

# TAILBENCH: A BENCHMARK SUITE AND EVALUATION METHODOLOGY FOR LATENCY- CRITICAL APPLICATIONS

**HARSHAD KASTURE**, DANIEL SANCHEZ

IISWC 2016

[tailbench.csail.mit.edu](http://tailbench.csail.mit.edu)



**Massachusetts  
Institute of  
Technology**



# Executive Summary

- Latency-critical applications have stringent performance requirements → low datacenter utilization
  - ▣ Wastes billions of dollars in energy and equipment annually
  
- Research in this area hampered by the lack of a comprehensive benchmark suite
  - ▣ Few latency-critical applications → limited coverage
  - ▣ Complicated setup and configuration
  - ▣ Methodological issues

} Inaccurate latency measurements
  
- TailBench makes latency-critical applications easy to analyze
  - ▣ Varied application domains and latency characteristics
  - ▣ Standardized, statistically sound methodology
  - ▣ Supports simplified load-testing configurations

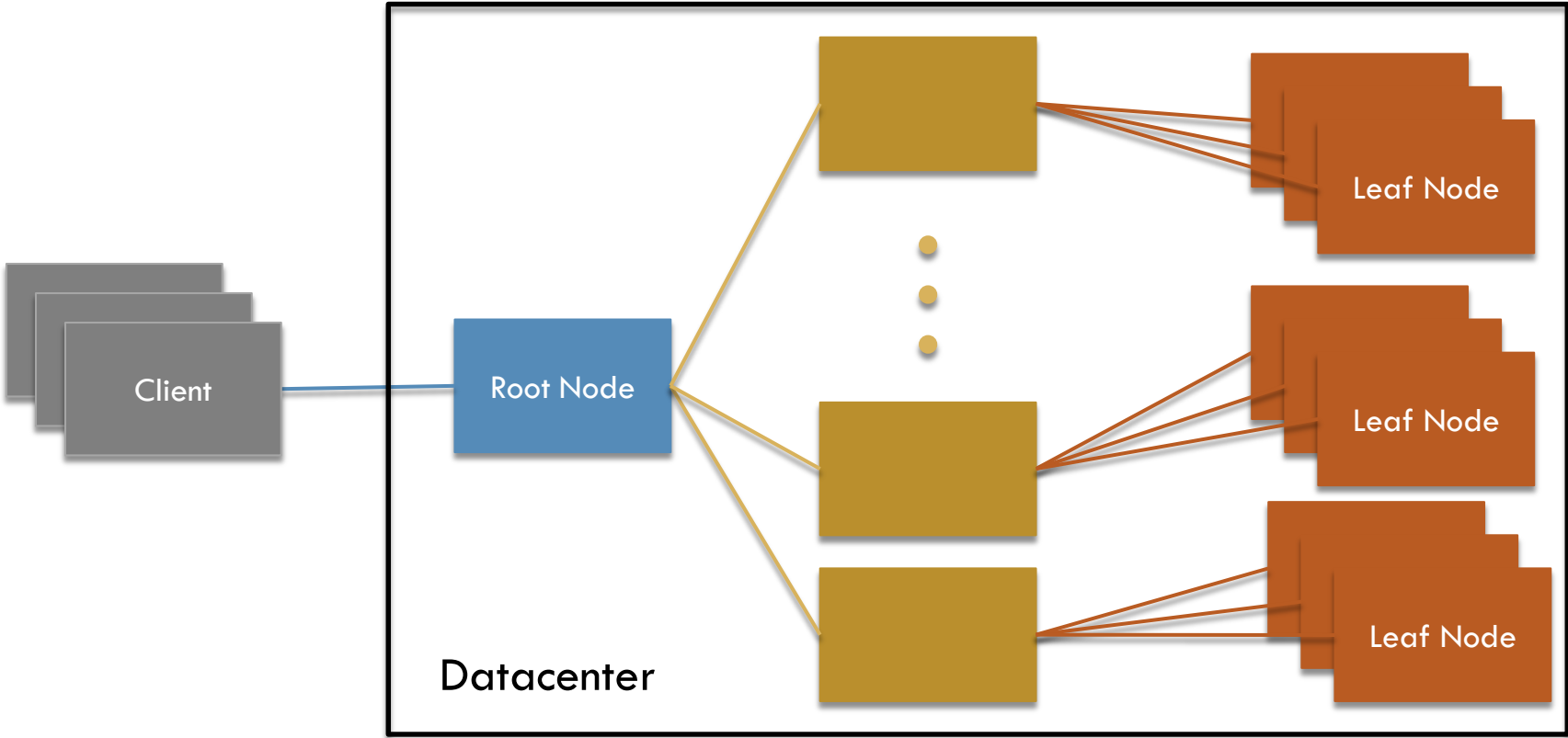
- Background and Motivation

- TailBench Applications

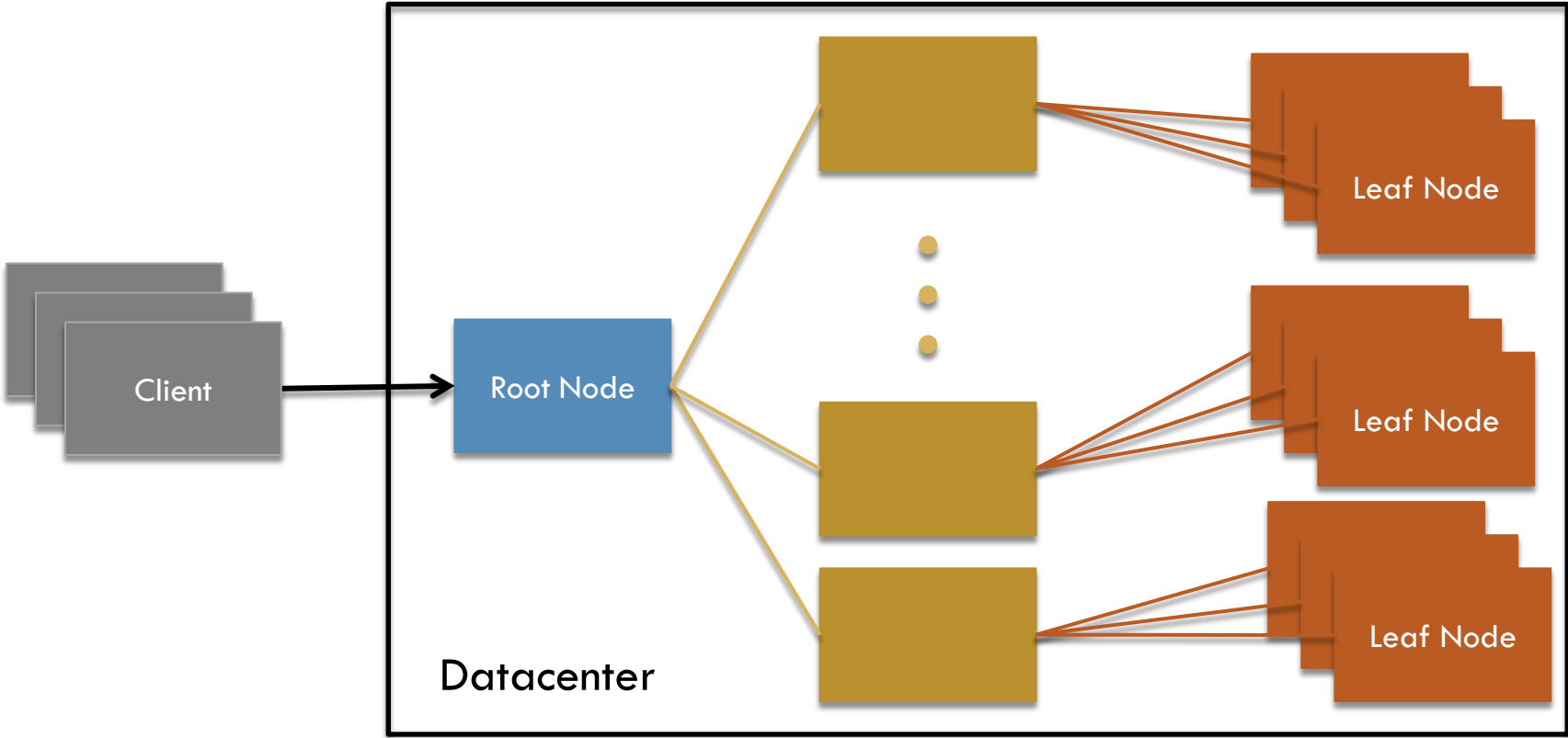
- TailBench Harness

- Simplified Configurations

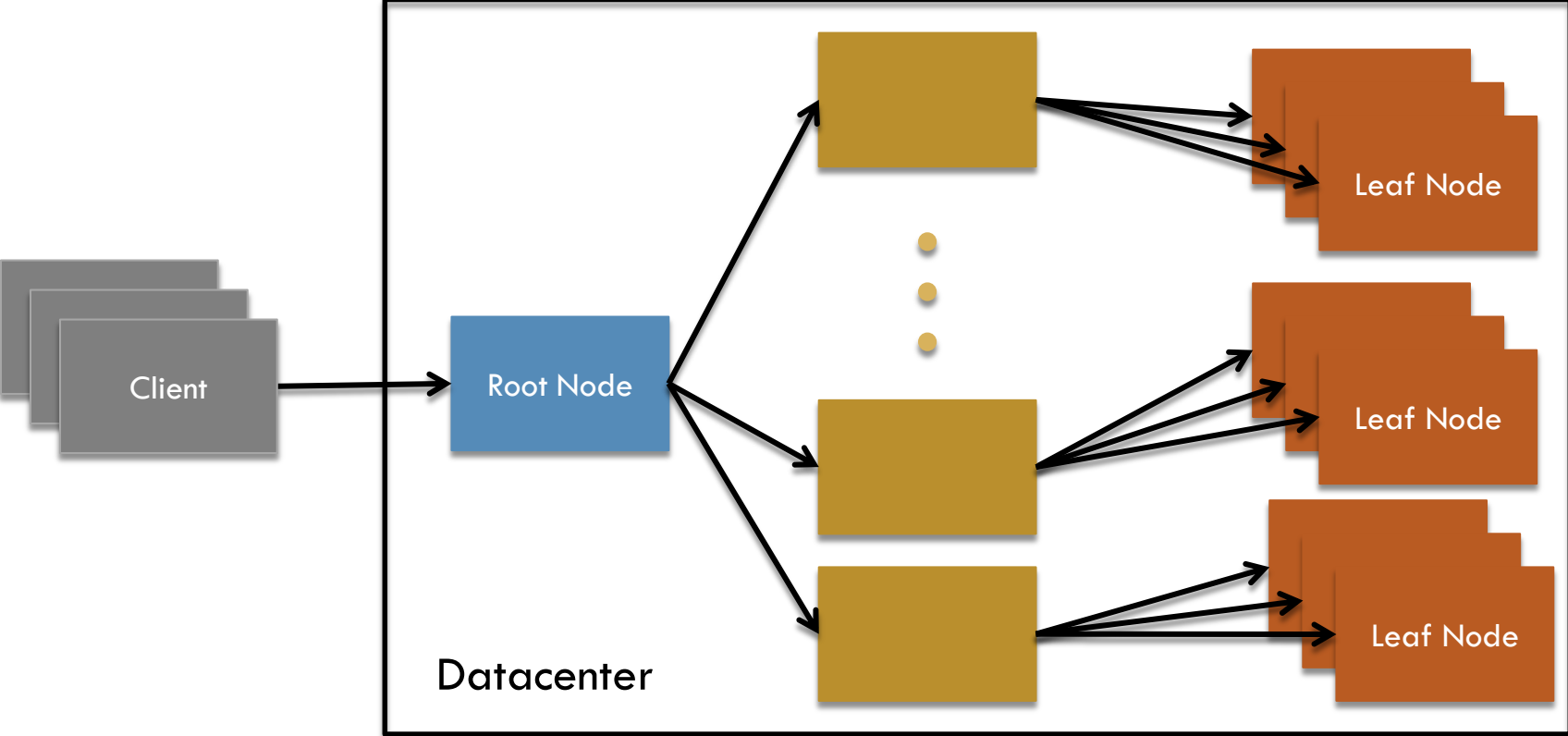
# Understanding Latency-Critical Applications <sub>4</sub>



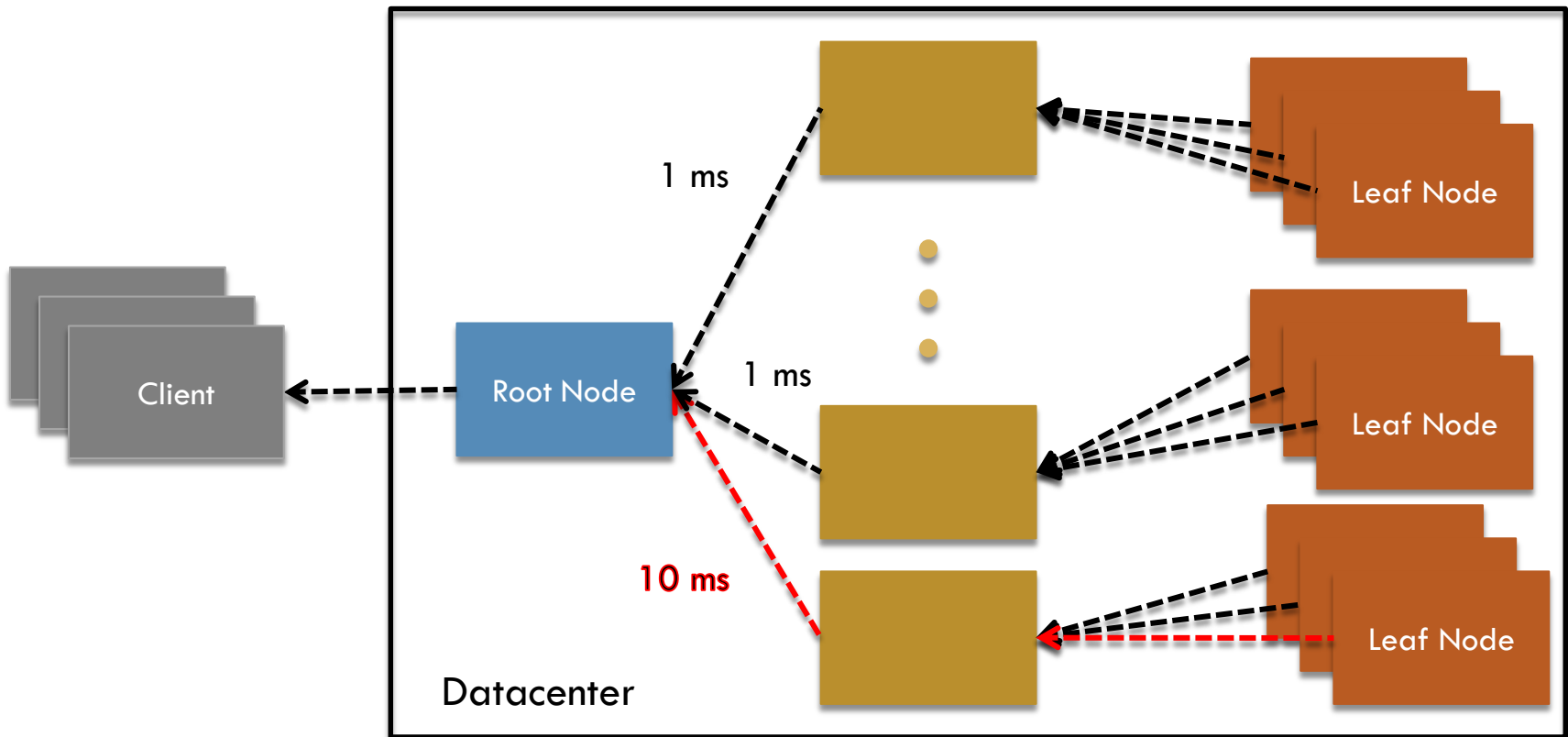
# Understanding Latency-Critical Applications <sub>5</sub>



# Understanding Latency-Critical Applications <sub>6</sub>

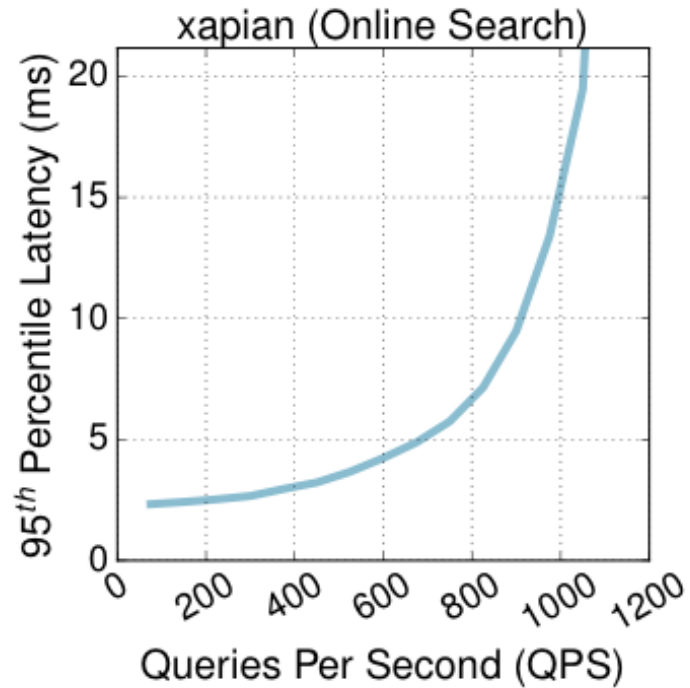


# Understanding Latency-Critical Applications <sub>7</sub>



- The few slowest responses determine user-perceived latency
  - ▣ Tail latency (e.g., 95<sup>th</sup> / 99<sup>th</sup> percentile), not mean latency, determines performance

# Latency Requirements Cause Low Utilization



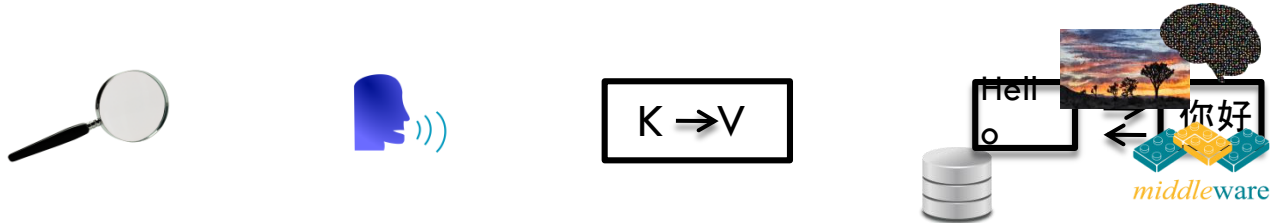
- End-to-end latency increases rapidly with load
  - ▣ Must keep utilization low to keep latency within reasonable bounds
- Traditional resource management techniques (e.g., colocation) often cannot be used since they degrade latency
- Low resource utilization wastes billions of dollars in energy and equipment
  - ▣ Sparked research in latency-critical systems



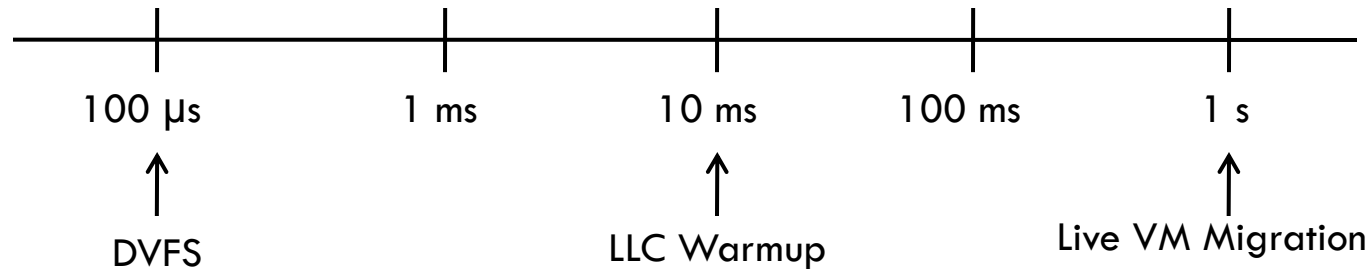
# Benchmark Suite Design Goals

9

- Applications from a diverse set of domains



- Applications with diverse tail latency characteristics



- Easy to set up and run

- Support different measurement scenarios

- Robust latency measurement methodology

- Background and Motivation

- TailBench Applications

- TailBench Harness

- Simplified Configurations

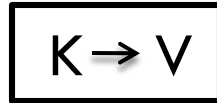
# TailBench Applications

xapian



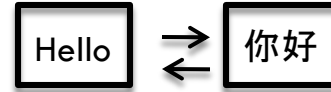
Online Search

masstree



Key-Value Store

moses



Statistical Machine Translation

sphinx



Speech Recognition

img-dnn



Image Recognition

specjbb



Java Middleware

silo



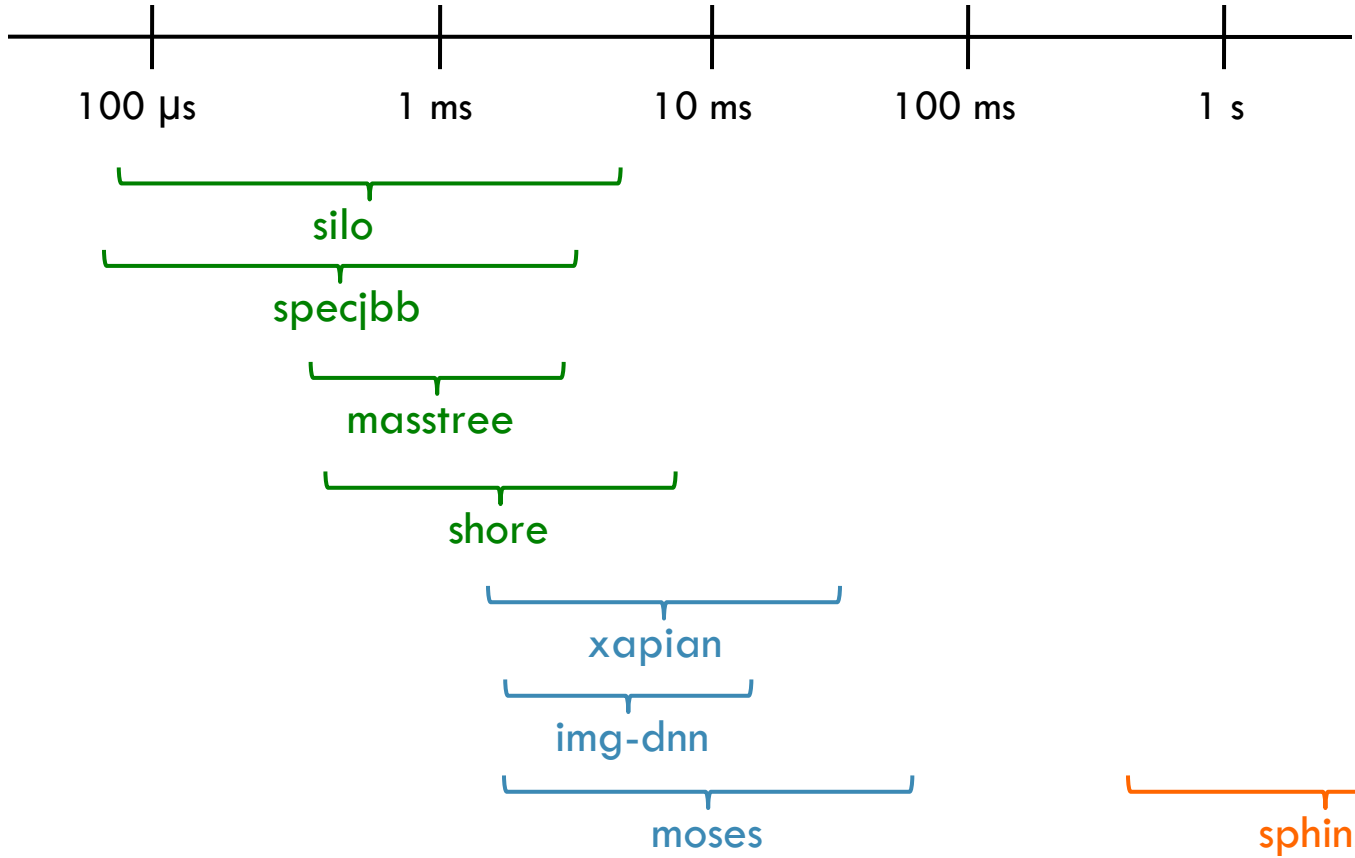
In-memory Database

shore



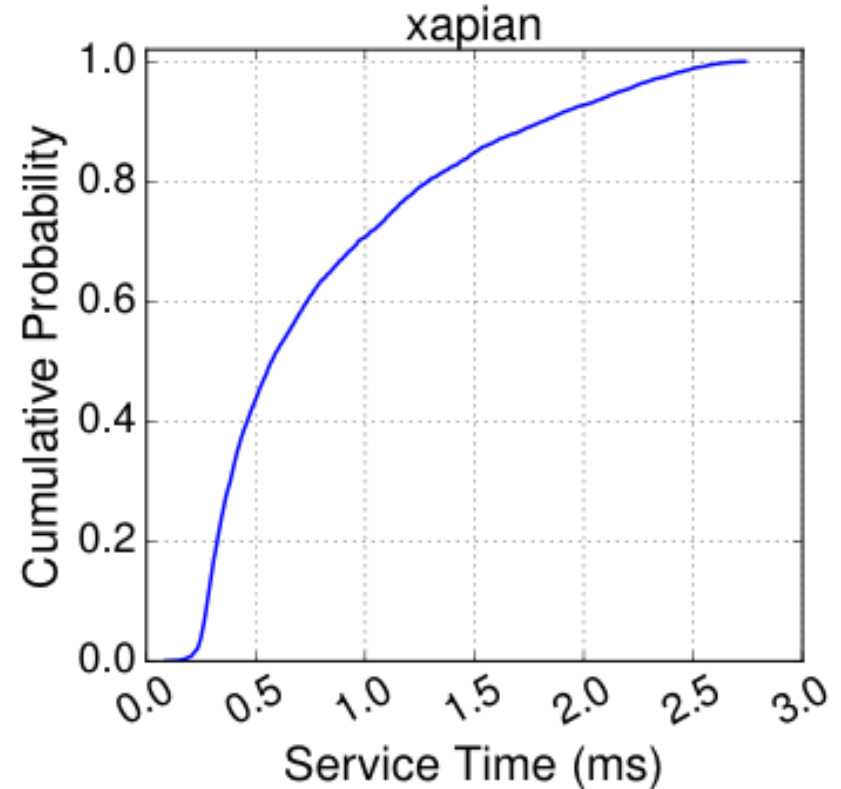
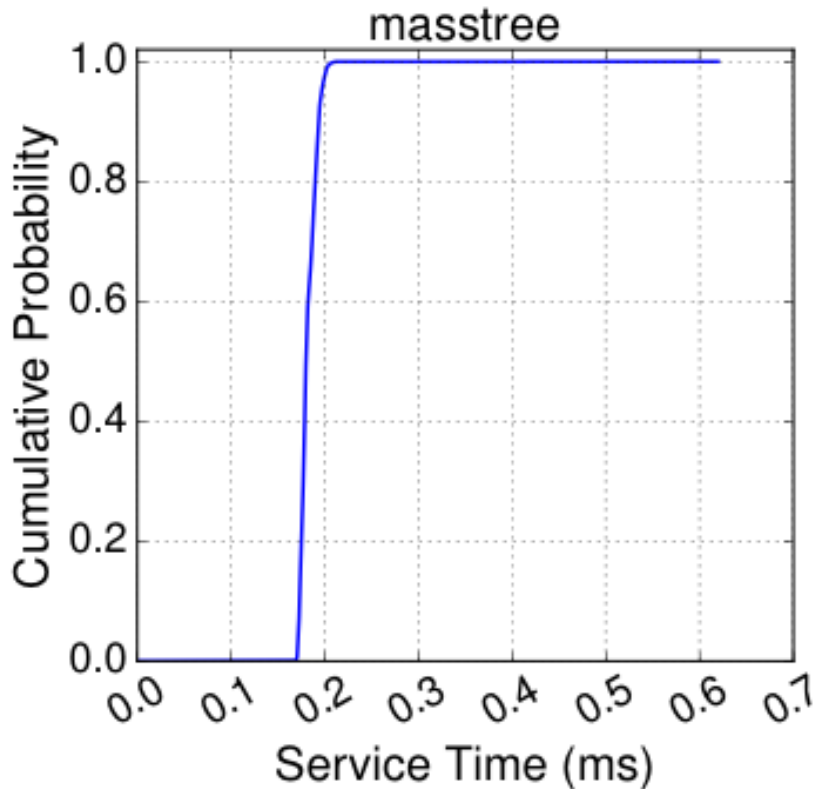
On-disk Database

# Wide Range of End-to-End Latencies



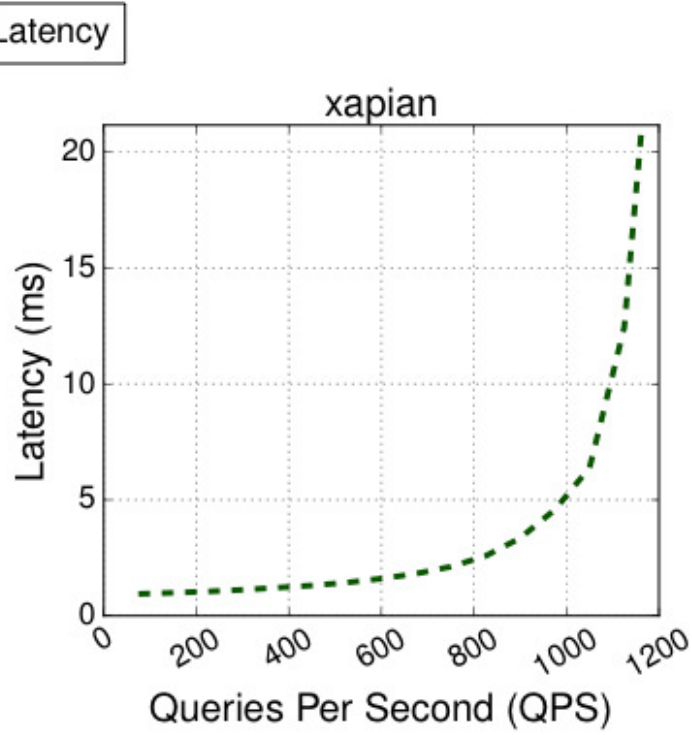
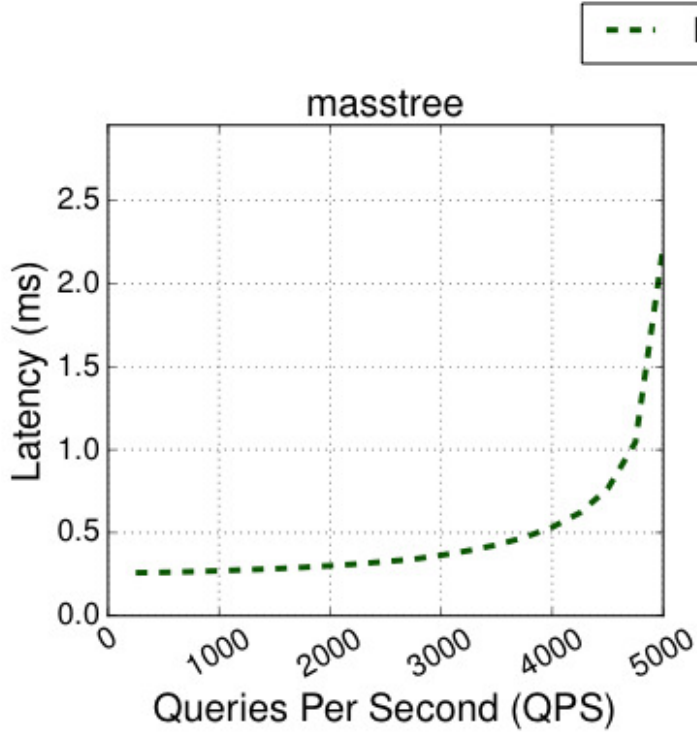
# Varied Service Time Characteristics

13

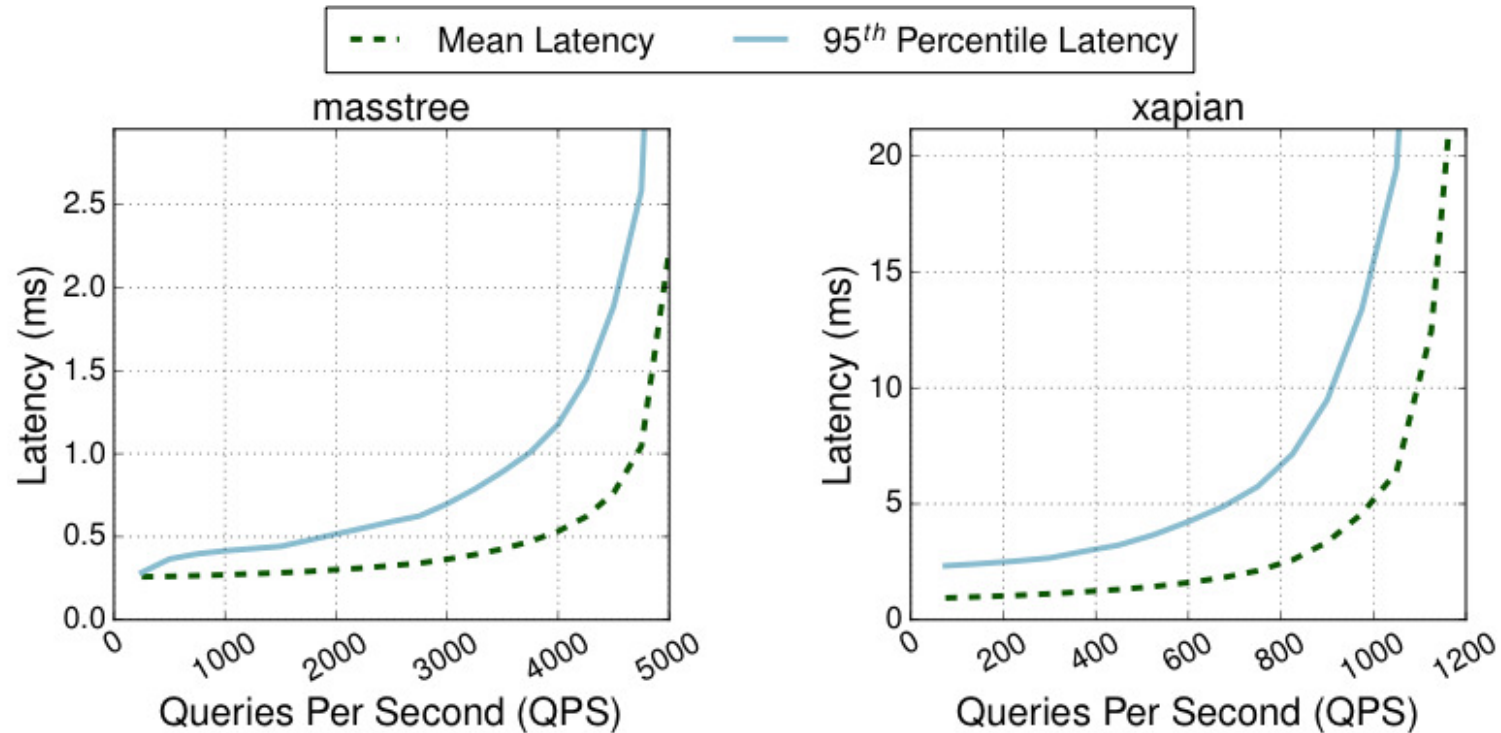


- masstree service times are more tightly distributed
- xapian service times are more loosely distributed

# End-to-End Latency vs. Load

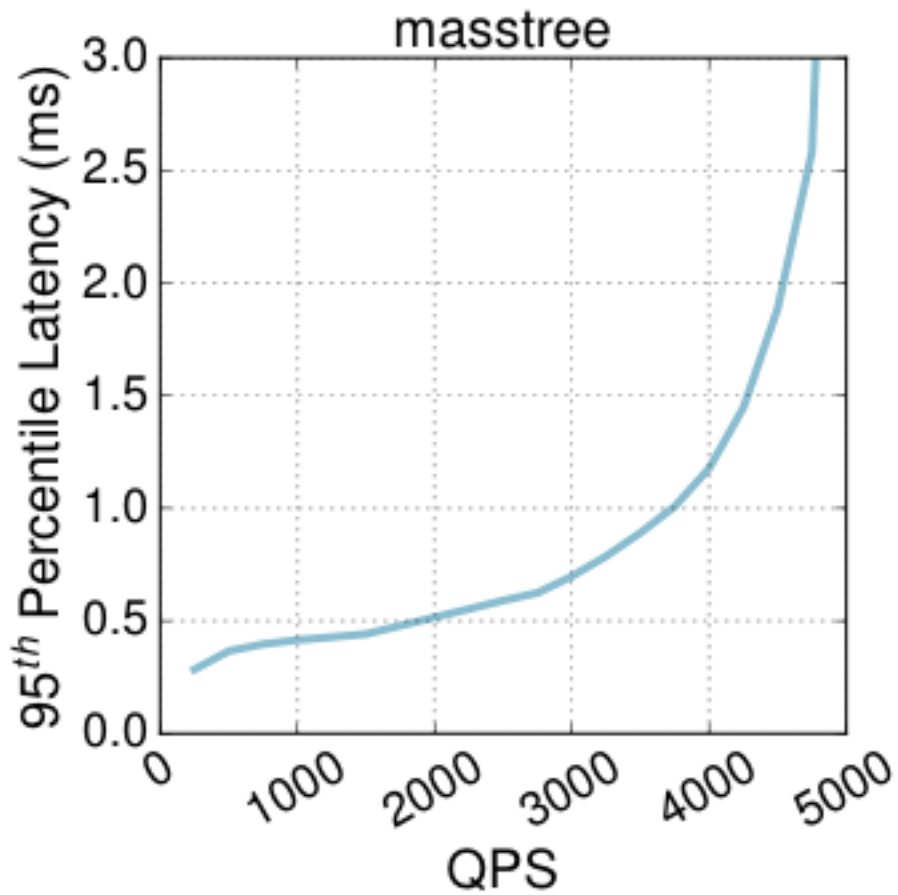


# Tail $\neq$ Mean



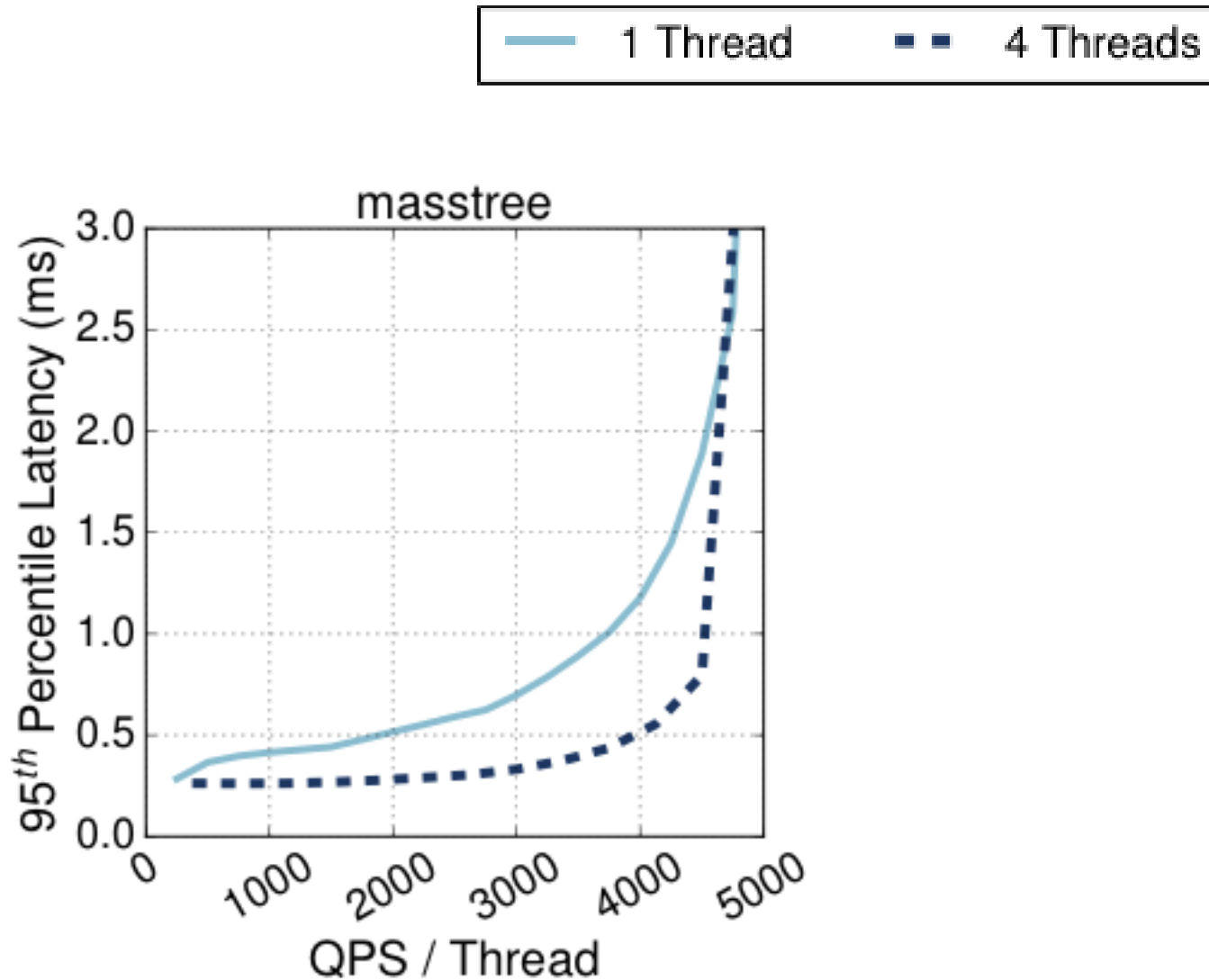
- Tail latency increases more rapidly with load than mean latency
- Relationship between mean and tail latencies is hard to predict

# Impact of Parallelism

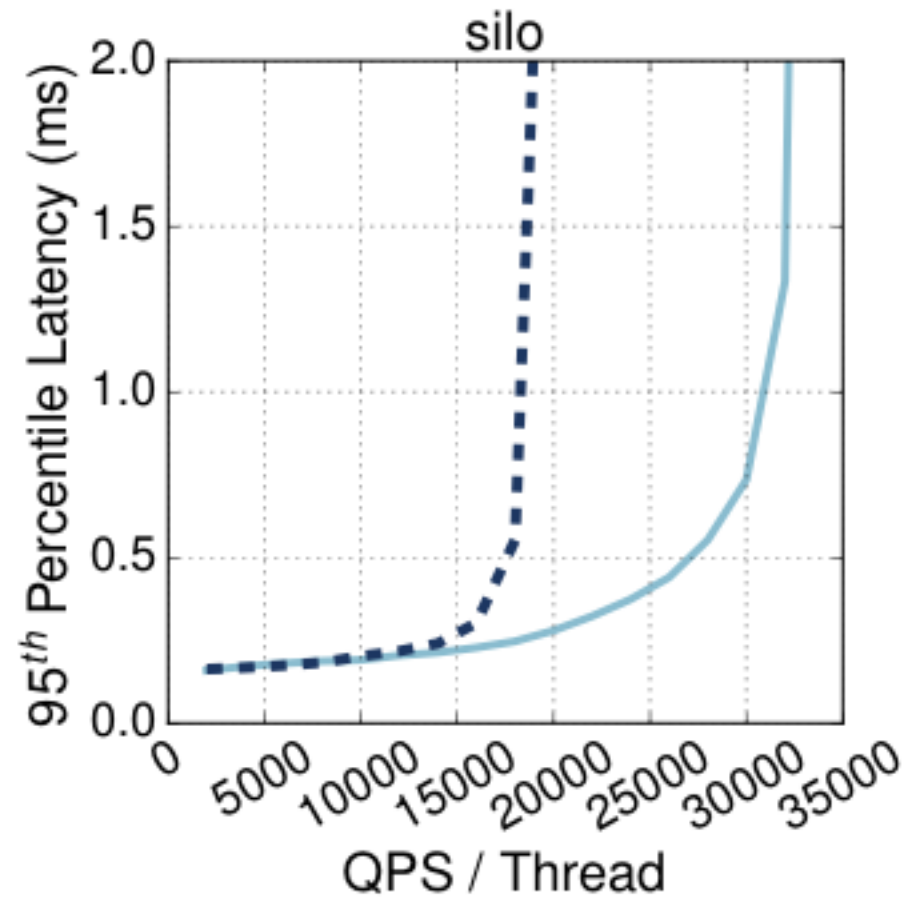
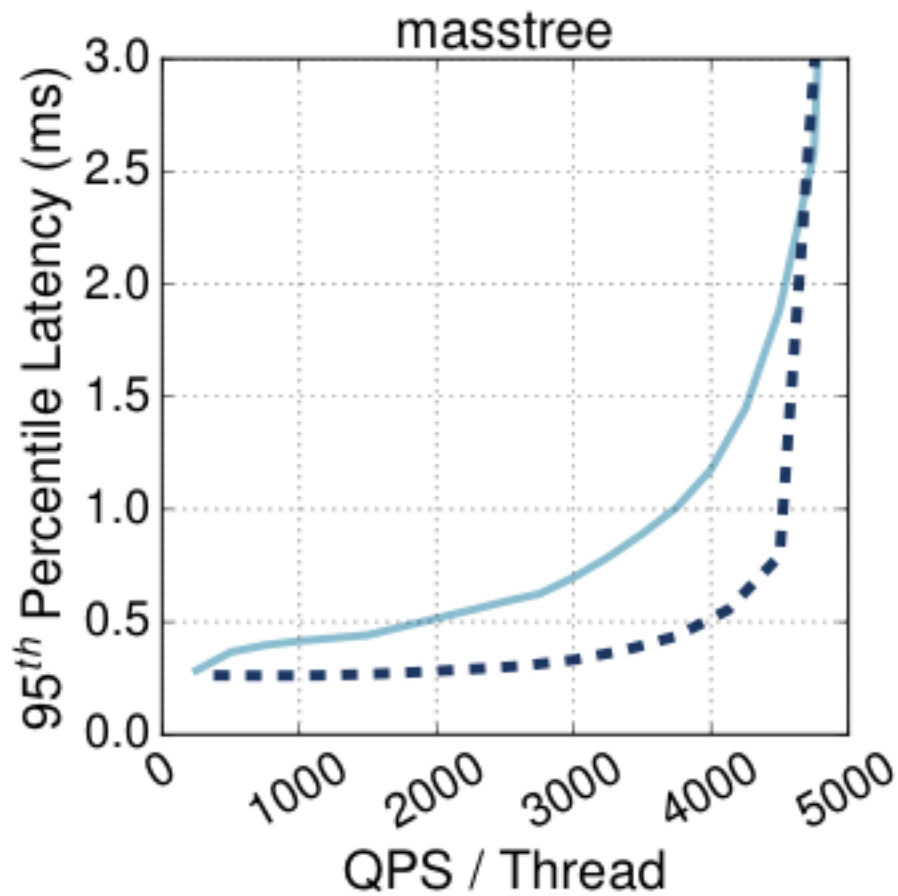




# Parallelism Helps Some Applications



# ...But Hurts Others



- Background and Motivation

- TailBench Applications

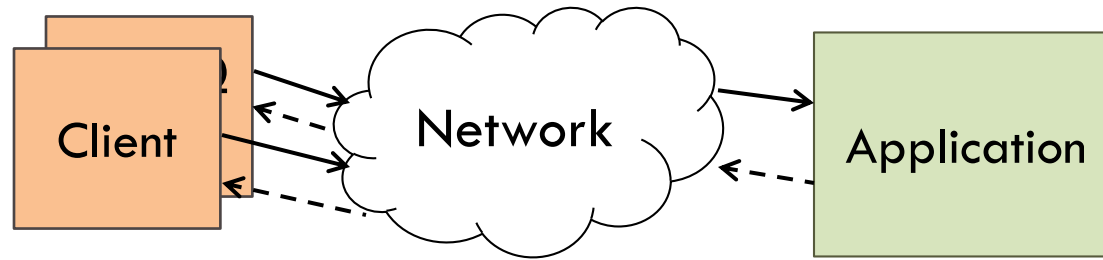
- TailBench Harness

- Simplified Configurations

- Measuring tail latency accurately is complicated
  - ▣ Load generation, statistics aggregation, warmup periods...
  
- Harness encapsulates most of the complexity
  
- Harness makes TailBench easily extensible
  - ▣ New benchmarks reuse existing harness functionality
  
- Simplified harness configurations enable different measurement scenarios
  - ▣ Trade off some accuracy for reduced setup complexity

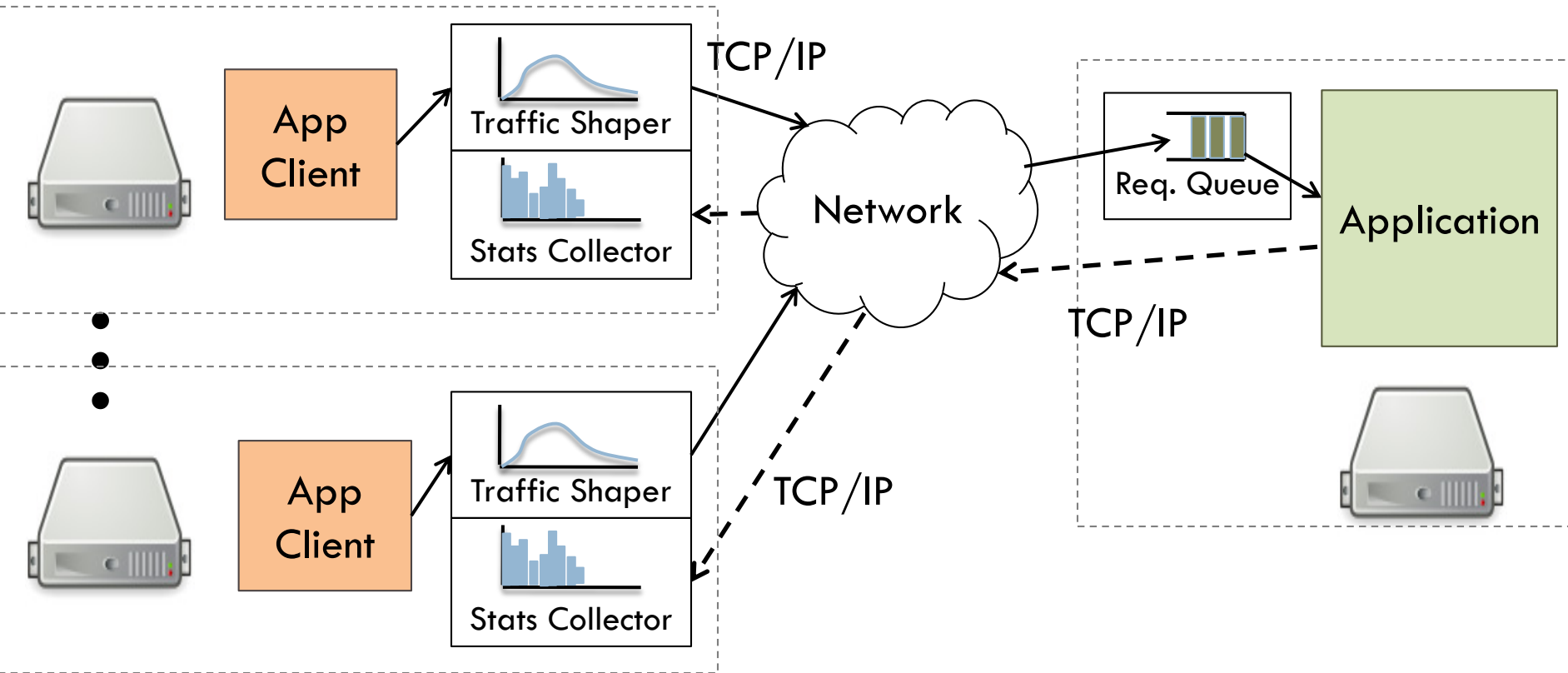
# Example: Open- vs. Closed-Loop Clients

21

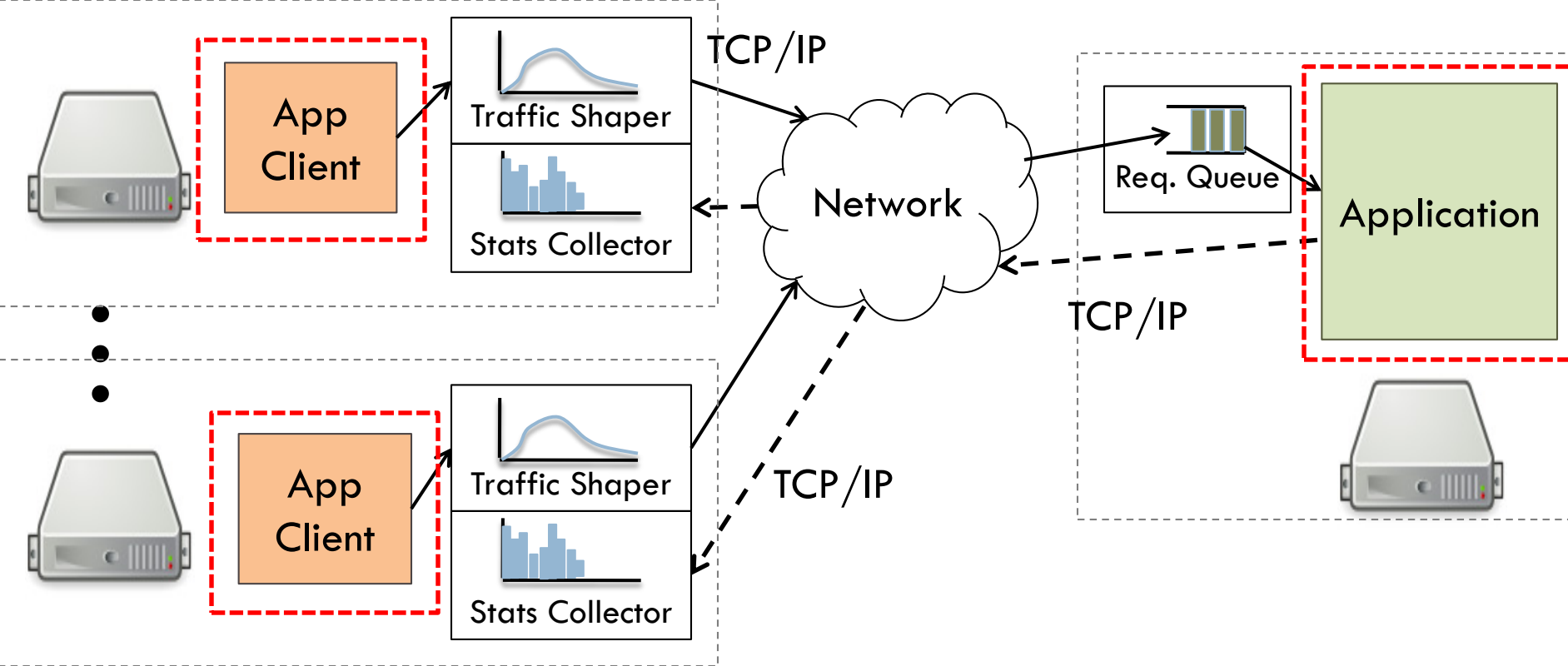


- Many popular load testers use closed-loop clients
  - ▣ Clients wait for response before submitting next request
  - ▣ Increase in application load throttles client request rate
- Latency-critical applications typically service a large number of independent clients
  - ▣ Request rate independent of application load
  - ▣ Better modeled by open-loop clients
- Closed-loop clients can underestimate latency by orders of magnitude [Tene LLS 2013, Zhang ISCA 2016]

# Networked Harness Configuration



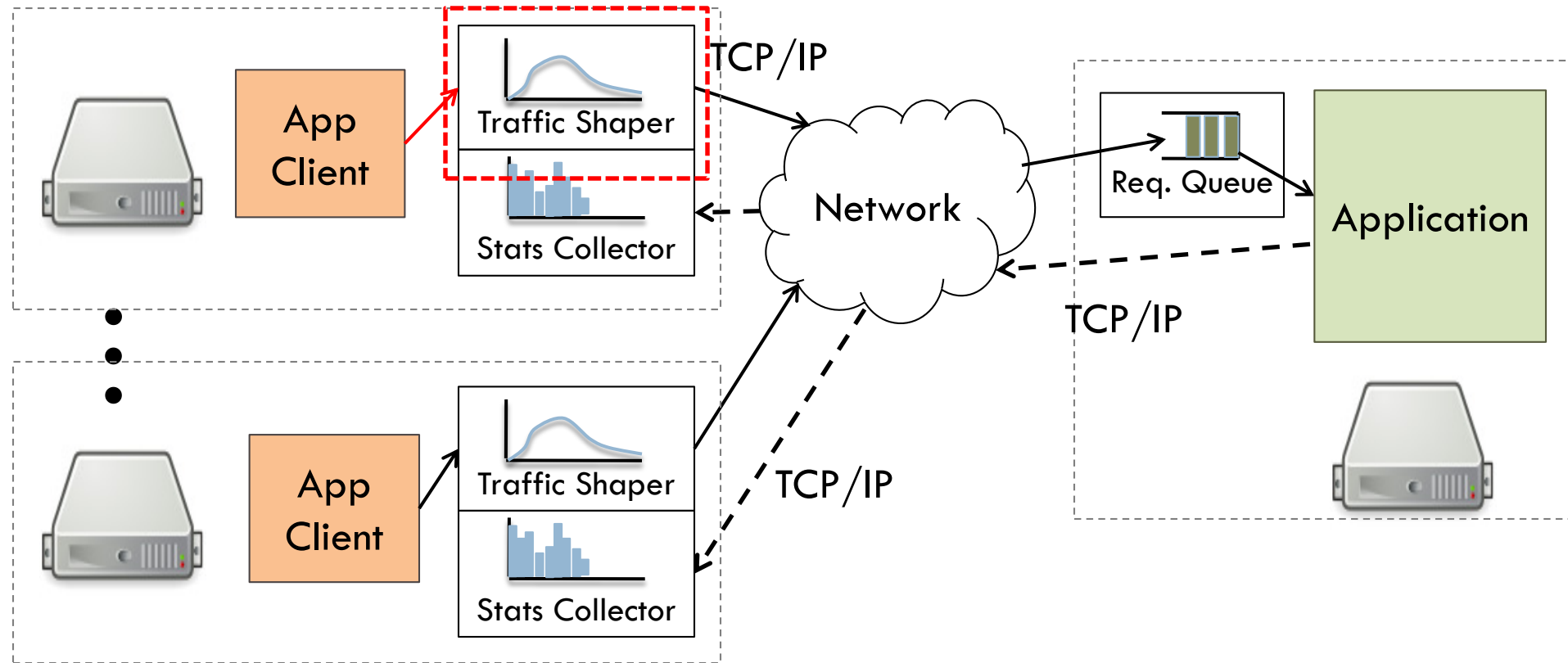
# Networked Harness Configuration



- Application and the clients run on separate machines

# Networked Harness Configuration

24

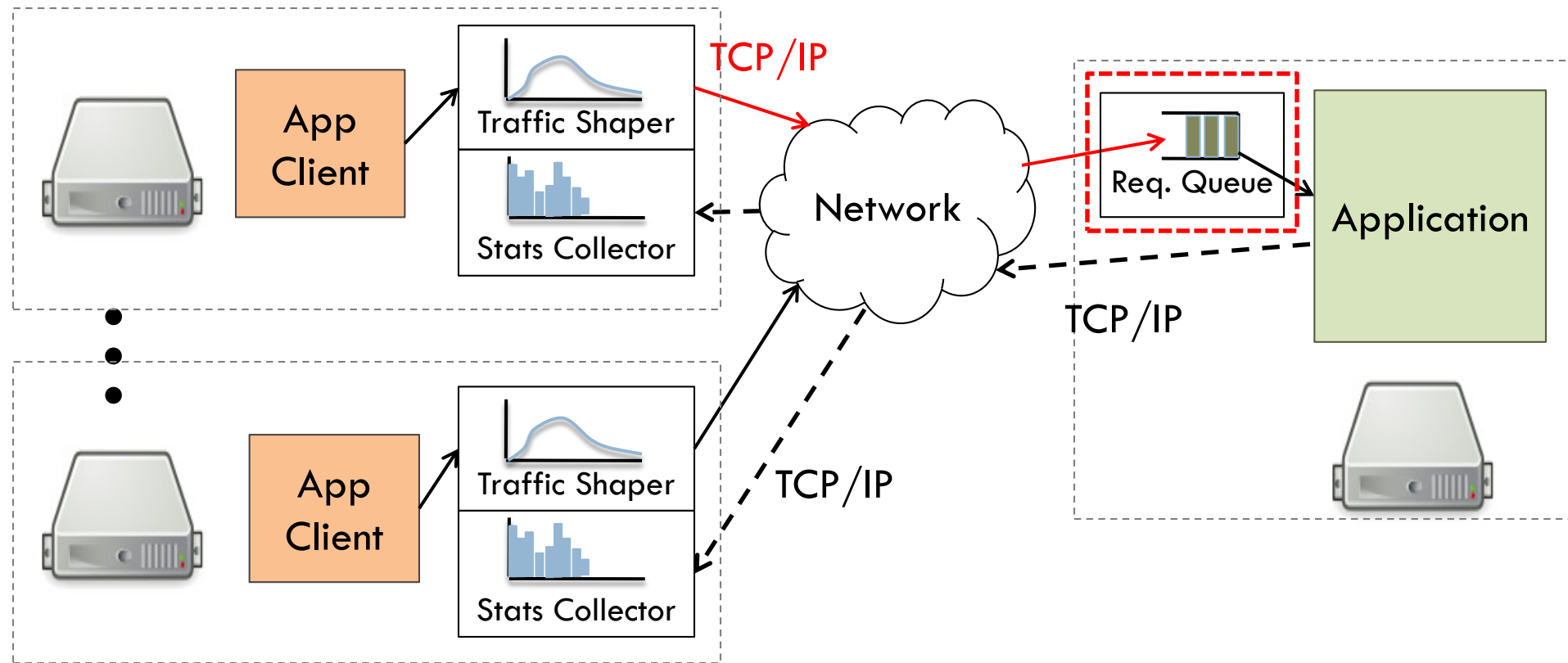


- Application and the clients run on separate machines
- Traffic Shaper inserts inter-request delays to model load



# Networked Harness Configuration

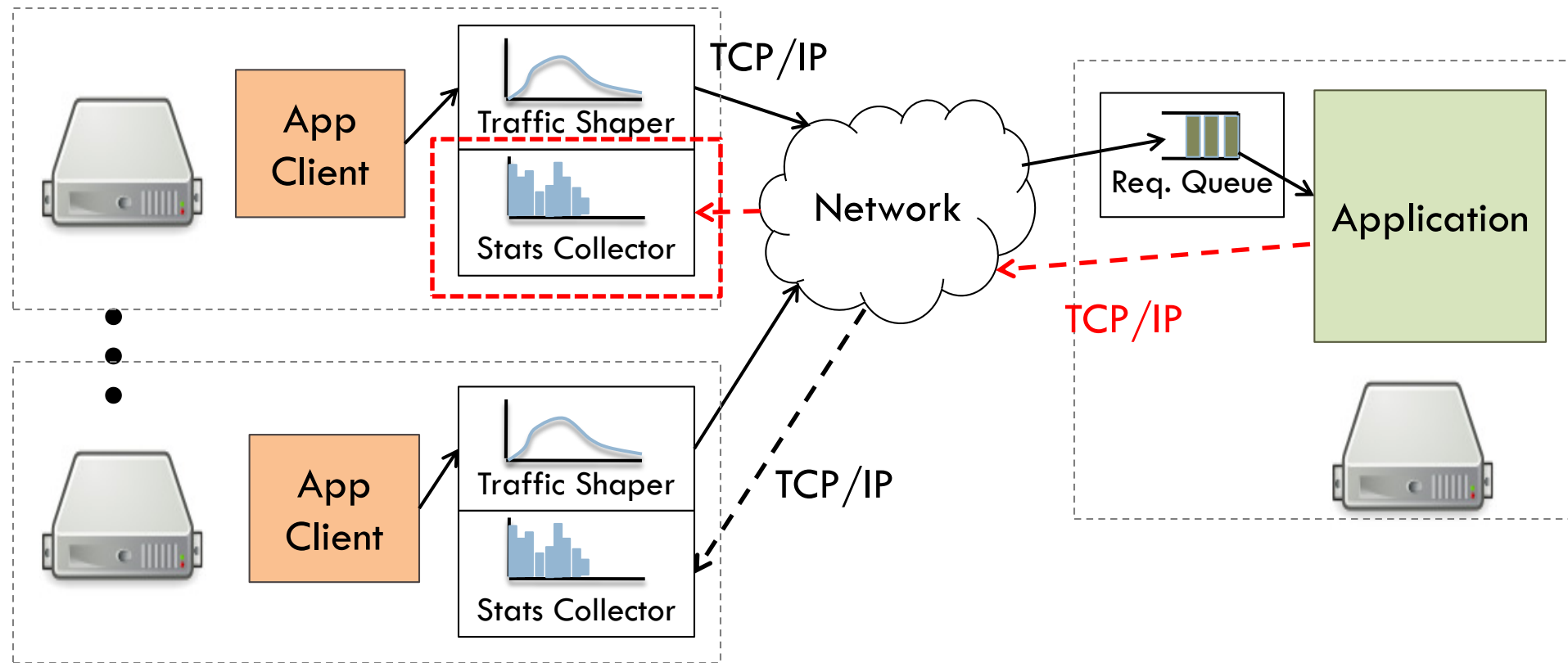
25



- Application and the clients run on separate machines
- Traffic Shaper inserts inter-request delays to model load
- Request Queue enqueues incoming requests and measures service times and queuing delays

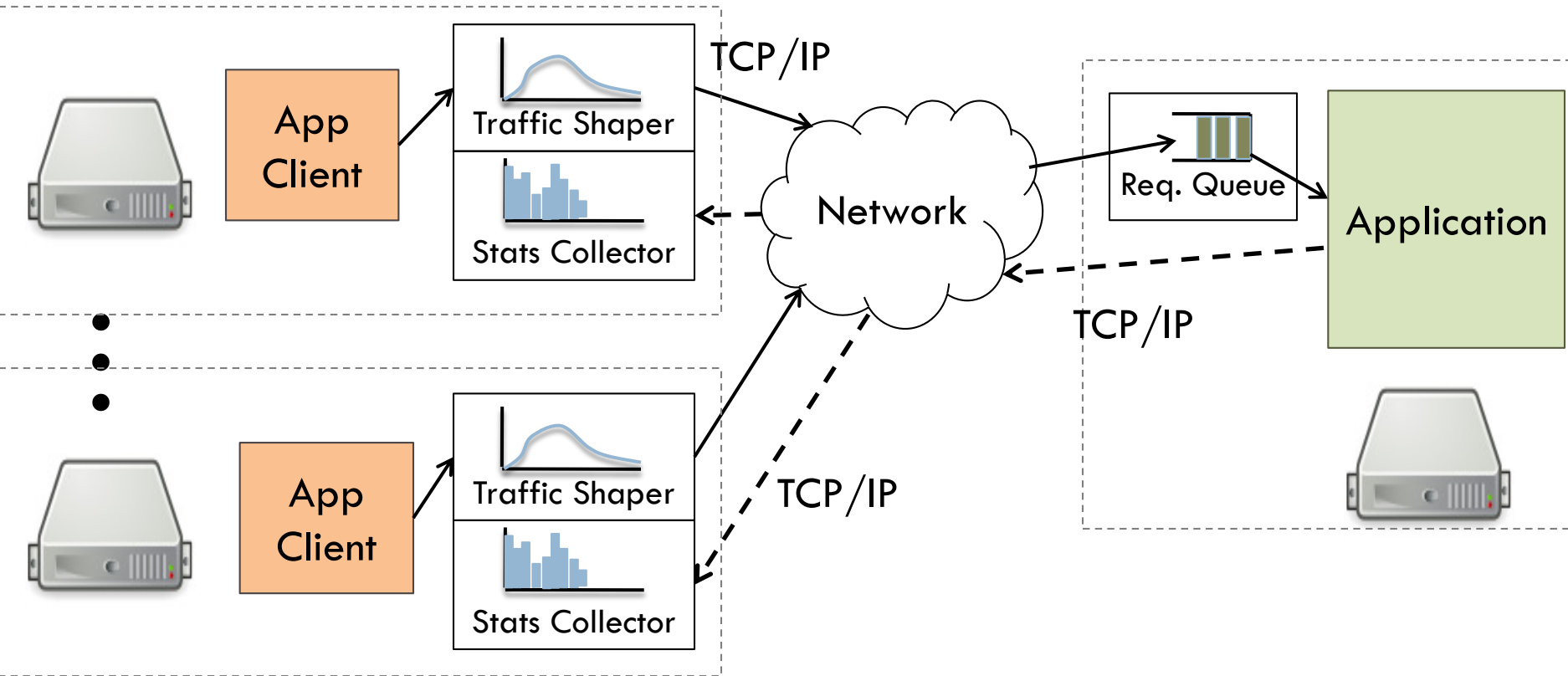
# Networked Harness Configuration

26



- Application and the clients run on separate machines
- Traffic Shaper inserts inter-request delays to model load
- Request Queue enqueues incoming requests and measures service times and queuing delays
- Statistics Collector aggregates latency data

# Networked Harness Configuration

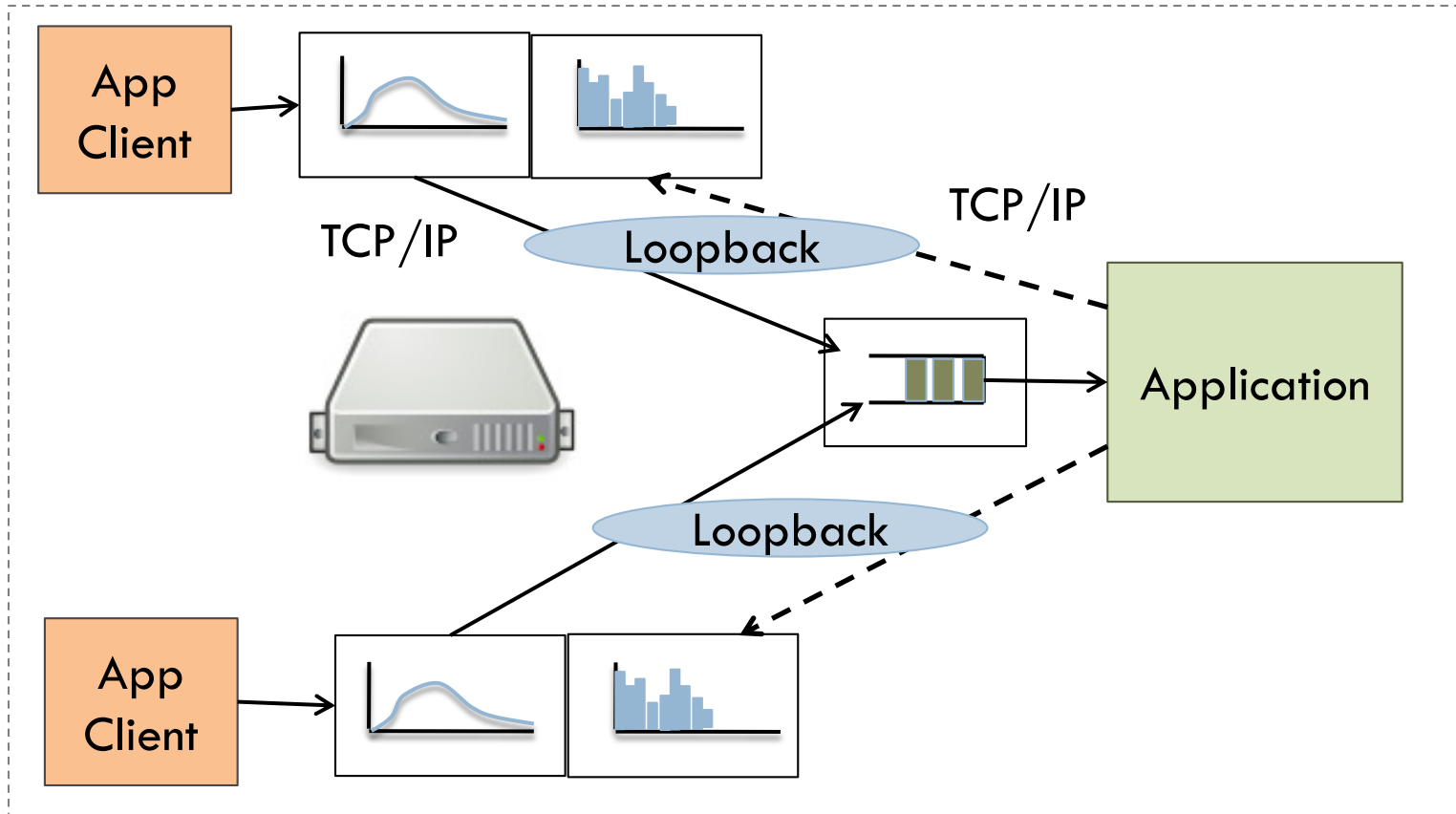


✓ Faithfully captures all sources of overhead

x Difficult to configure and deploy

- Background and Motivation
- TailBench Applications
- TailBench Harness
- Simplified Configurations

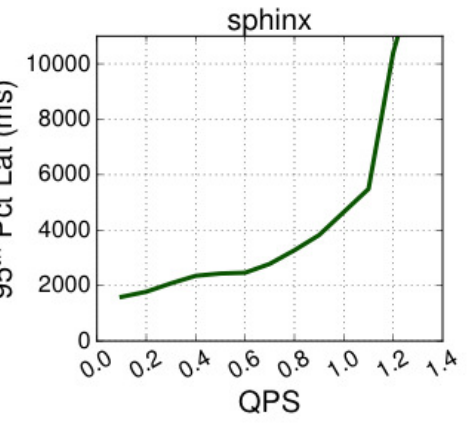
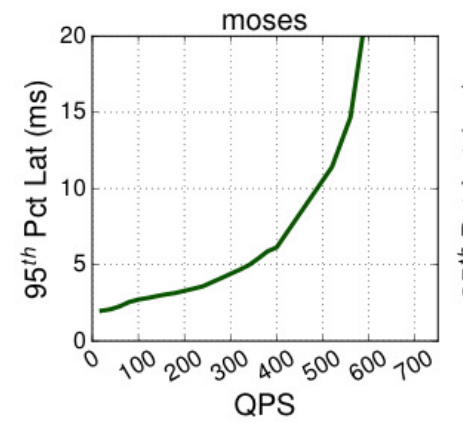
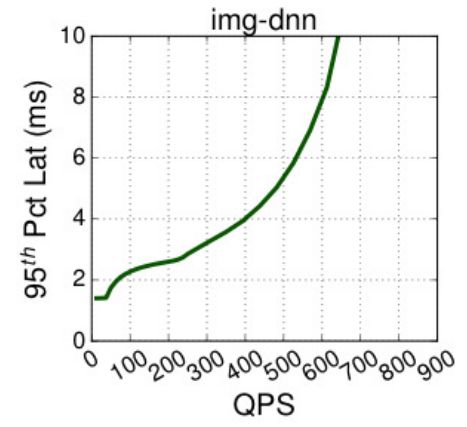
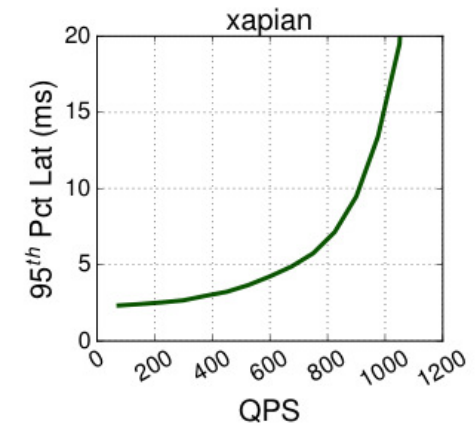
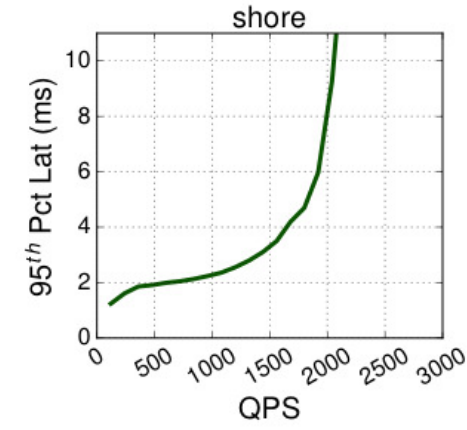
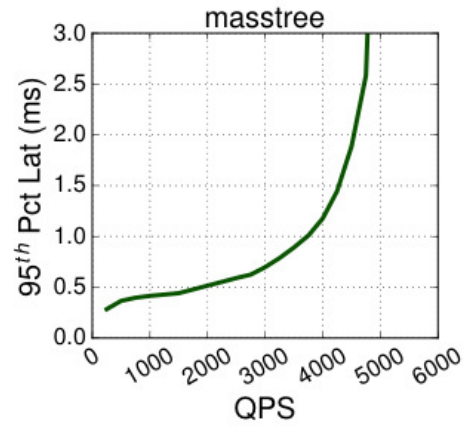
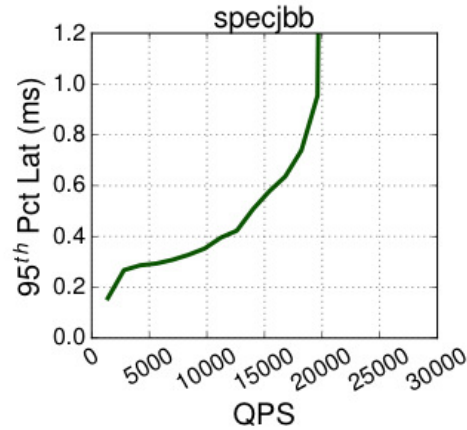
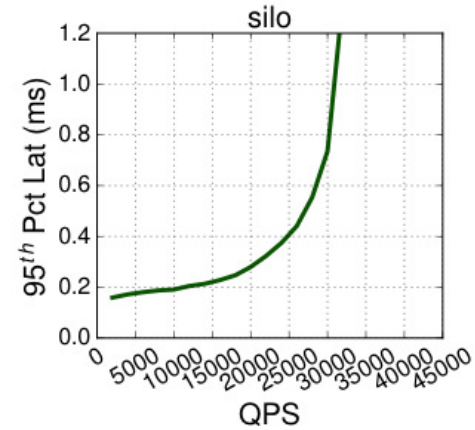
# Loopback Harness Configuration



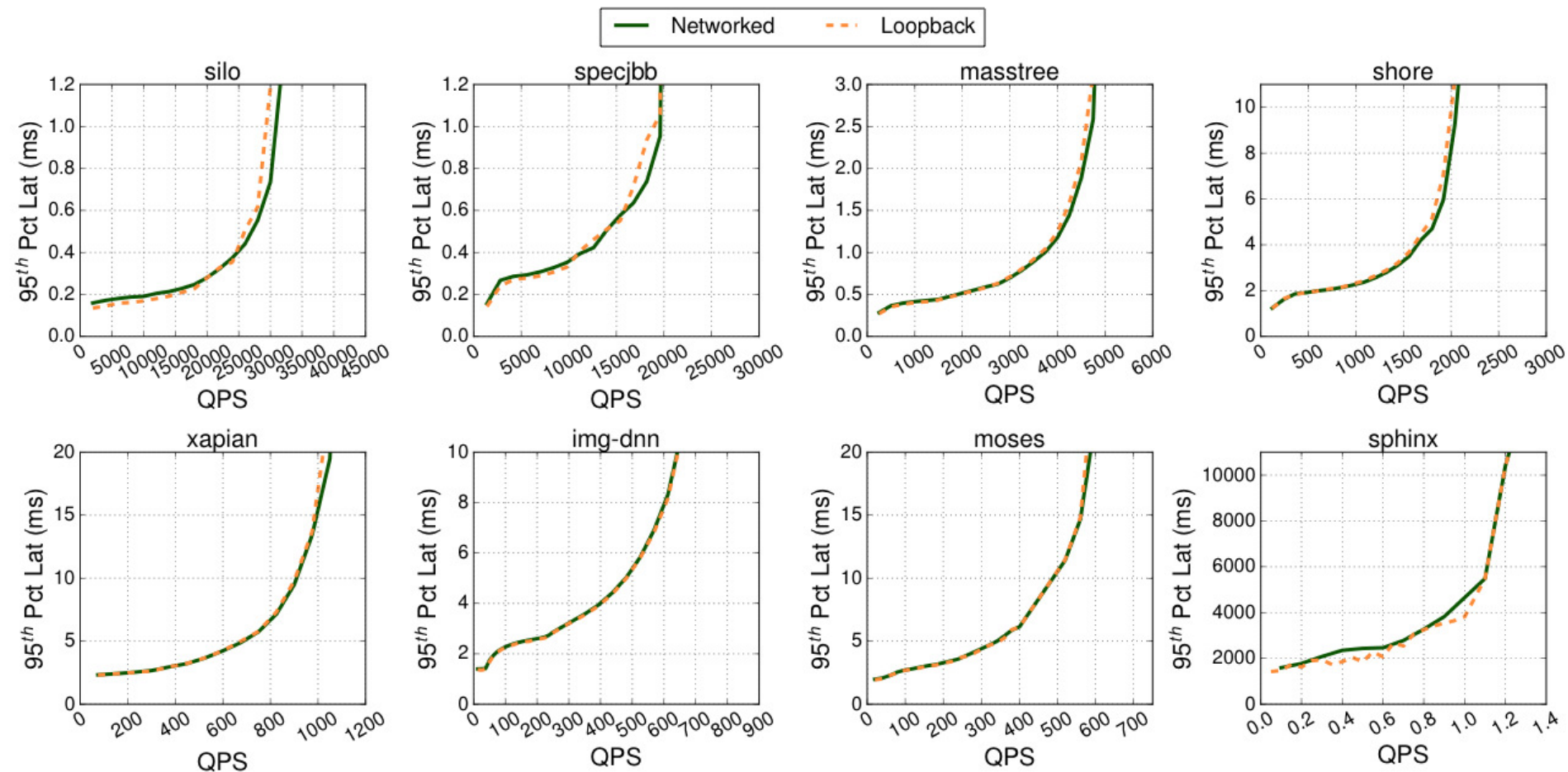
- ❑ Application and clients reside on the same machine
- ✓ Reduced setup complexity
- ✓ Highly accurate in many cases

# Load-Latency for Networked Configuration 30

— Networked



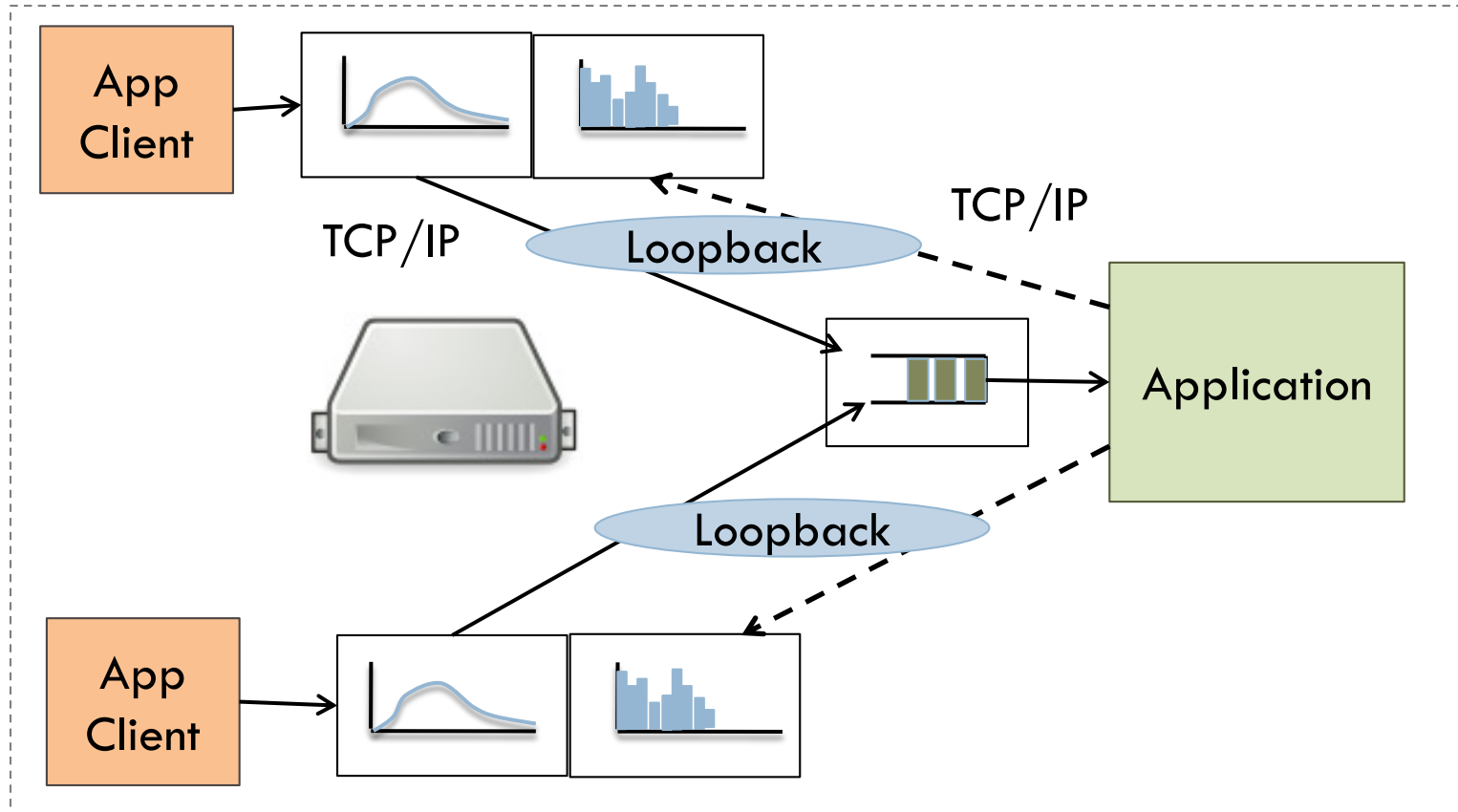
# Loopback Configuration Highly Accurate



- Loopback and Networked configurations have near-identical performance
  - Networking delays minimal in our setup

# Loopback Harness Configuration

32

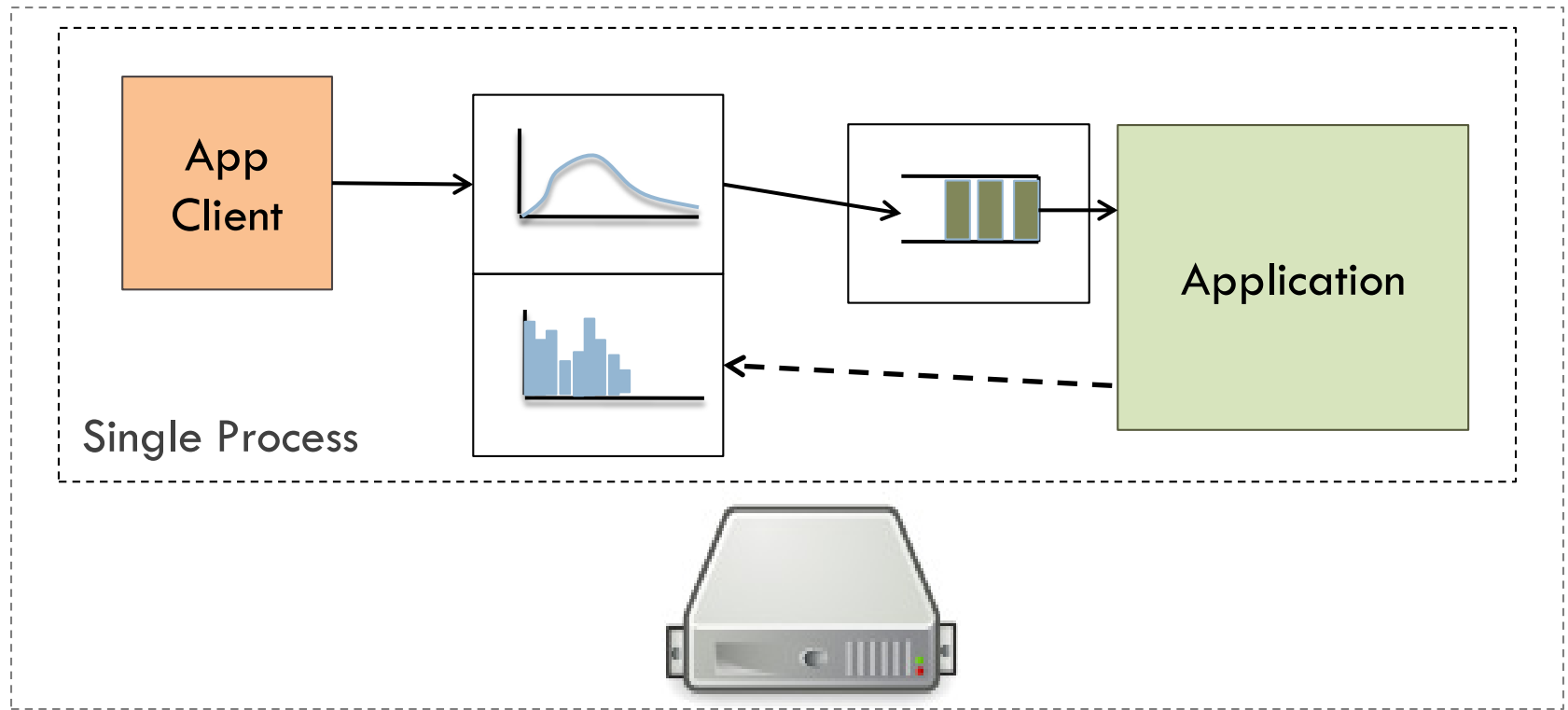


- Application and clients reside on the same machine
- ✓ Reduced setup complexity
- ✓ Highly accurate in many cases
- ✗ Still difficult to simulate



# Integrated Harness Configuration

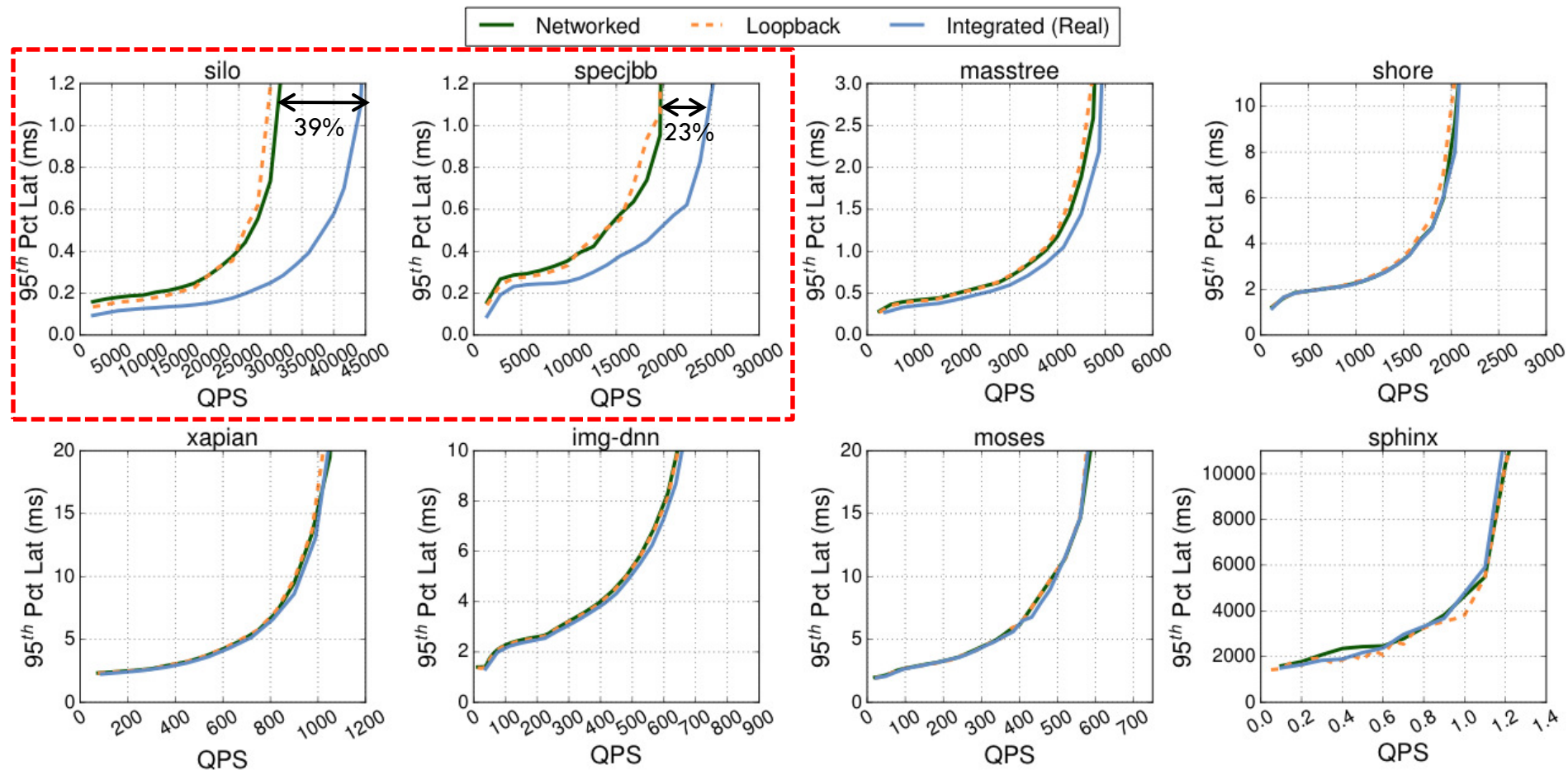
33



- Application and client integrated into a single process
- ✓ Easy to setup
- ✗ Some loss of accuracy

# Integrated Configuration Validation

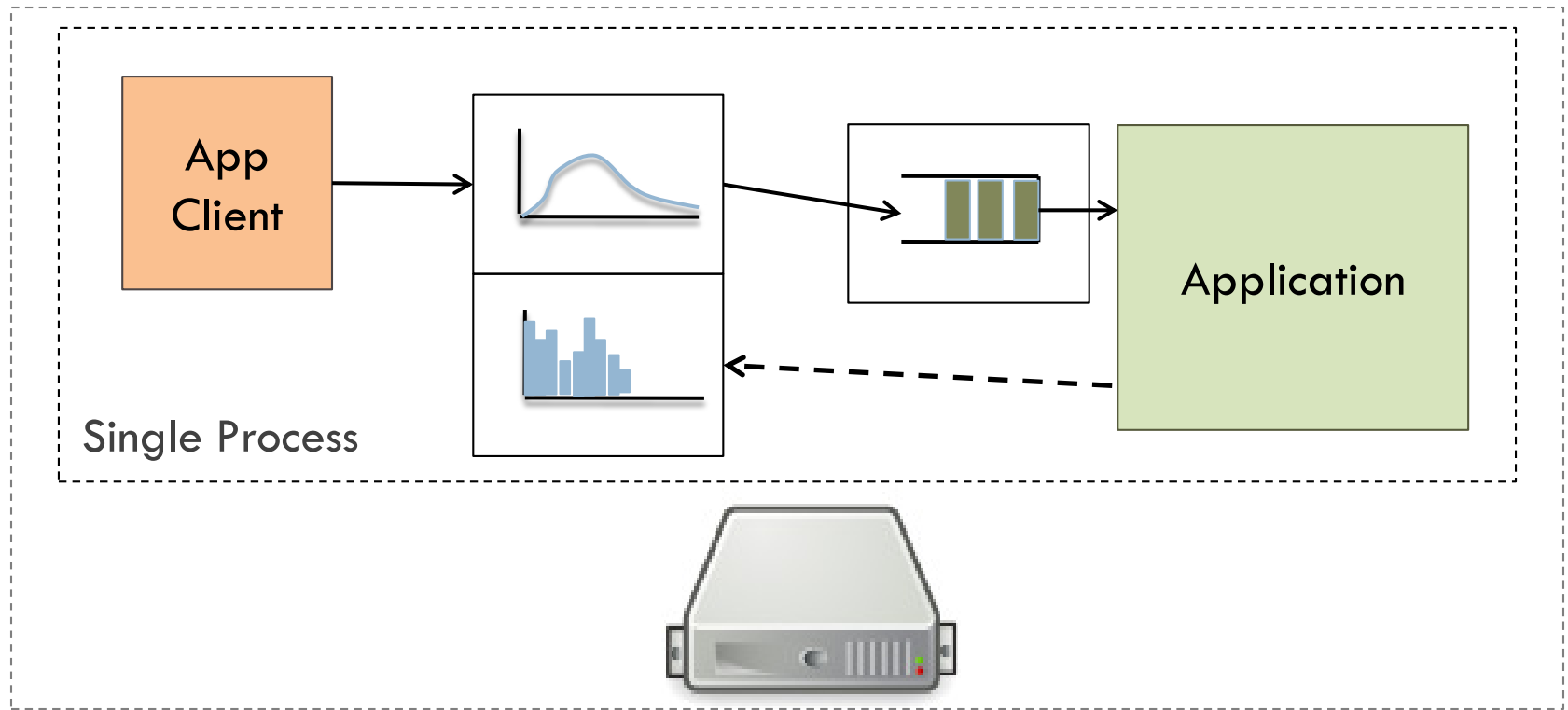
34



- Networked/Loopback configurations saturate earlier for applications with short requests (silo, specjbb)
- TCP/IP processing overhead a significant fraction of request

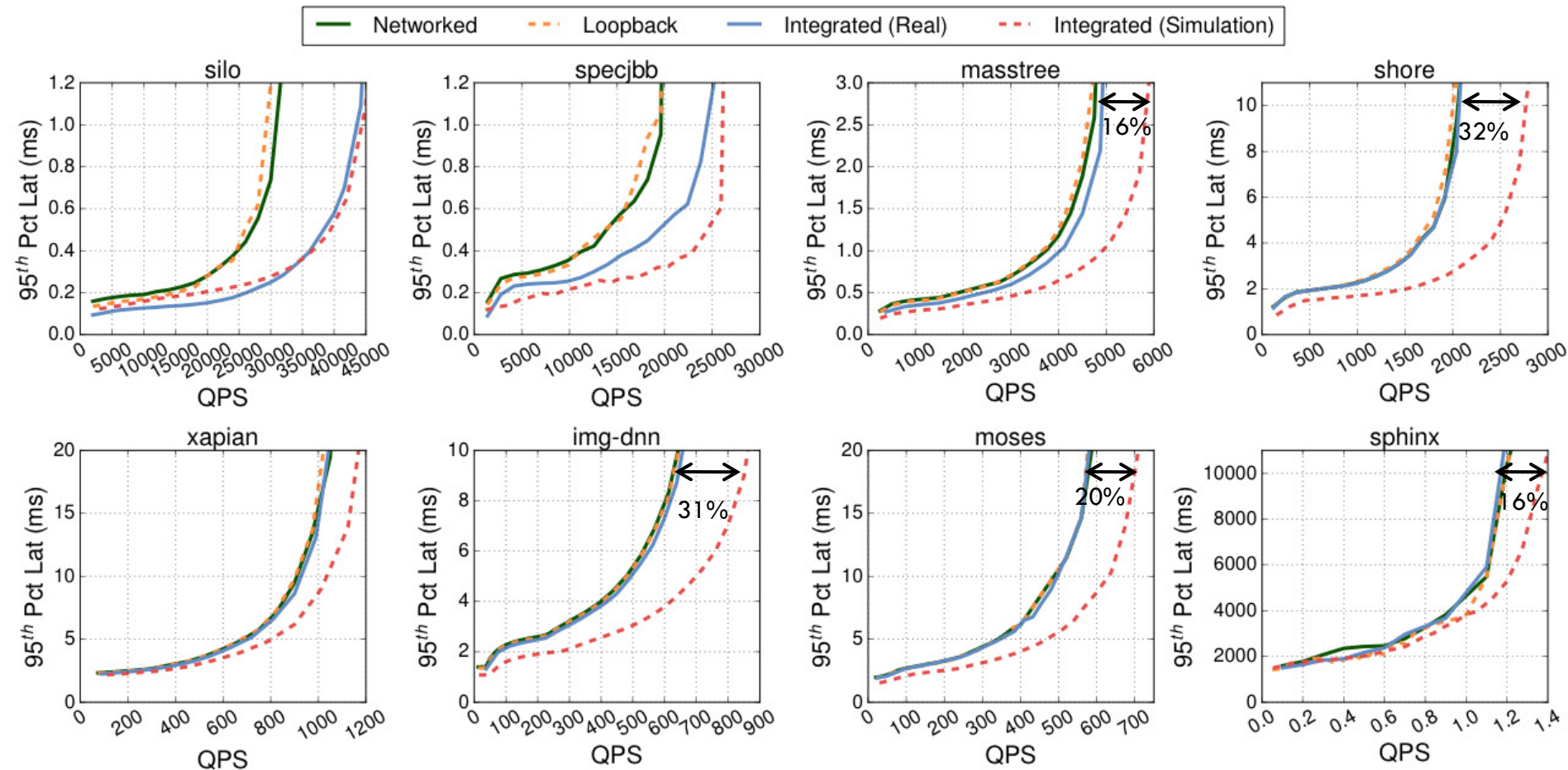
# Integrated Harness Configuration

35



- ❑ Application and client integrated into a single process
- ✓ Easy to setup
- ✗ Some loss of accuracy
- ✓ Enables user-level simulations

# Simulation vs. Real System



Performance difference between real and simulated systems well within usual simulation error bounds

- Average absolute error in saturation QPS: 14%

- zsim IPC error for SPEC CPU2006 applications: 8.5 – 21%

- TailBench includes a diverse set of latency-critical applications with varied latency characteristics
- TailBench harness implements a statistically sound experimental methodology to achieve accurate results
- Various harness configurations allow trading off configuration complexity for some accuracy
  - ▣ Our results show that the integrated configuration is highly accurate for six of our eight benchmarks

**THANKS FOR YOUR ATTENTION!**

**QUESTIONS?**

[tailbench.csail.mit.edu](http://tailbench.csail.mit.edu)



**Massachusetts  
Institute of  
Technology**

