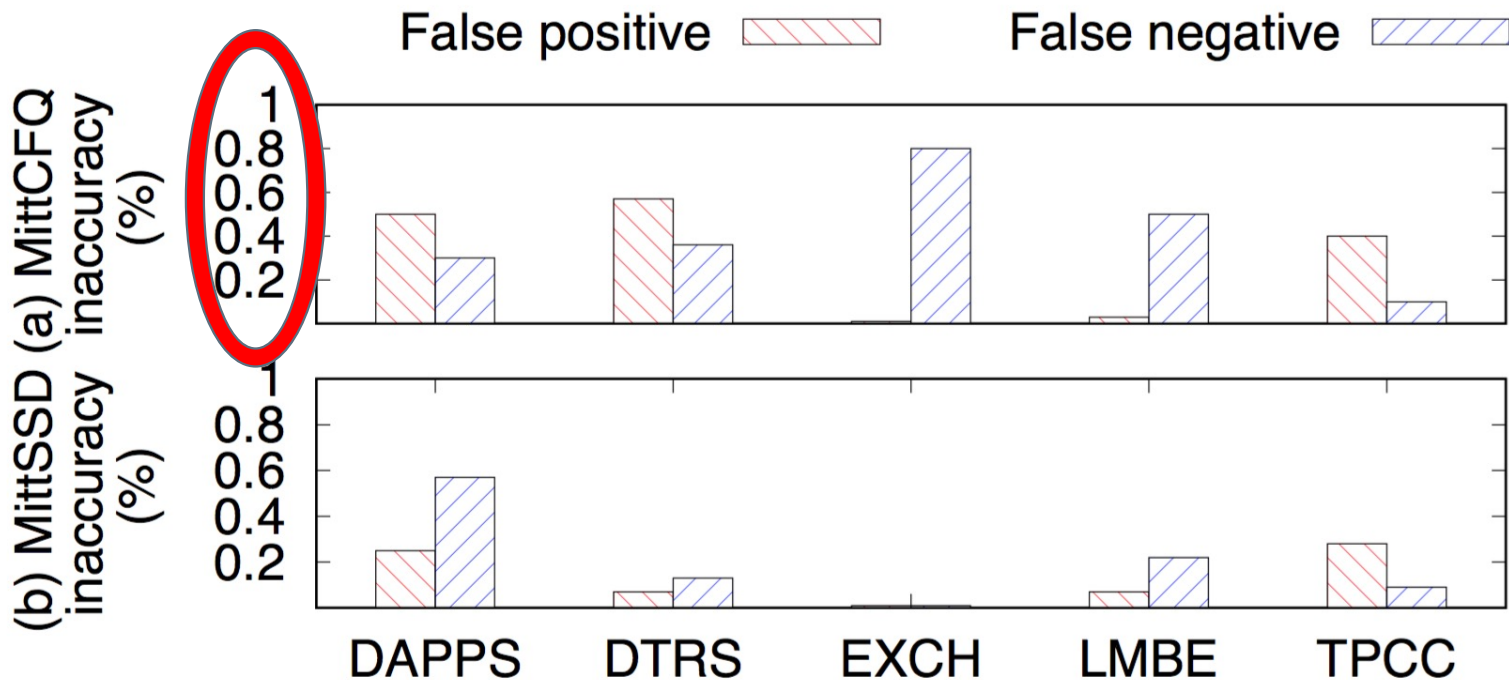
Accuracy Evaluation



Accuracy Evaluation

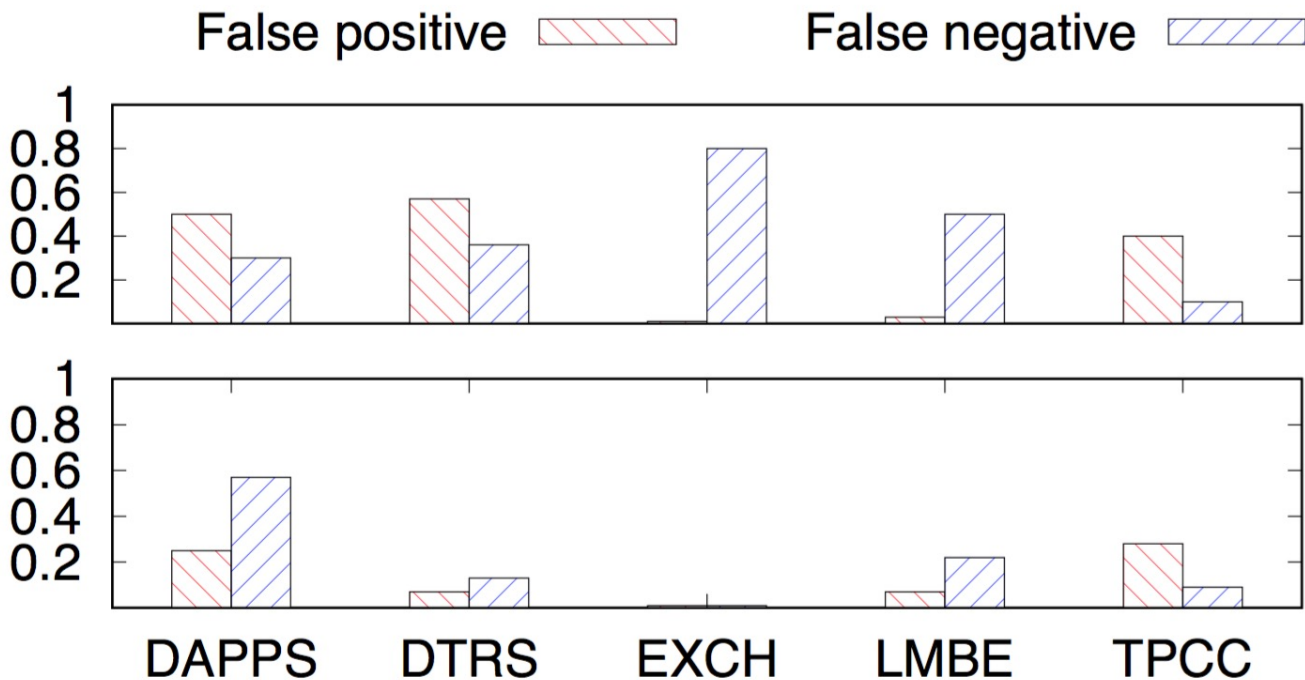


Accuracy Evaluation



Accuracy Evaluation

(b) MittSSD (a) MittCFQ
inaccuracy inaccuracy (%) (%)

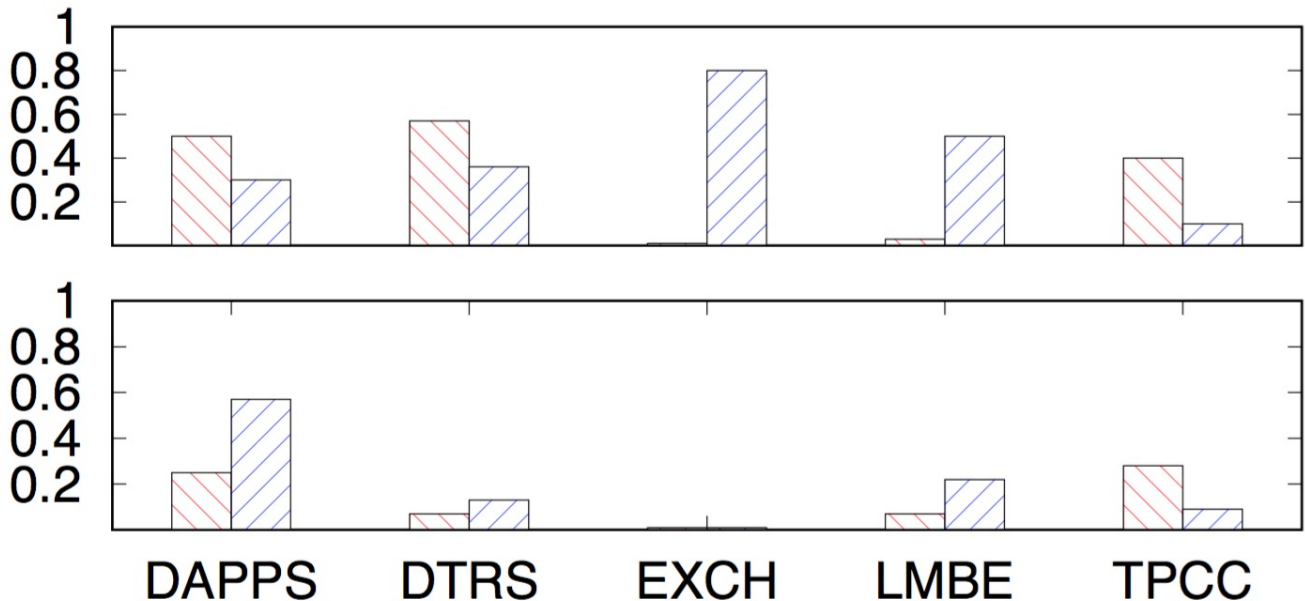


Accuracy Evaluation

Only <1% inaccuracy!

False positive  False negative 

(a) MittCFQ inaccuracy (%)
(b) MittSSD inaccuracy (%)



Accuracy Evaluation

Only <1% inaccuracy!

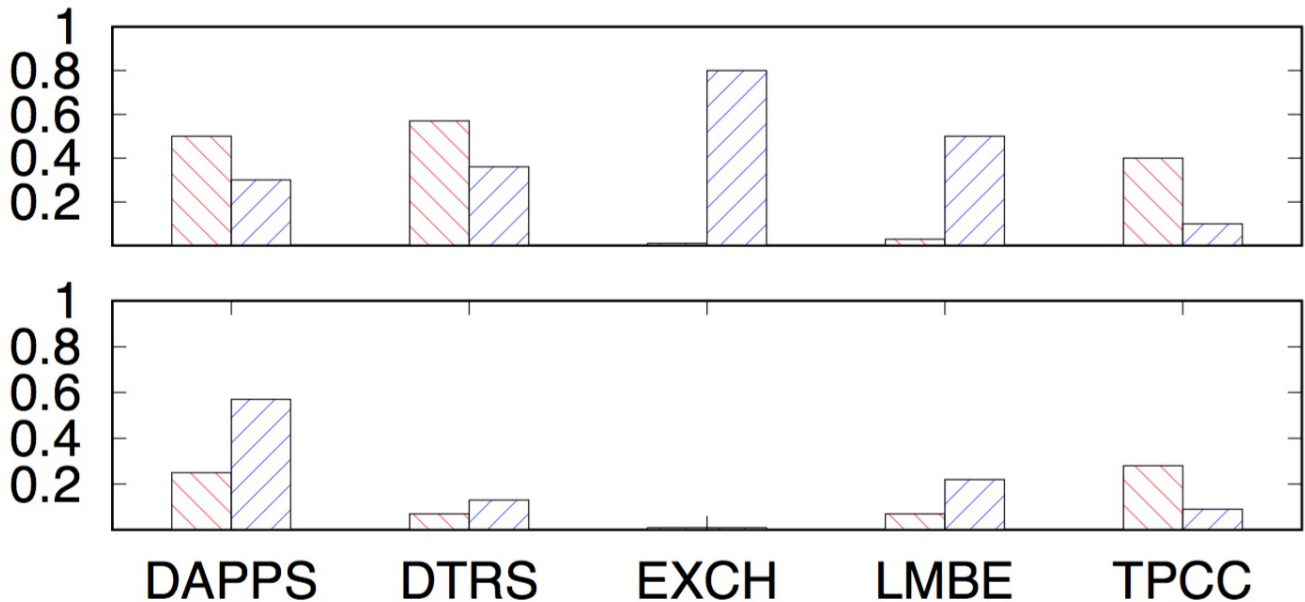
Among *incorrect* cases:

MittCFQ:
< 3ms diff

MittSSD:
< 1ms diff

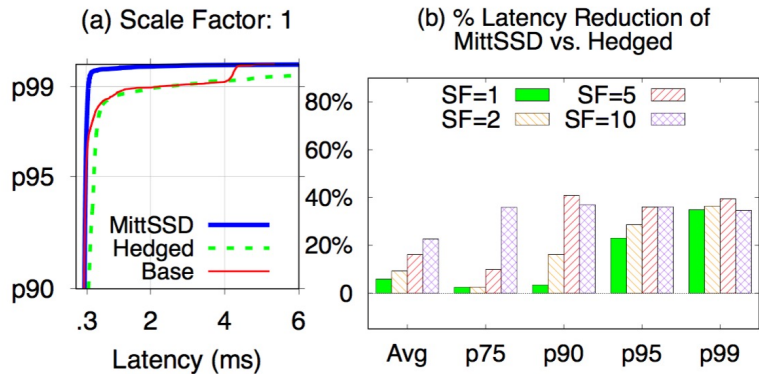
False positive  False negative 

(a) MittCFQ inaccuracy (%)
(b) MittSSD inaccuracy (%)

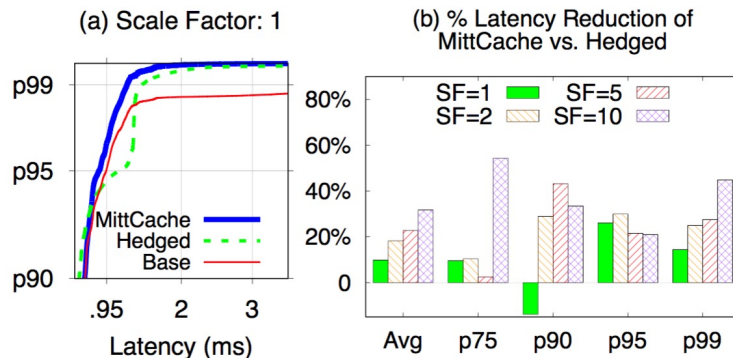




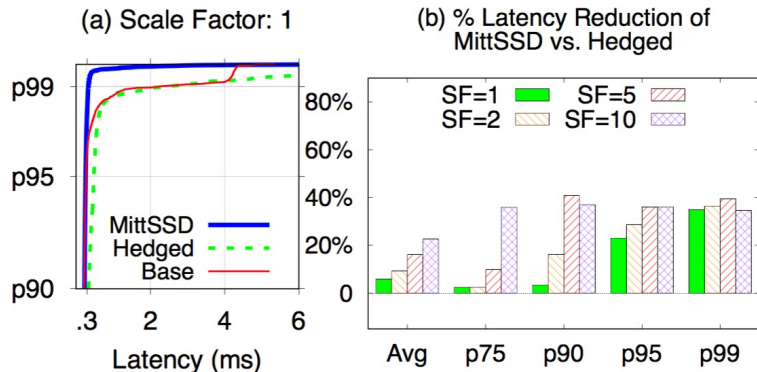
MittSSD



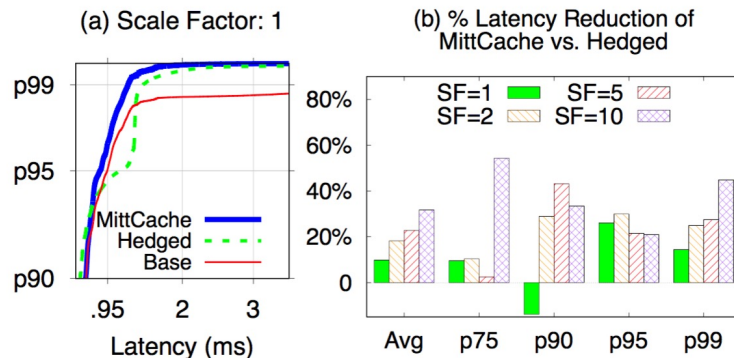
MittCache



MittSSD

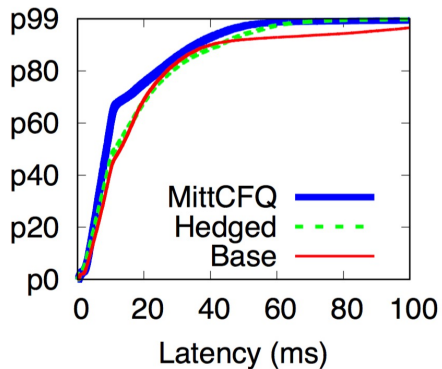


MittCache

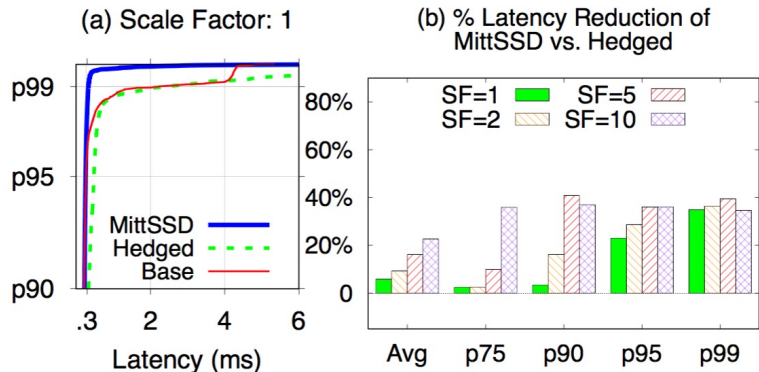


MongoDB + Filebench + Hadoop

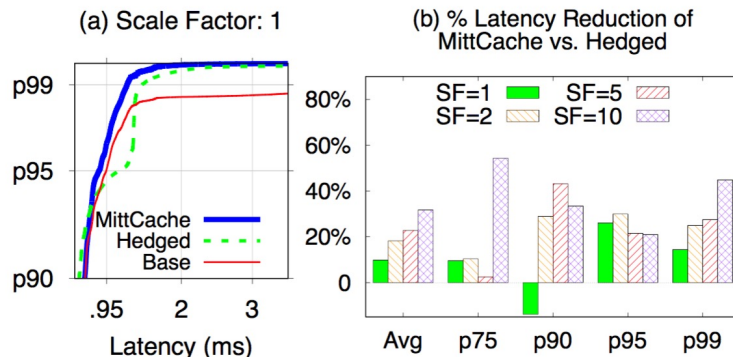
(a) Latency CDF of MongoDB



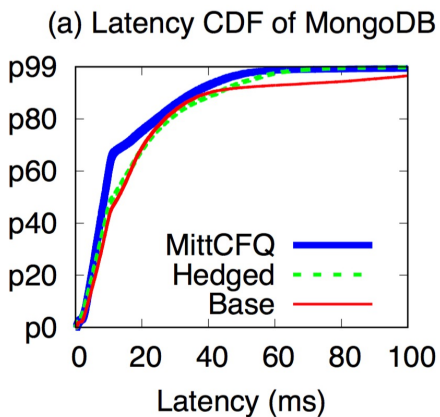
MittSSD



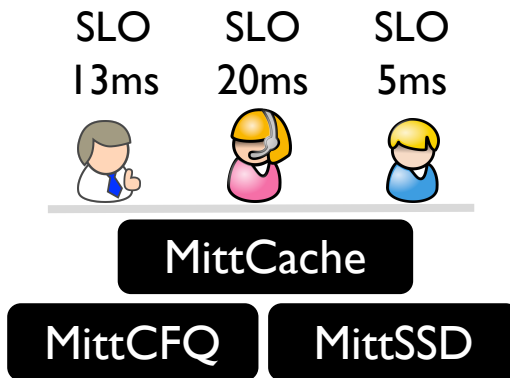
MittCache



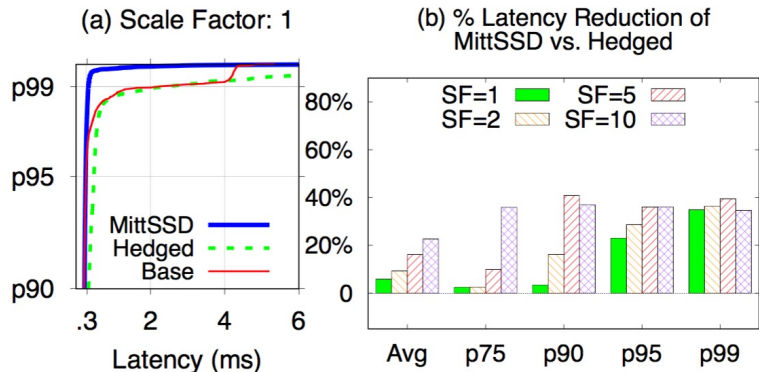
MongoDB + Filebench + Hadoop



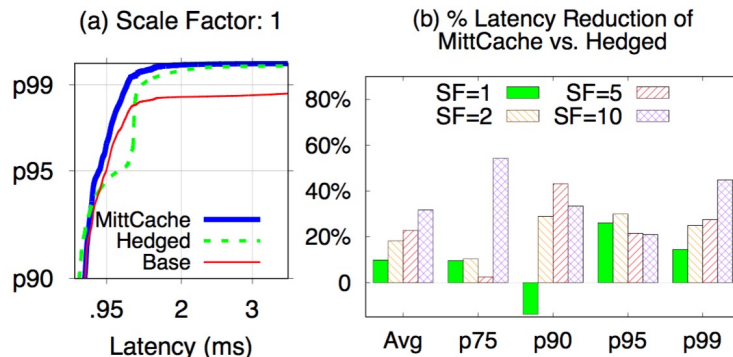
All in one



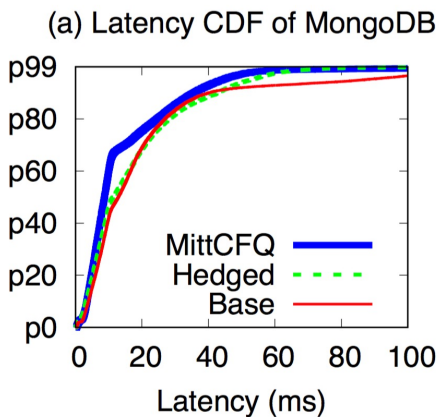
MittSSD



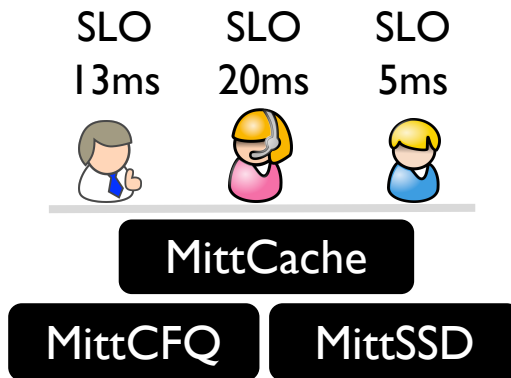
MittCache



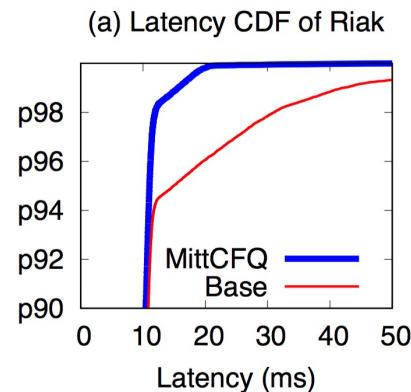
MongoDB + Filebench + Hadoop



All in one

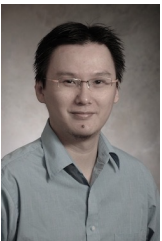


Riak

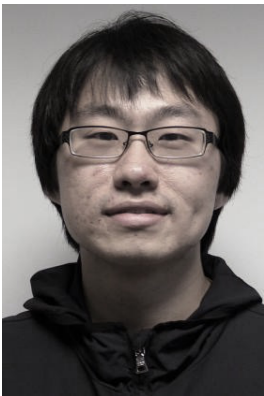




Conclusion

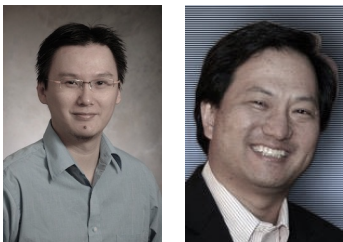


Do X





Conclusion



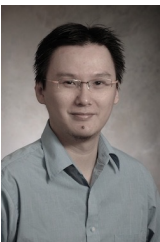
Do X



I'm busy!



Conclusion



Do X



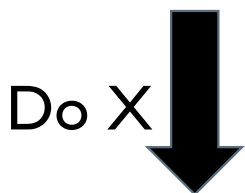
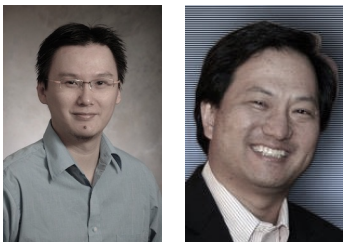
Reject!



I'm busy!



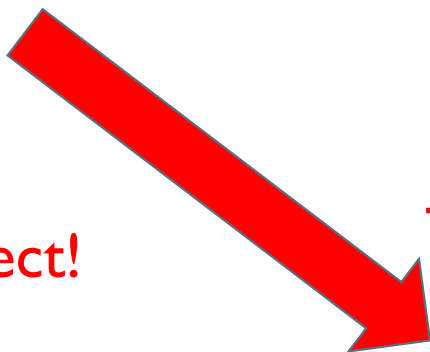
Conclusion



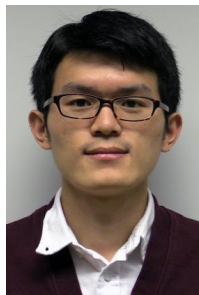
Reject!



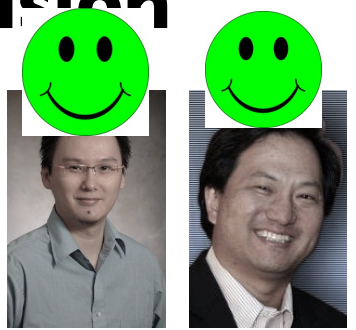
I'm busy!



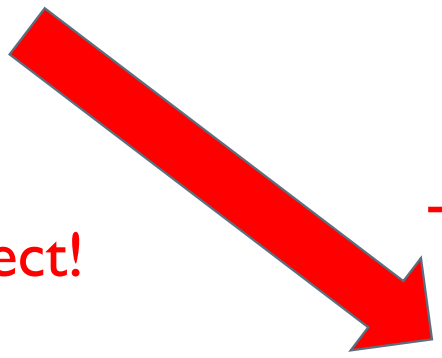
Try other students!



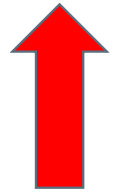
Conclusion



No wait!



Do X



Reject!

Try other students!



I'm busy!





Conclusion

Fast Reject (No-wait) Interface

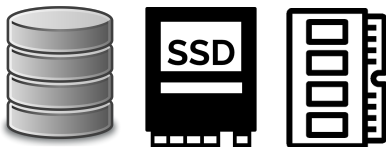
MittOS

Conclusion

Fast Reject (No-wait) Interface

MittOS

Latency Predictions





Conclusion

MittOS-
powered
apps



MongoDB



Fast Reject (No-wait) Interface

MittOS

Latency Predictions



Conclusion

MittOS-powered apps



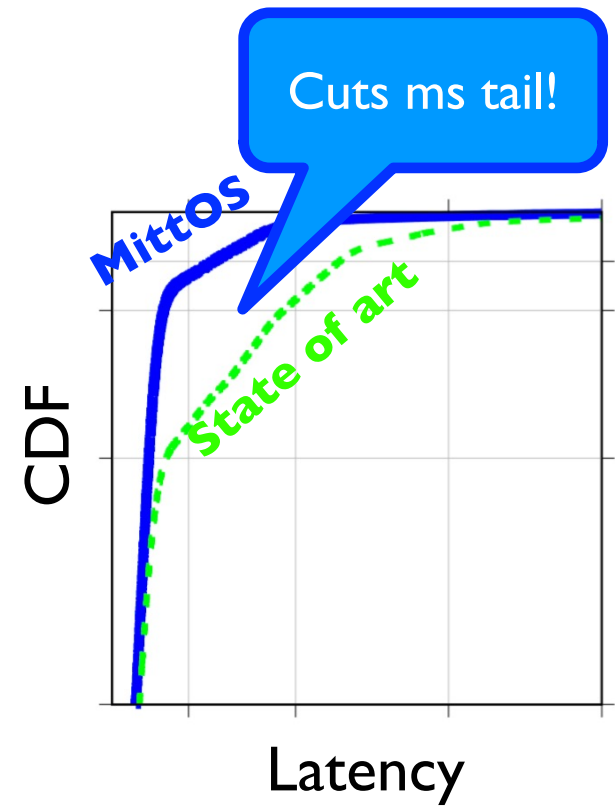
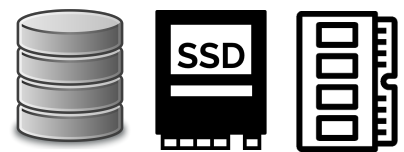
MongoDB



Fast Reject (No-wait) Interface

MittOS

Latency Predictions



Conclusion

MittOS-powered apps



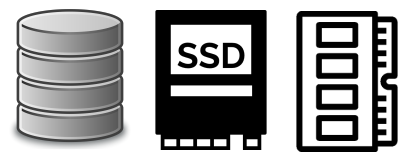
MongoDB



Fast Reject (No-wait) Interface

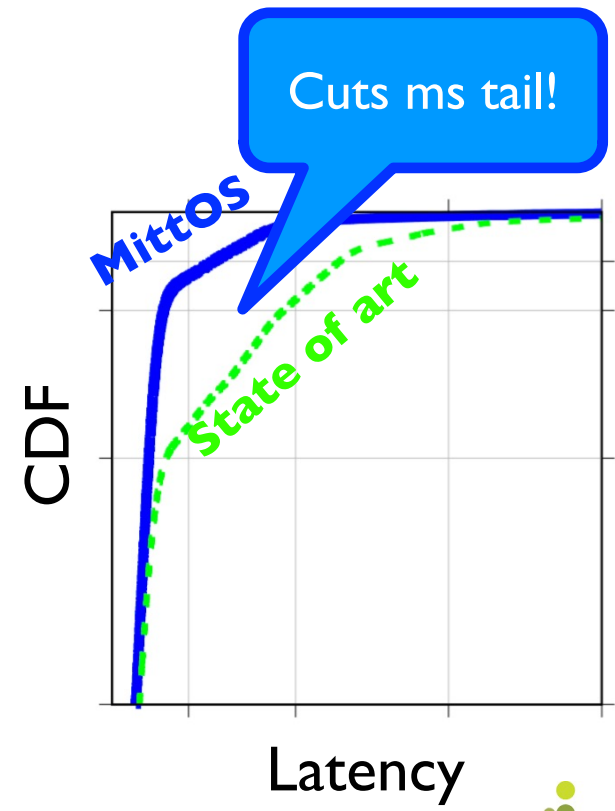
MittOS

Latency Predictions



<http://ucare.cs.uchicago.edu>

Thank you! Questions?



<https://ceres.uchicago.edu>