# SSD Performance Profiling
ECSE 4320
Ben Herman

# Content

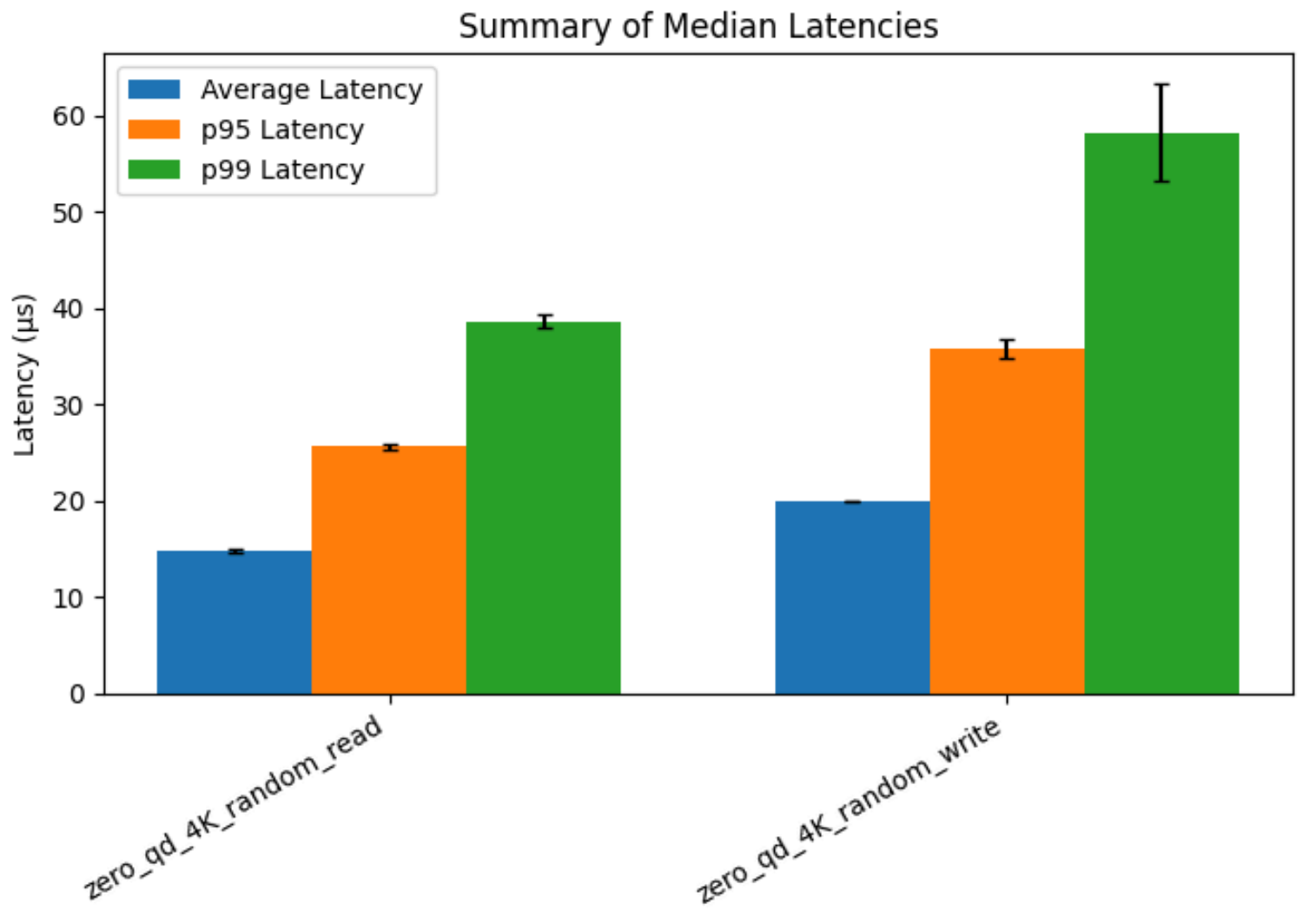# Experiment Setup

**Timing Measurement:**

- Execution time is measured using `mach_absolute_time()`.

**Conditions:**

- Model: M2 Mac
- OS: Sequoia 15.6
- Powersource: Wall outlet
- Ram: 16 GB

# Zero-Queue Baselines

Zero-queue latency for 4 KiB random read and write.



## Summary of Median Latencies

Legend:
- Average Latency
- p95 Latency
- p99 Latency

Y-axis: Latency (µs)

X-axis categories:
- zero_qd_4K_random_read
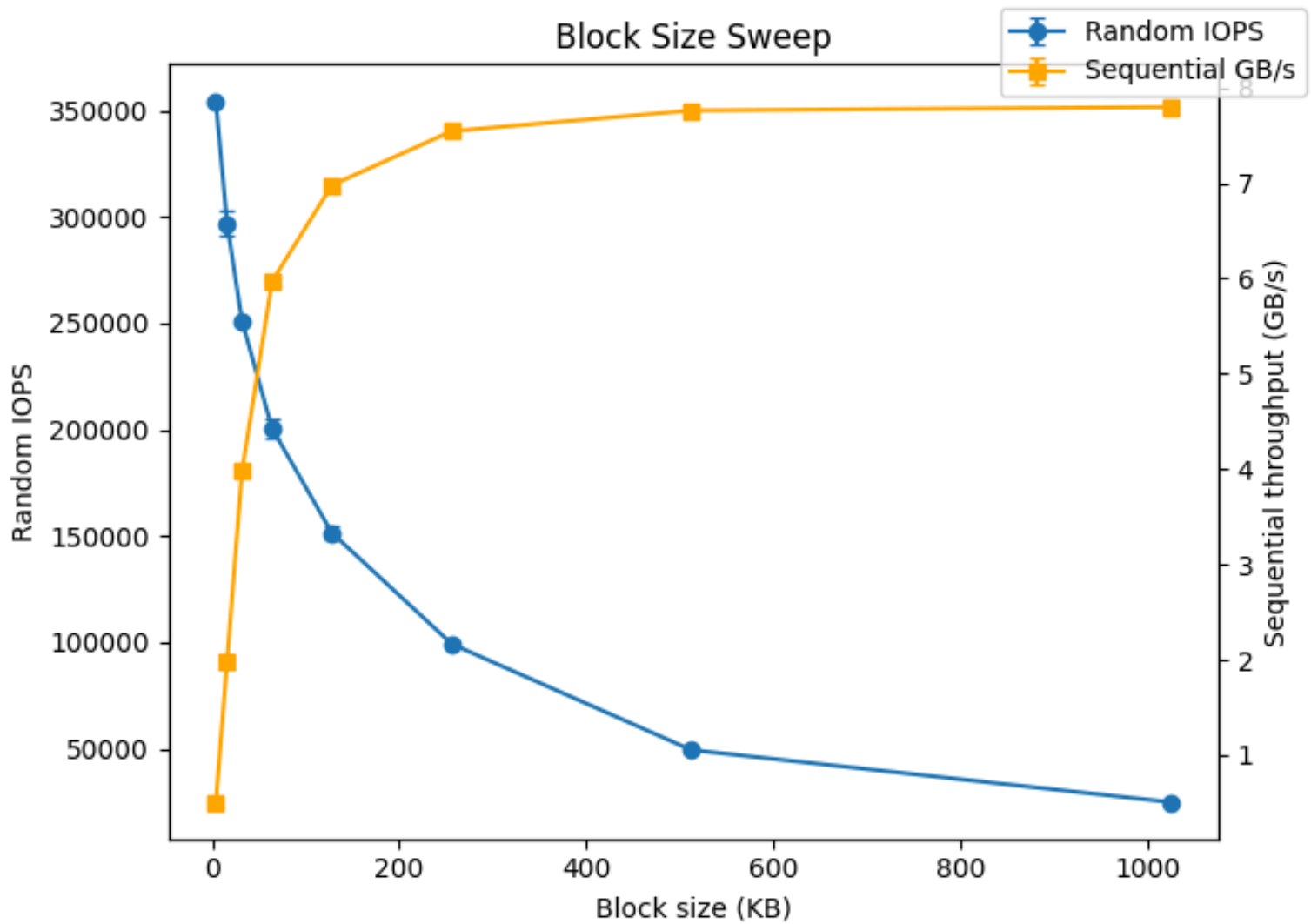- zero_qd_4K_random_write

# Block-Size Sweep

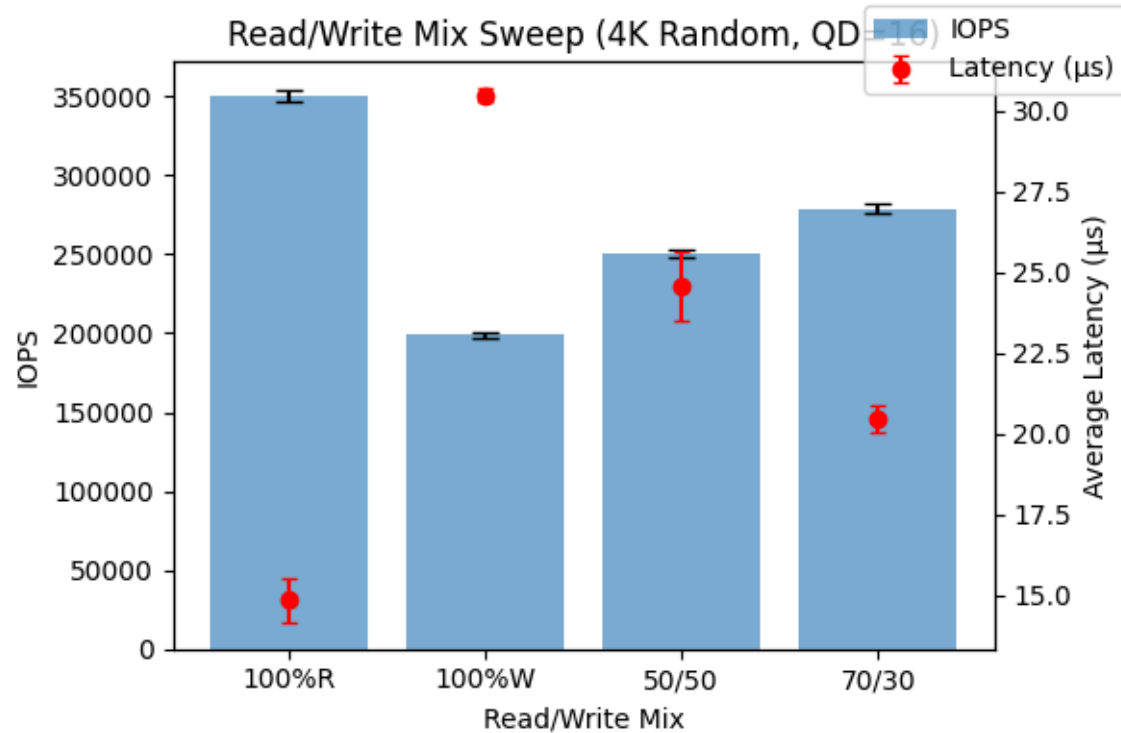Impact of block size on random IOPS and sequential throughput.

We see in the graph:

- Small blocks less than 192 KB are throughput limited by IOPS
- Large blocks more than 192 KB are throughput limited by the PCIe because its saturated with too many requests.

# Read/Write Mix Sweep

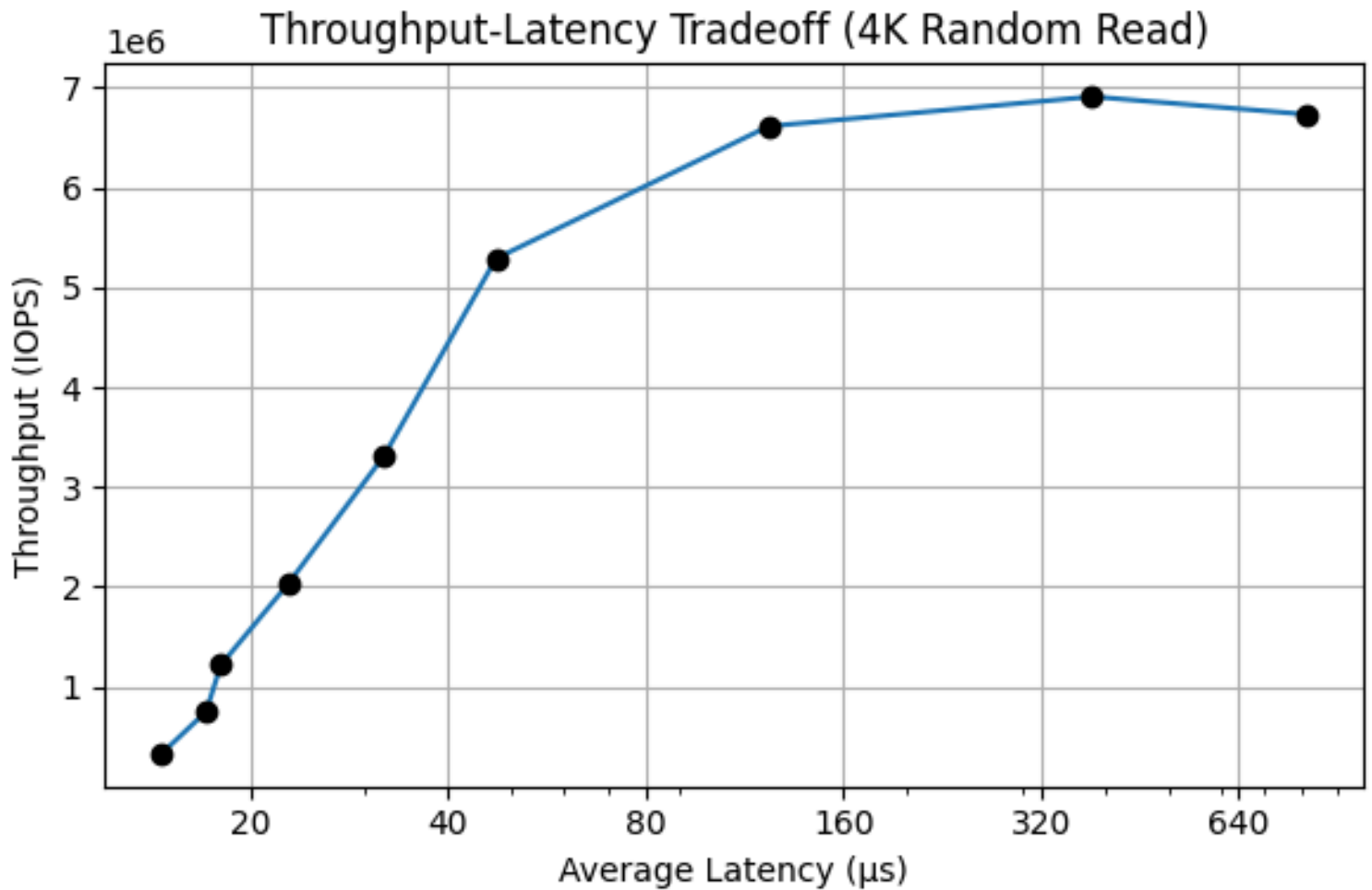Effect of varying read/write ratio at fixed block size (4 KiB random).

- 100% read yields highest IOPS and lowest latency and the opposite for 100 write
- Increasing write fraction increases latency and decreases IOPS
- We see the other ratios also follow these trends



Read/Write Mix Sweep (4K Random, QD=16)

# Queue Depth Sweep

Throughput-latency trade-off curve for 4 KiB random reads.
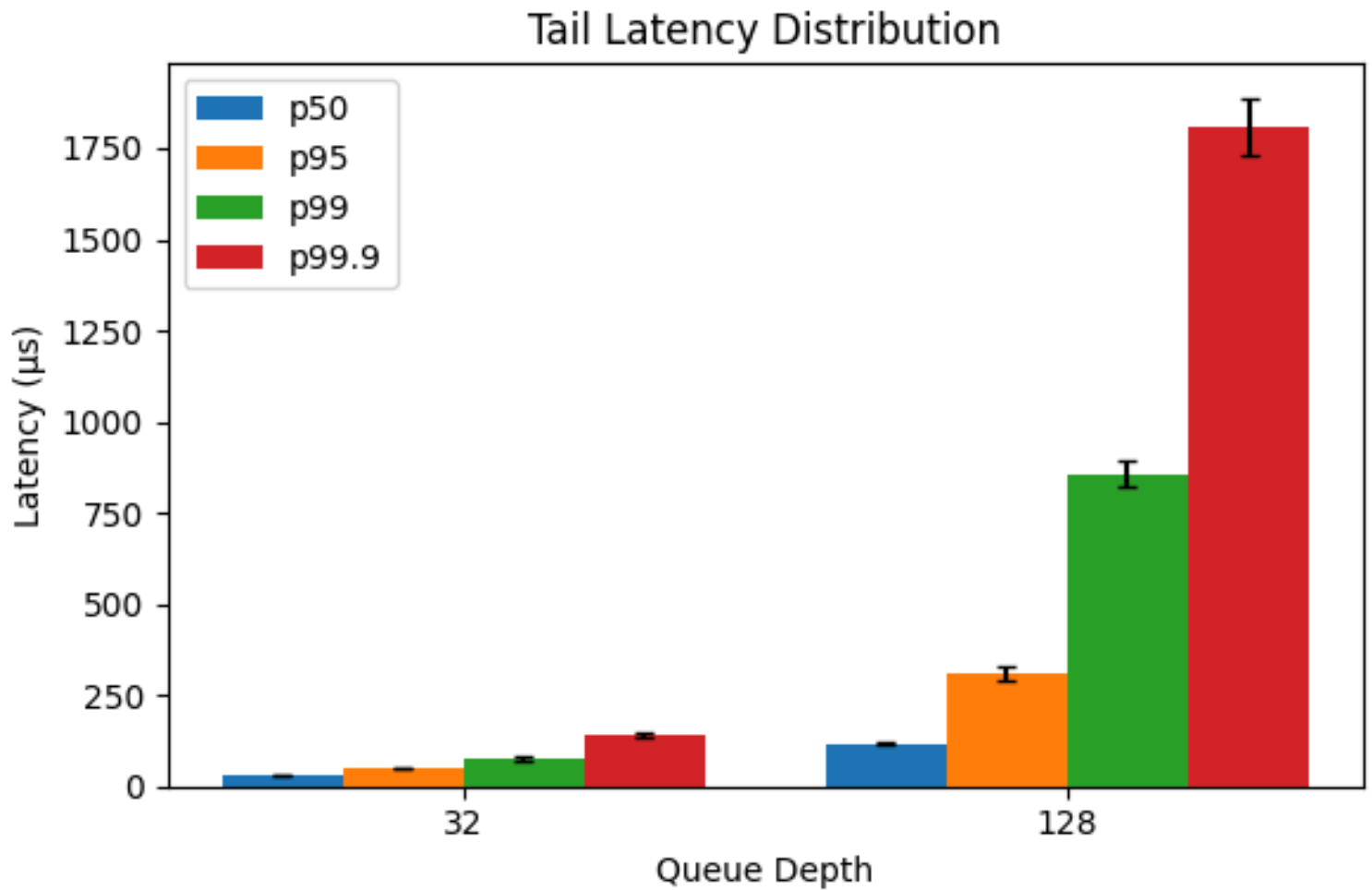
- Throughput rises with QD until saturation ( QD 32-64)
- Latency grows sharply past the knee
- We can see Little's Law holds because throughput and latency are inversely proportional



Throughput-Latency Tradeoff (4K Random Read)

# Tail Latency

Tail latency distribution (p50/p95/p99/p99.9) at different QDs.
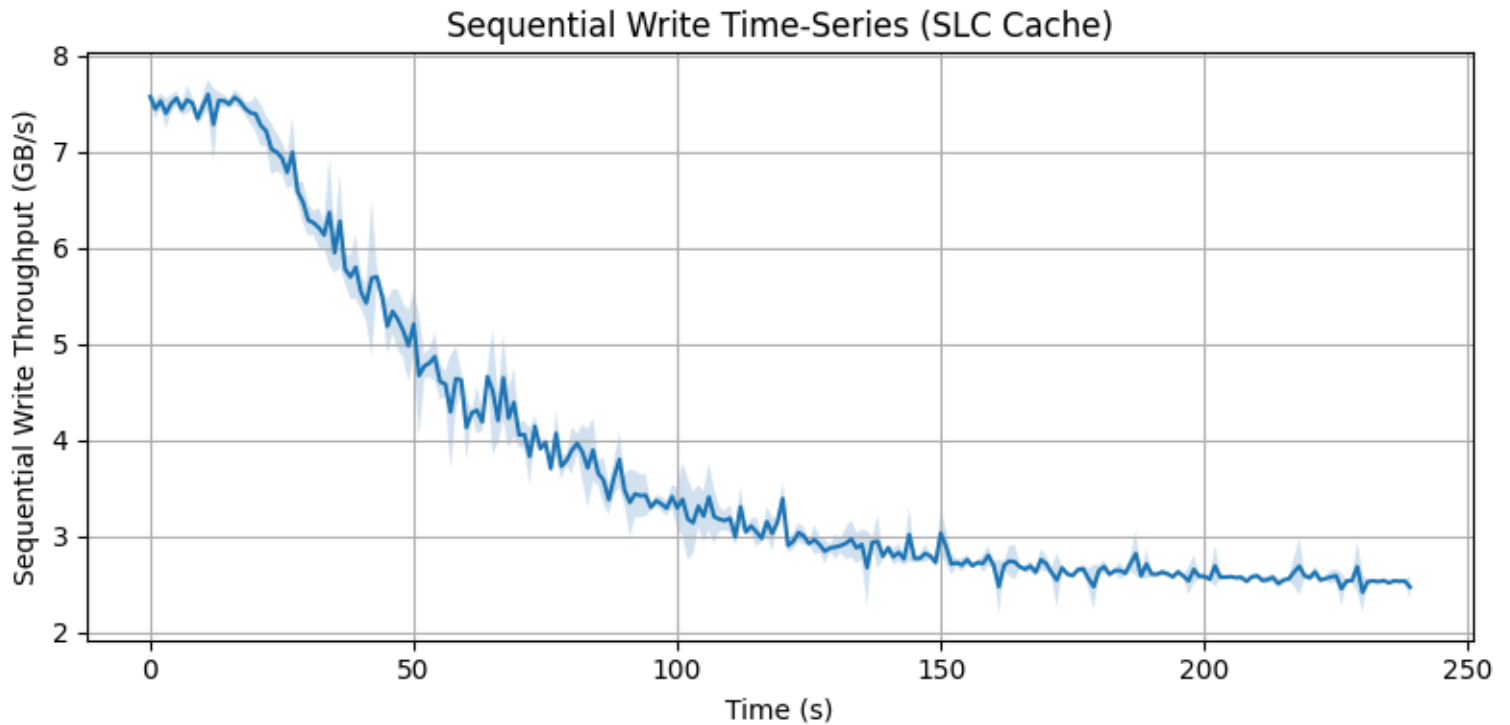
- p99.9 latency spikes significantly at high queue depth
- Important for SLA-sensitive workloads
- Highlights worst-case latency scenarios beyond average

# Sequential Write Time-Series

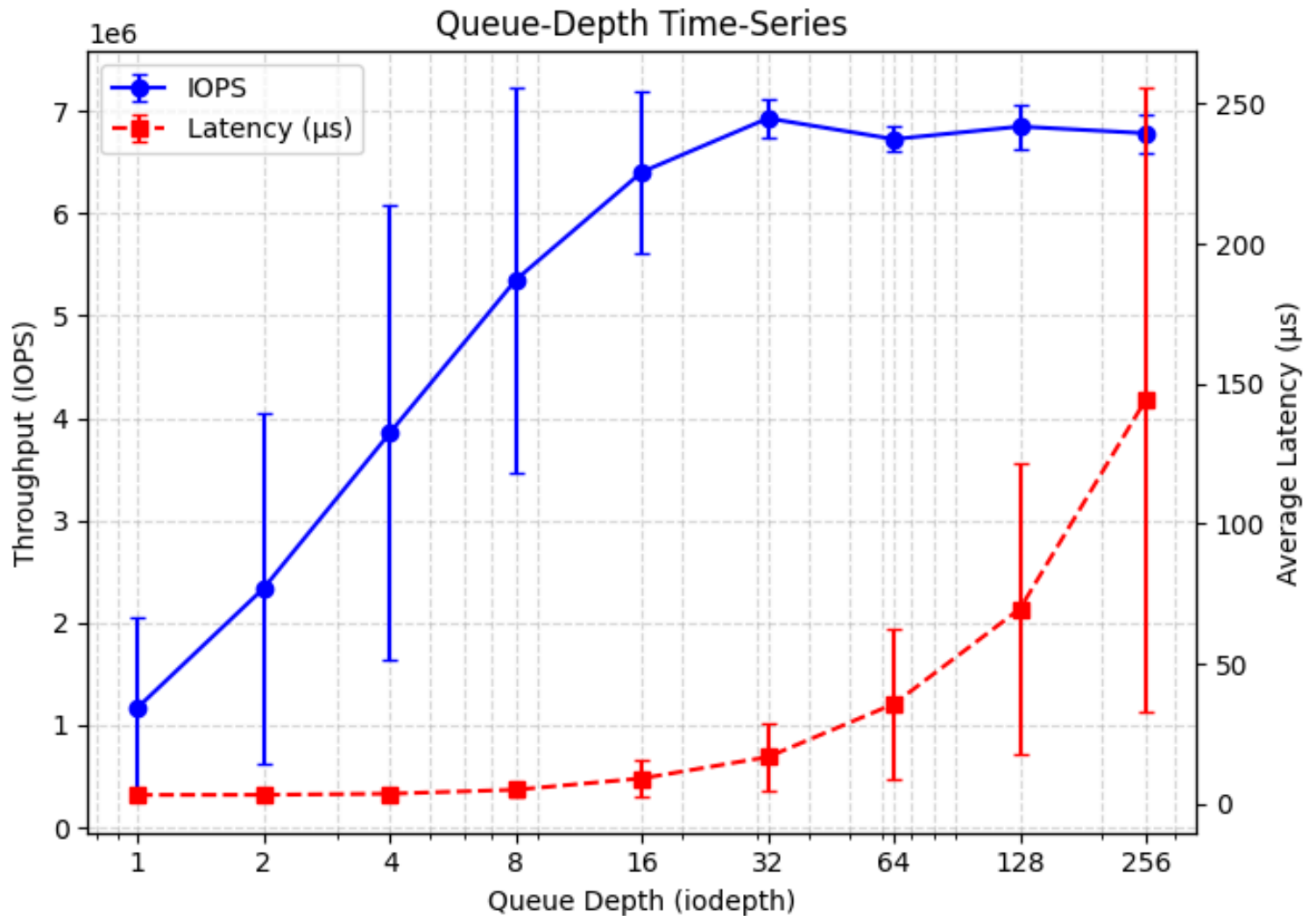Sequential write throughput over 240s, simulating SLC cache behavior.

We can see that the throughput starts out at 7.5GB/s and starts decreasing till it plateaus out at 2.5GB/s.



Sequential Write Time-Series (SLC Cache)

# Queue Depth Time-Series

IOPS and latency vs iodepth (1-256) time series.

Throughput increases as the queue depth decreases and plateaus out at an iodepth of 32. Average latency seems to increases exponentially with Queue Depth.

# Summary Overview

Median latency (avg/p95/p99) across experiments.

The random read generally had a smaller latency than its random write conterpart. For random read the 99th percentile of latencies was only just over being twice as slow as the average. For random write the 99th percentile of latencies was 3 times larger than the average.



Summary of Median Latencies