# 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,

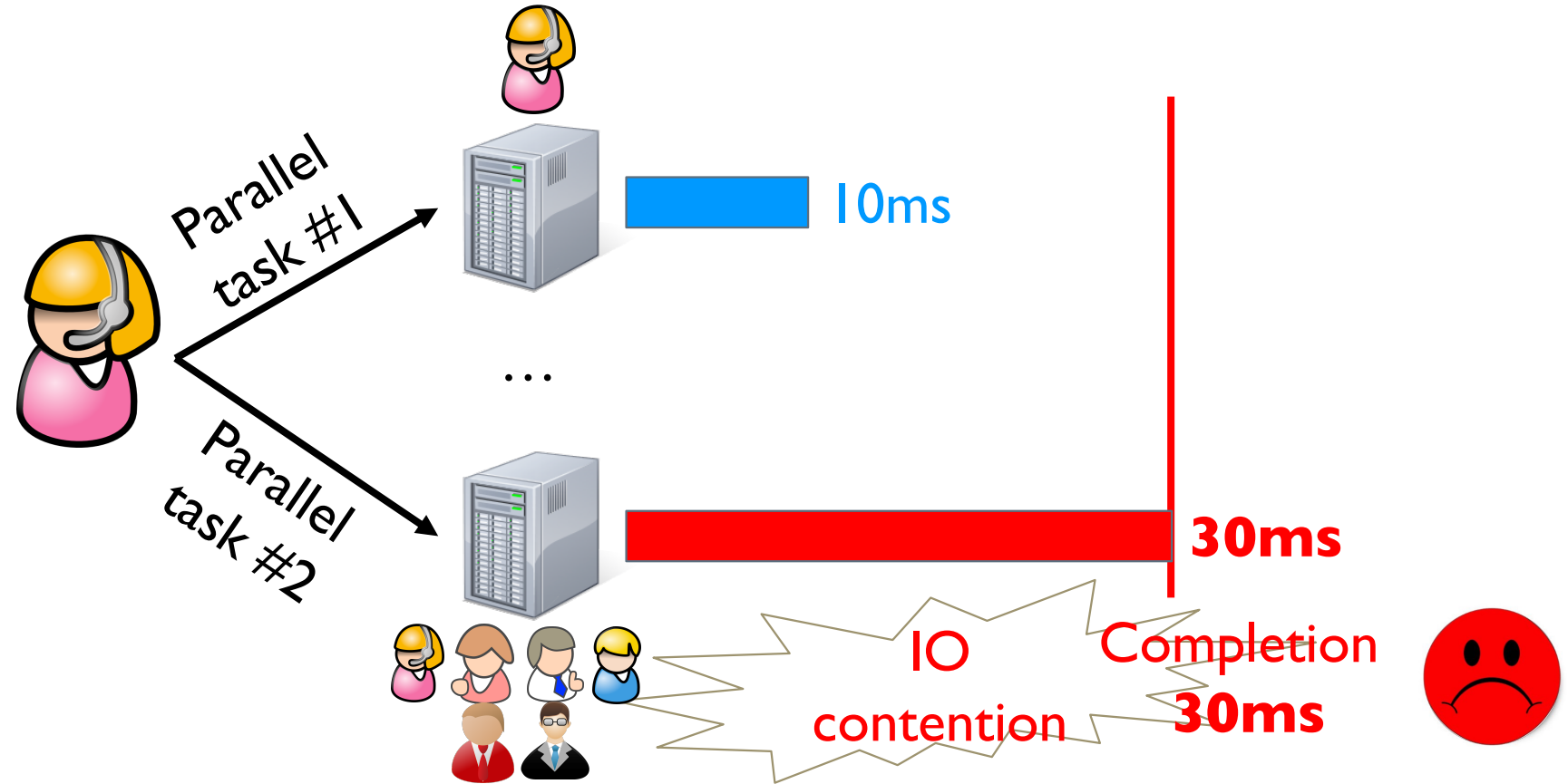Andrew A. Chien, and Haryadi S. Gunawi

THE UNIVERSITY OF CHICAGO

CERES
Center for Unstoppable Computing

# 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"
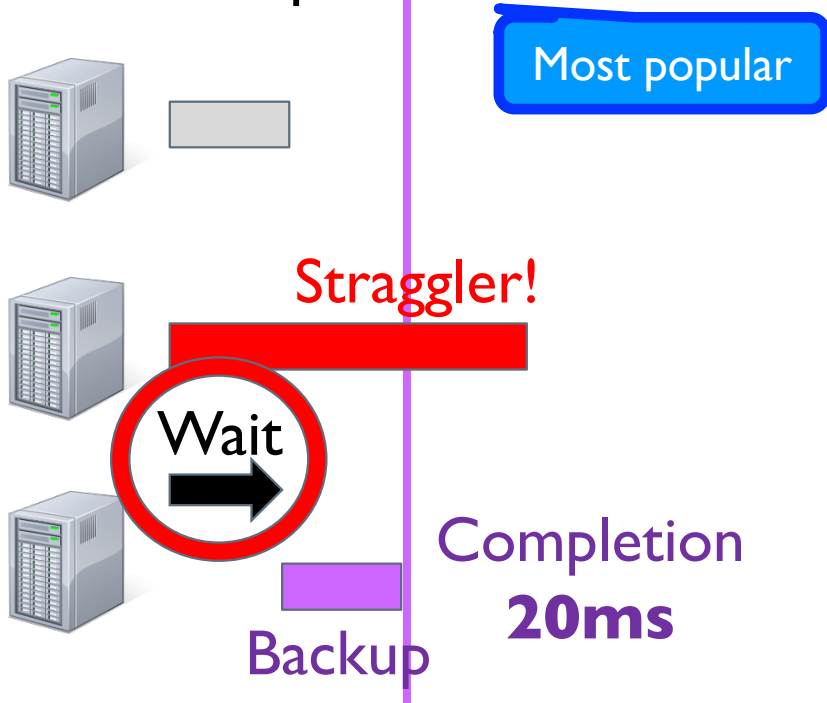
Google: "extra **500ms** in search page generation time dropped traffic by **20%**"

# Millisecond Tail Latency
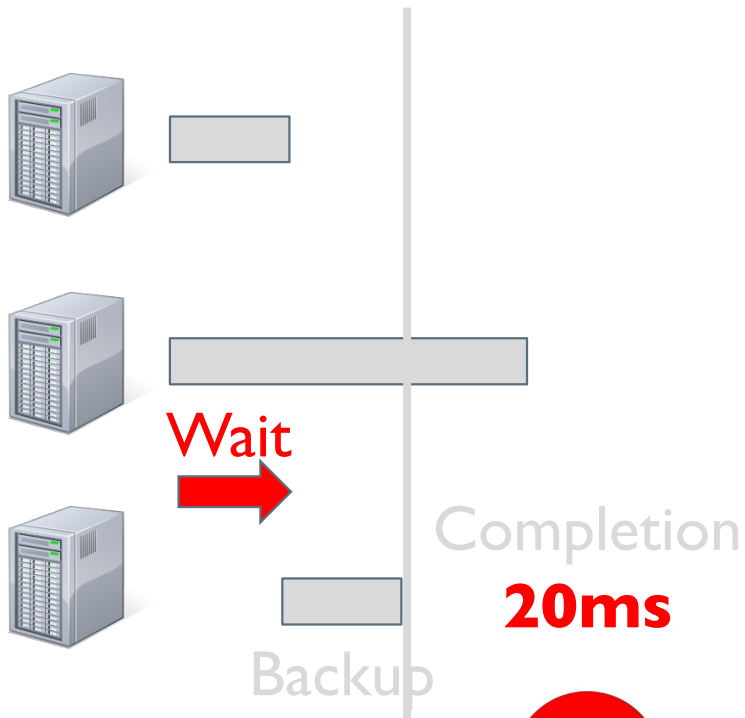
# Current Tail-Tolerant Mechanisms

1. Speculation

Most popular

Straggler!

Wait

Completion
**20ms**

Backup

2. Cloning
   - Introduces 2x workload

3. Snitching
   - Does not work when burstiness fluctuates in ms-level

# MittOS

**App** | **OS**

Disk Queue
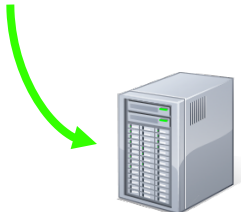
ret = read(.., < 20ms) →

if ( ret == Reject)
// failover
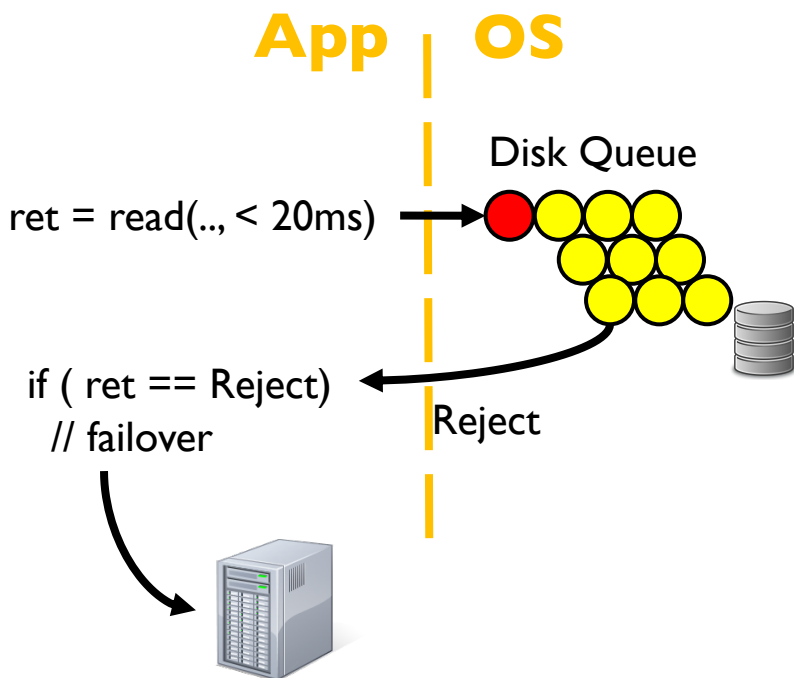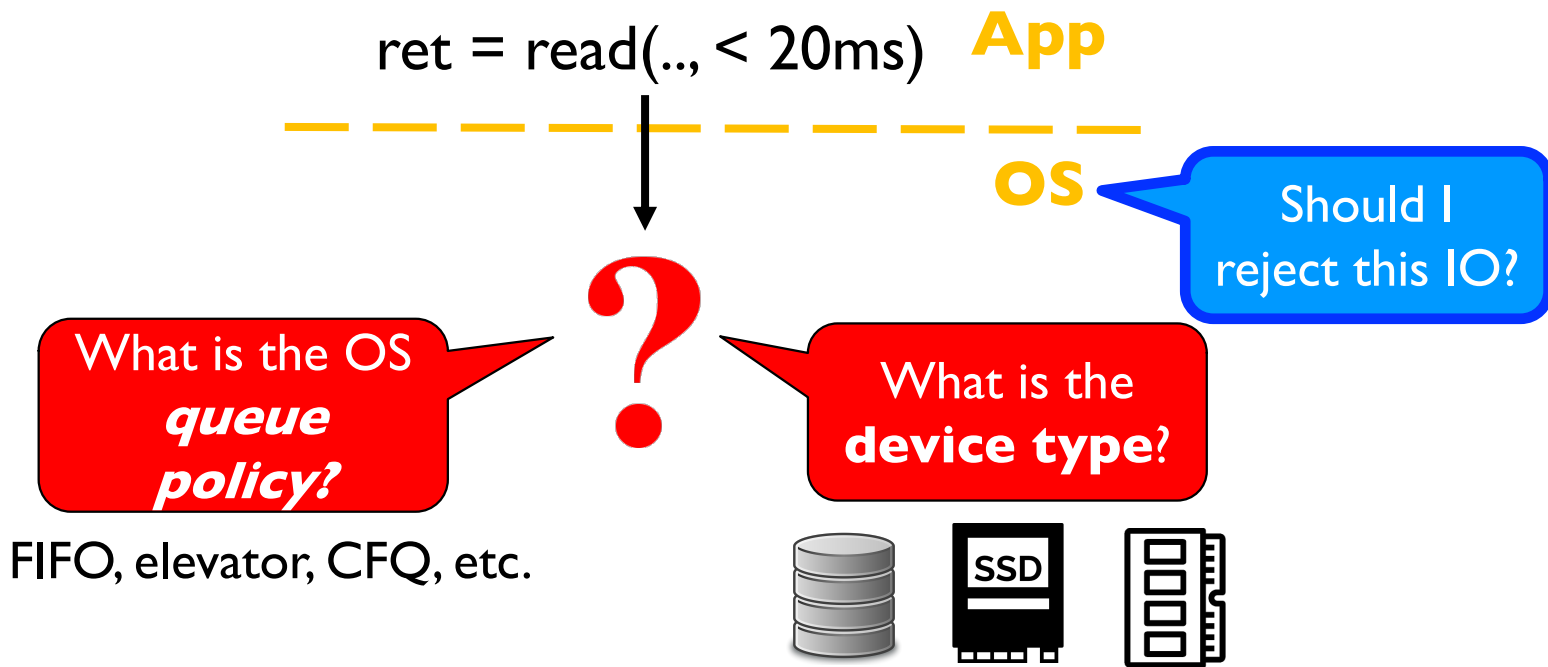
Reject

- 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

ret = read(.., < 20ms)   **App**

**OS**

Should I reject this IO?

What is the OS *queue policy?*

?

What is the **device type**?

FIFO, elevator, CFQ, etc.

SSD

Prediction depends on queue policy and device type

# Contribution     +50 LOC

MittOS-
powered
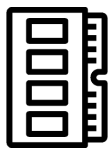


### Fast Reject Interface

## 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**

# Prediction

# Challenge #1: Modeling Queue Policy

# Challenge #2: Device Type



Reject / Accept depends on **device type**

FIFO

> 20ms

Reject

Single spindle

Disk

Parallel channels & chips

FIFO

< 20ms

Accept

SSD

...

# Challenge #2: Device Type

# Outline

❑ Introduction

❑ **Design**

   ▪ Challenges

   ▪ Solutions

❑ Evaluation

❑ Conclusion

# Reject/Latency Prediction

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

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

Simple type

Complicated type

SSD

Profiling is enough

White-box knowledge required

MittCFQ

MittSSD

# MittCFQ

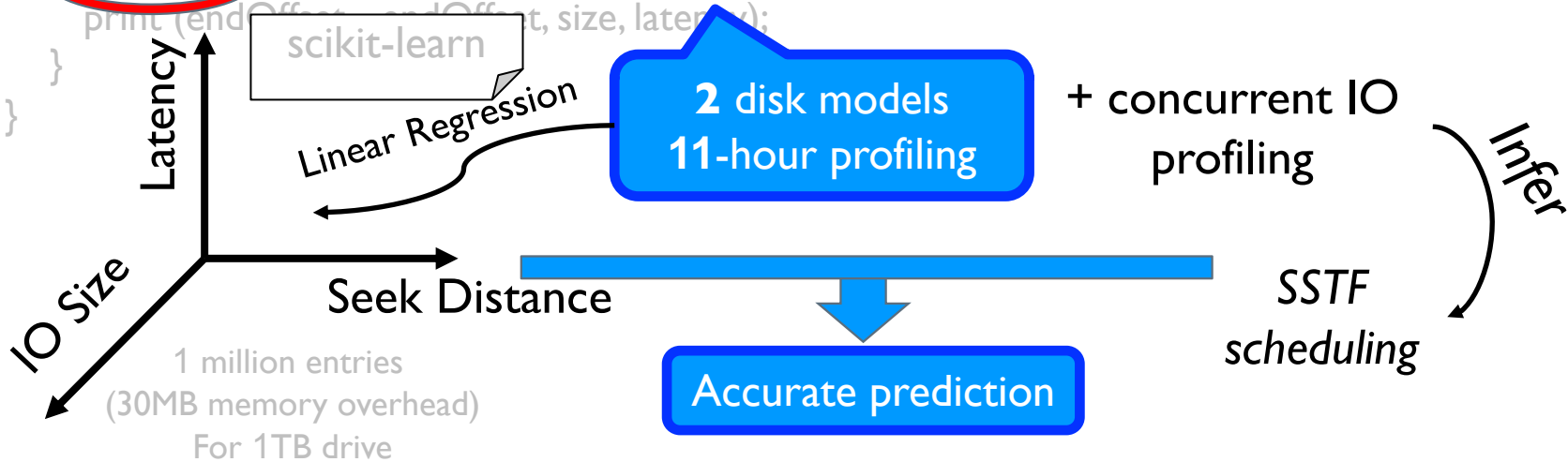# MittCFQ Profiling

```
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, endOffset, size, latency);
        }
    }
}
```
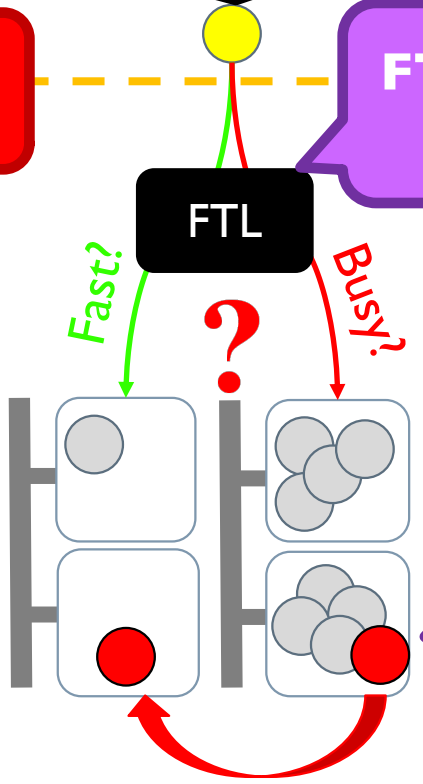
**Random seek**

**Random read**

**Collect latency**

scikit-learn

Latency

IO Size

Seek Distance

Linear Regression

**2** disk models
**11**-hour profiling

+ concurrent IO
profiling

*Infer*

*SSTF
scheduling*

1 million entries
(30MB memory overhead)
For 1TB drive

Accurate prediction

# 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

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

Reverse engineering based on source code

Simple type

Complicated type

SSD

Profiling is enough

MittCFQ ~**1800 LOC**

White-box knowledge required

MittSSD ~**1400 LOC**

LightNVM + Open-Channel SSD

# Other Solved Challenges

- Prediction overhead optimizations
  - Avoids going through every IO in the queue
  - Reduces overhead from O(n) to roughly O(1)
  - Shows < 5µs overhead for MittCFQ prediction
  - < 300ns for MittSSD prediction

- MittCache
  - Prediction for OS Cache

# Outline

❑ Introduction

❑ Design

❑ **Evaluation**

   ▪ Tail reduction

   ▪ Latency prediction accuracy

❑ Conclusion

# MittCFQ-powered MongoDB



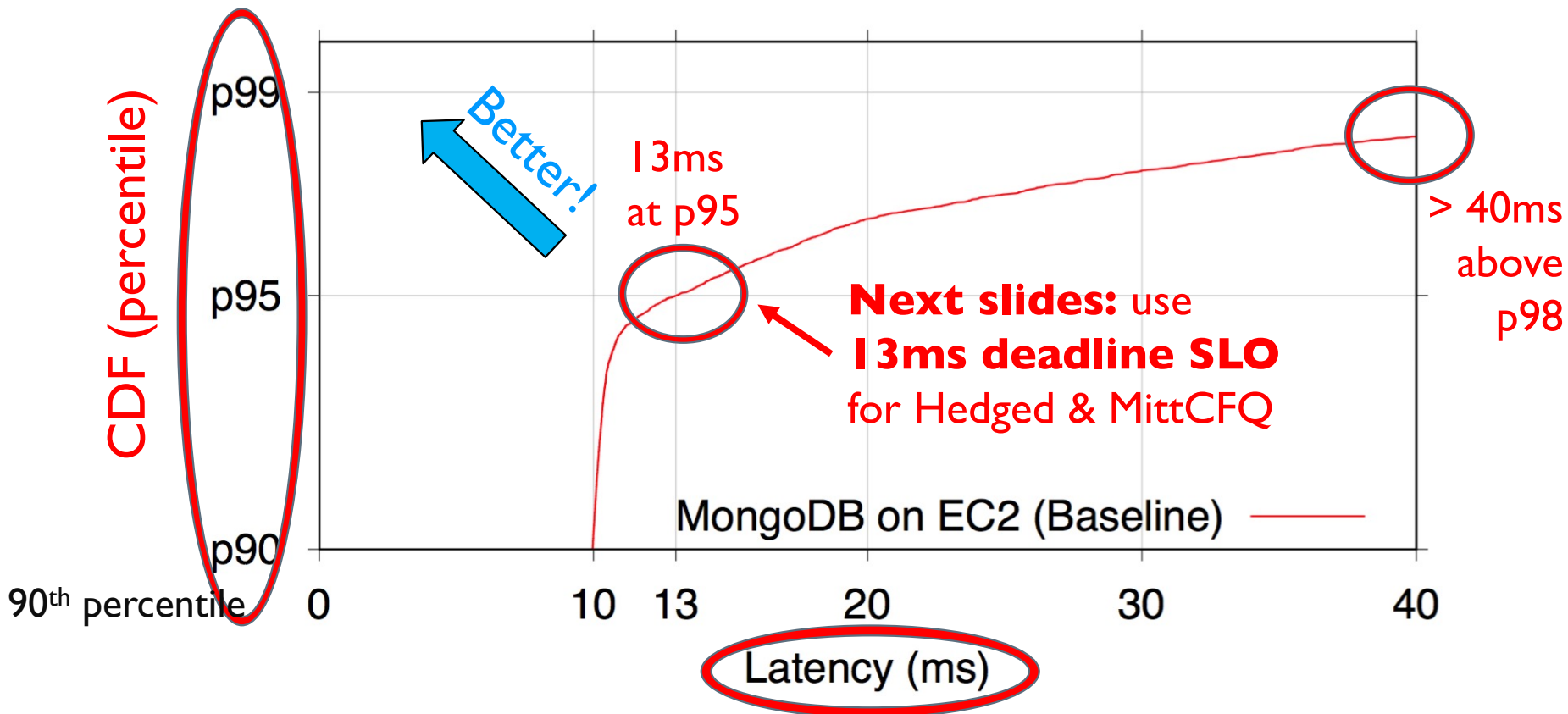Remote YCSB client #1

get()

Can failover 3 replicas

Physical node #1

Client #2 · · · Client #20

get()

get()

Node #2

· · ·

Node #20

Noisy neighbors based on EC2 data

**Metric**:
CDF of all **get()**
requests latencies
(total 6 million
data points)

# Baseline



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

CDF (percentile)

p99

p95

p90

90th percentile

Better!

13ms at p95

> 40ms above p98

**Next slides:** use **13ms deadline SLO** for Hedged & MittCFQ

MongoDB on EC2 (Baseline)

0    10  13    20    30    40

Latency (ms)

# Clone



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

Tail reduction

Worse performance < p95

Clone
"Baseline"

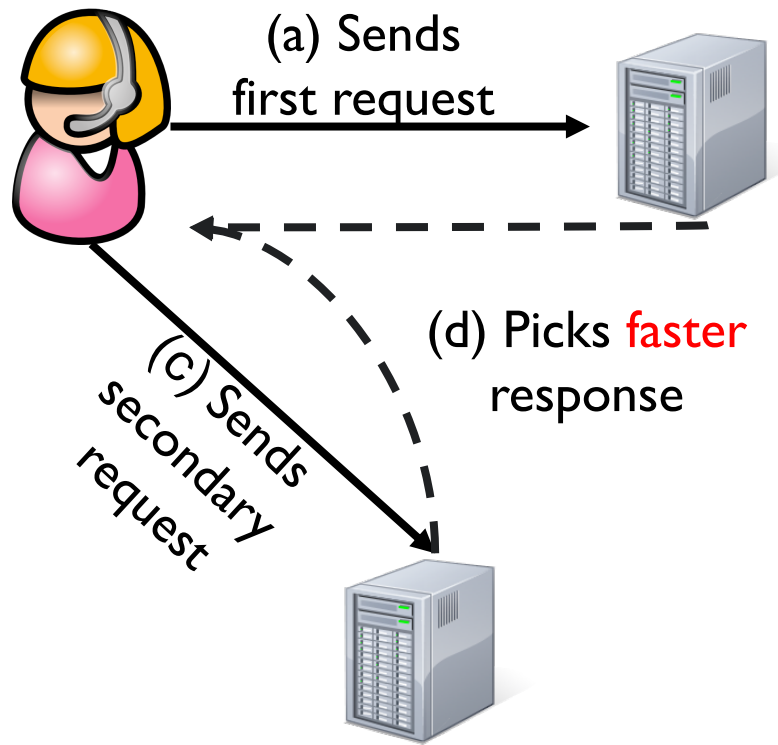Latency (ms)

# Hedged Requests

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

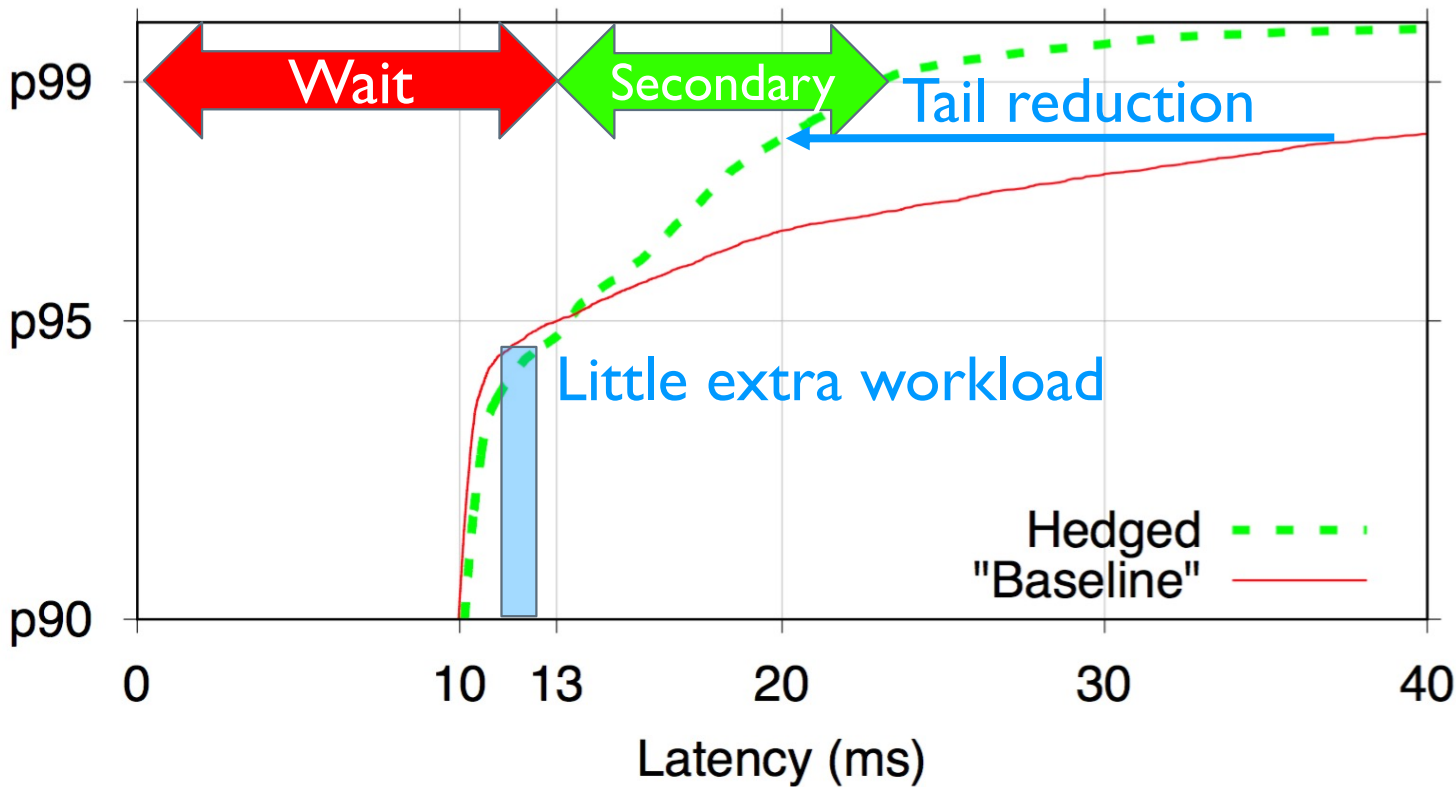BY JEFFREY DEAN AND LUIZ ANDRÉ BARROSO

# The Tail at Scale

*Communications of the ACM, vol. 56 (2013), pp. 74-80*
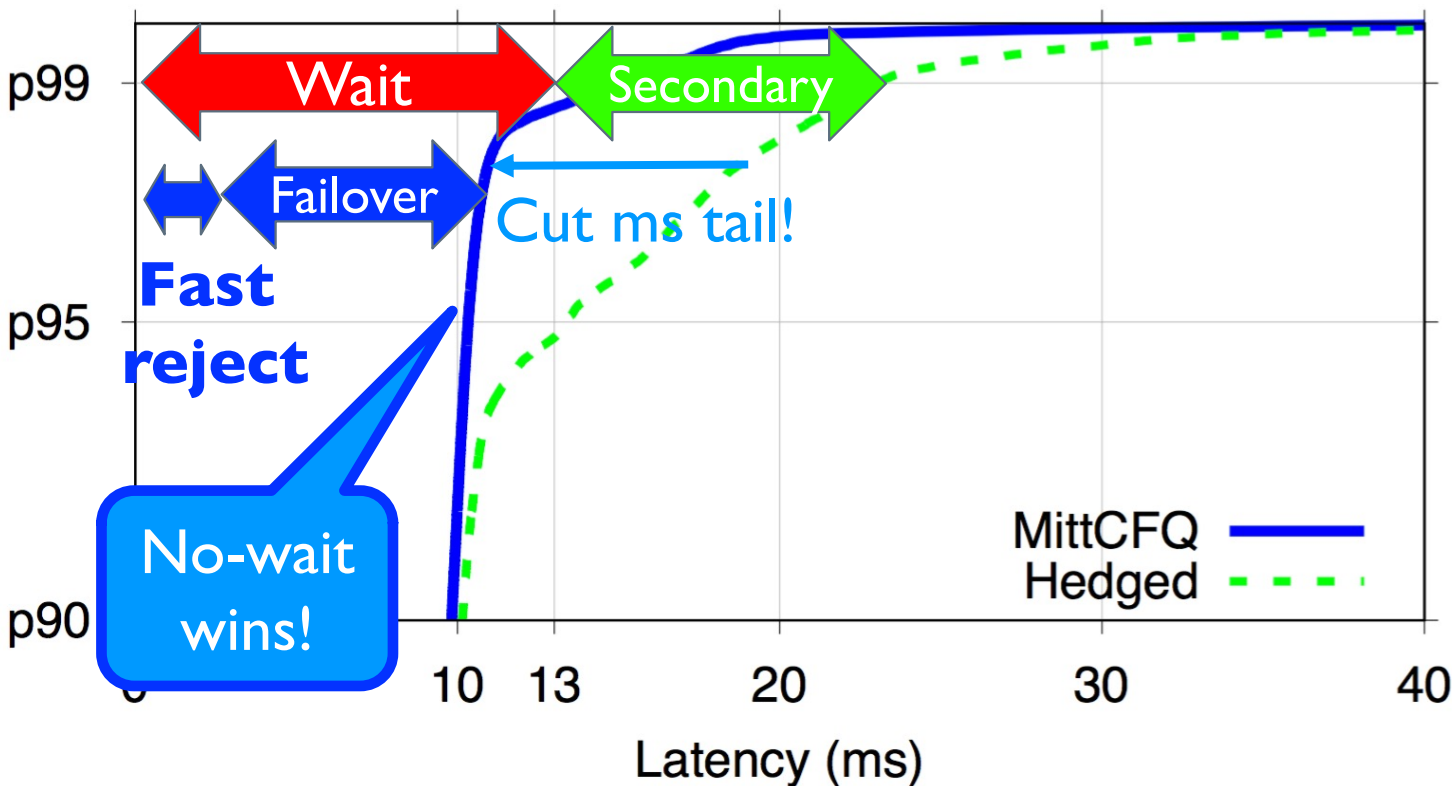
(b) Waits for 13ms **timeout**

(a) Sends first request

(c) Sends secondary request

(d) Picks **faster** response

THE UNIVERSITY OF CHICAGO



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

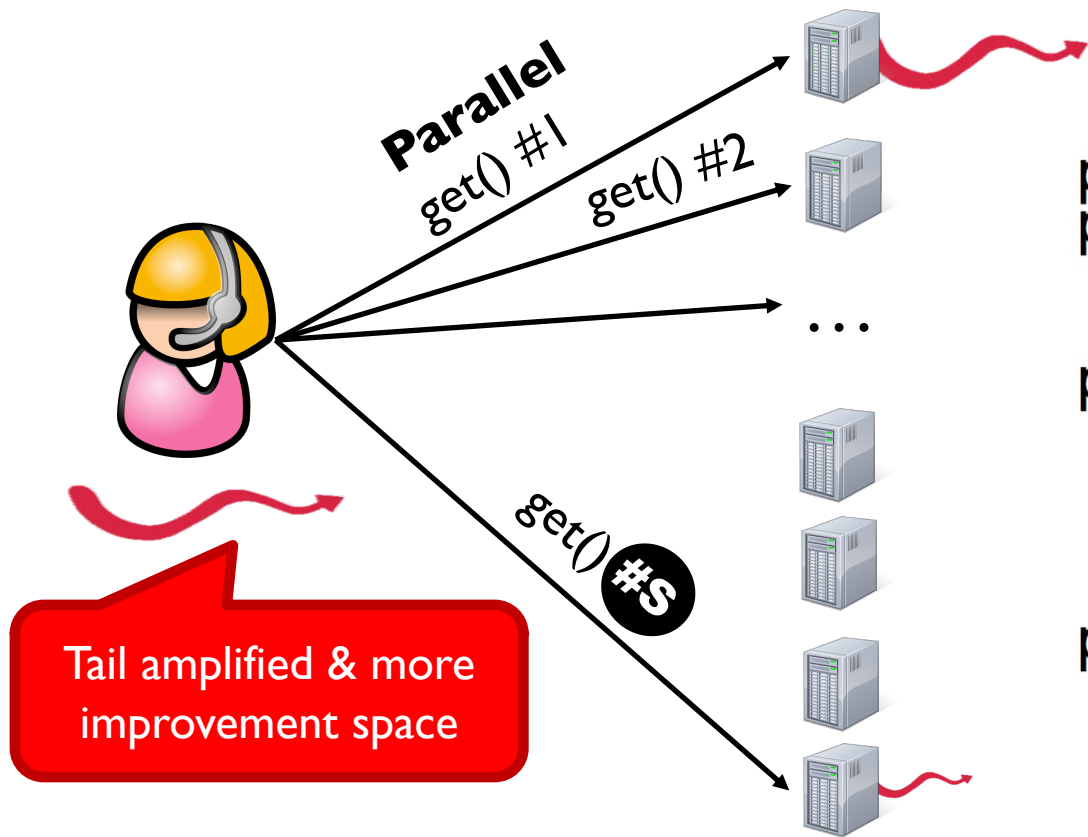# MittCFQ
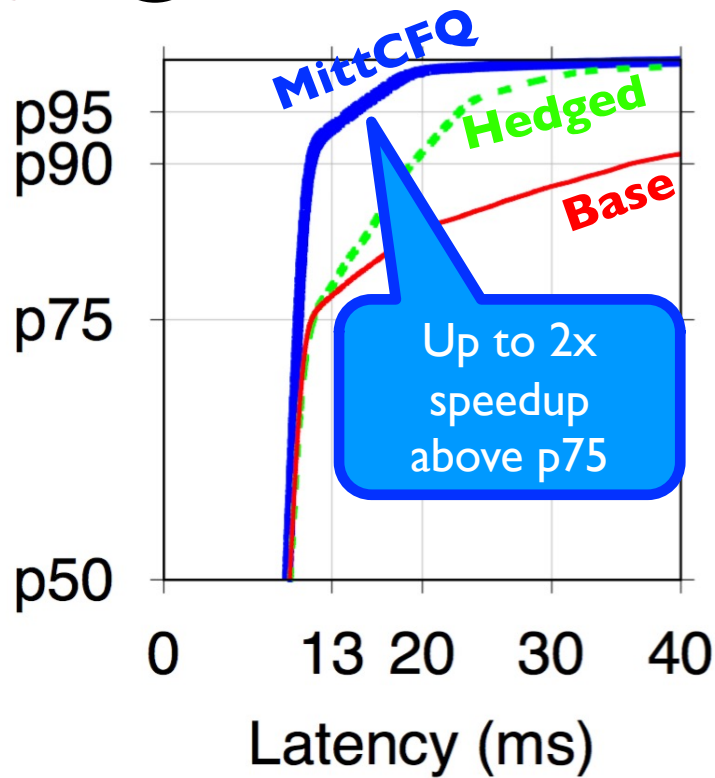


CDF of YCSB get() Latencies on 20-node MongoDB

# Tail amplified at Scale

# Accuracy Evaluation

**MittCFQ**     **MittSSD**

Disk

Open-Channel
SSD

**5** real-world block-level traces

DAPPS

DTRS                TPCC

EXCH    LMBE

Metrics:

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

# Accuracy Evaluation

Among *incorrect* cases:

MittCFQ:          MittSSD:
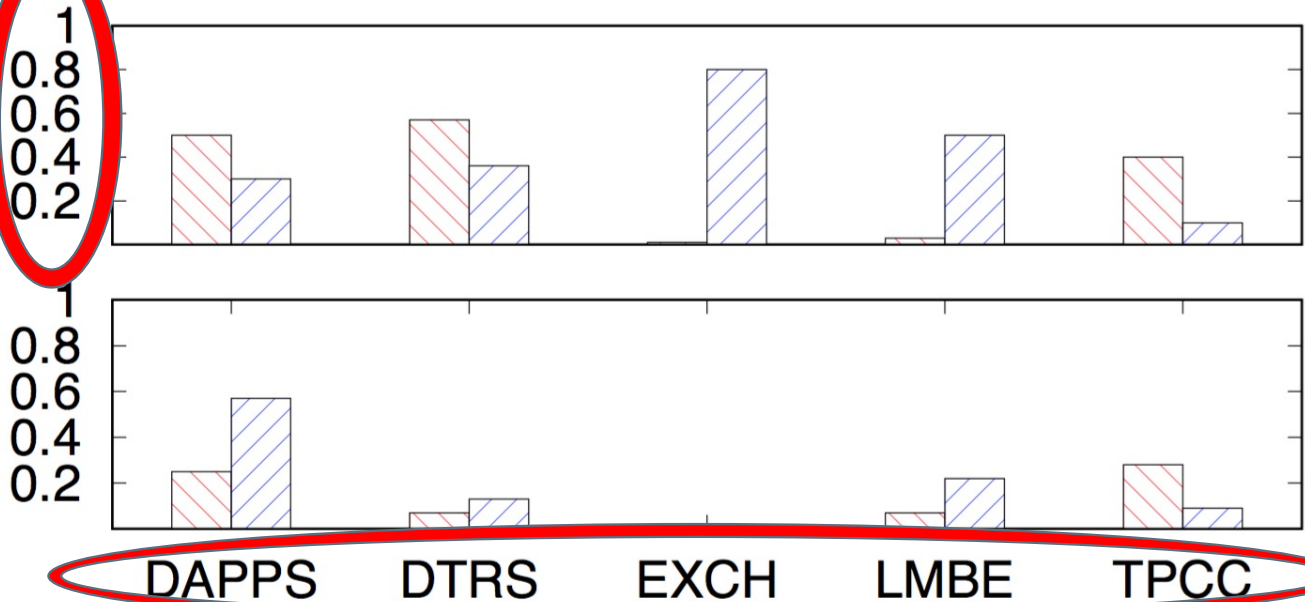< 3ms diff       < 1ms diff

Only **<1%** inaccuracy!



False positive          False negative

# MittSSD



(a) Scale Factor: 1

(b) % Latency Reduction of MittSSD vs. Hedged

# MittCache



(a) Scale Factor: 1

(b) % Latency Reduction of MittCache vs. Hedged

## MongoDB + Filebench + Hadoop



(a) Latency CDF of MongoDB

## All in one

SLO 13ms    SLO 20ms    SLO 5ms

MittCache

MittCFQ    MittSSD

## Riak



(a) Latency CDF of Riak

# Conclusion

# Conclusion

MittOS-
powered
**apps**

Fast Reject (No-wait) Interface

**MittOS**   Latency Predictions

Cuts ms tail!

MittOS

State of art

CDF

Latency

**Thank you! Questions?**