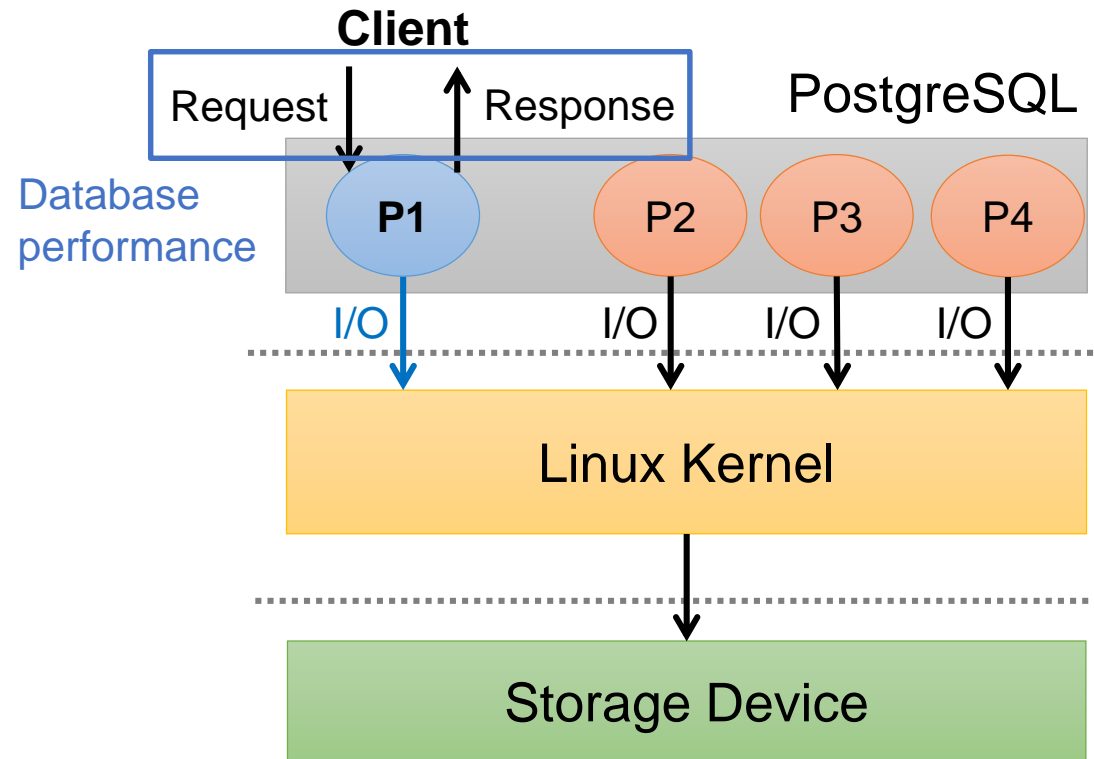


# Taming Performance Variability in PostgreSQL

Shawn S. Kim

apposha

# PostgreSQL Execution Model

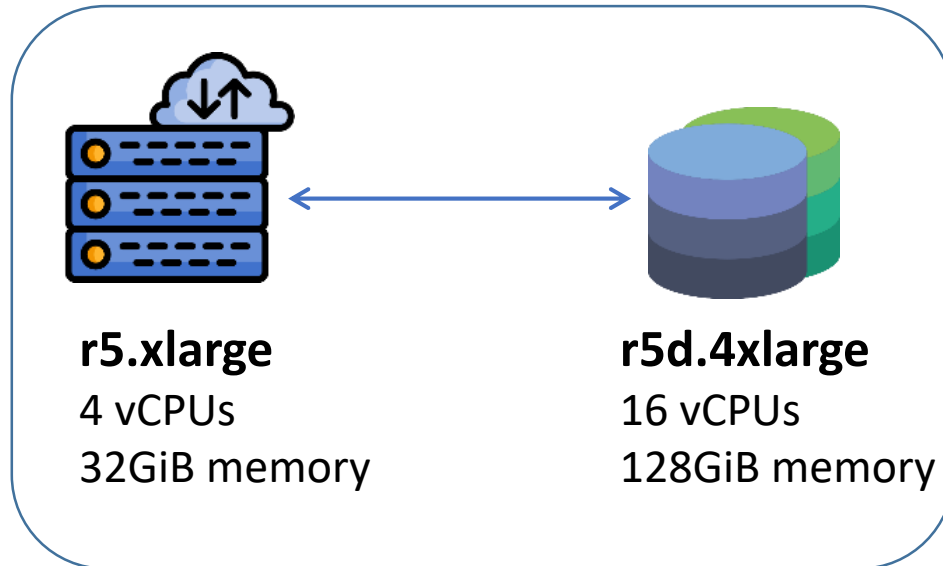


## \* PostgreSQL's processes

- Backend (foreground)
- Checkpointer
- Autovacuum workers
- WAL writer
- Writer
- ...

# Impact of Background Tasks

**SysBench 1.0.15**  
**--oltp-table-size=**  
**10000000**  
**--oltp-tables-count=**  
**24**  
**(50GB dataset)**



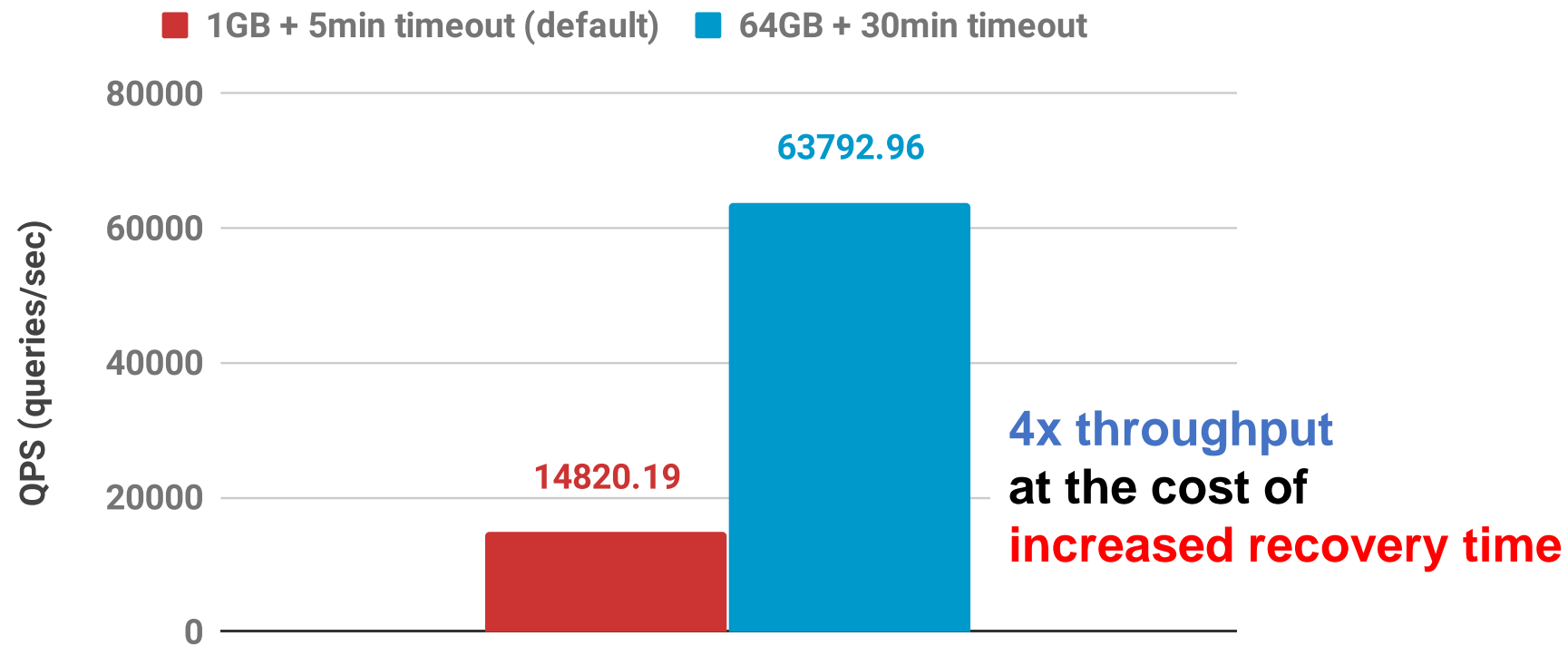
**AWS Seoul region**

**PostgreSQL 11.3**  
**32GB shared\_buffer**  
**64GB effective\_cache\_size**  
...  
**Local NVMe SSD**

# Impact of Background Tasks

## Checkpoint tuning\* improves average throughput

sysbench oltp-write-only throughput, 50 GiB, 100 clients



\* Basics of Tuning Checkpoints (<https://www.2ndquadrant.com/en/blog/basics-of-tuning-checkpoints>)

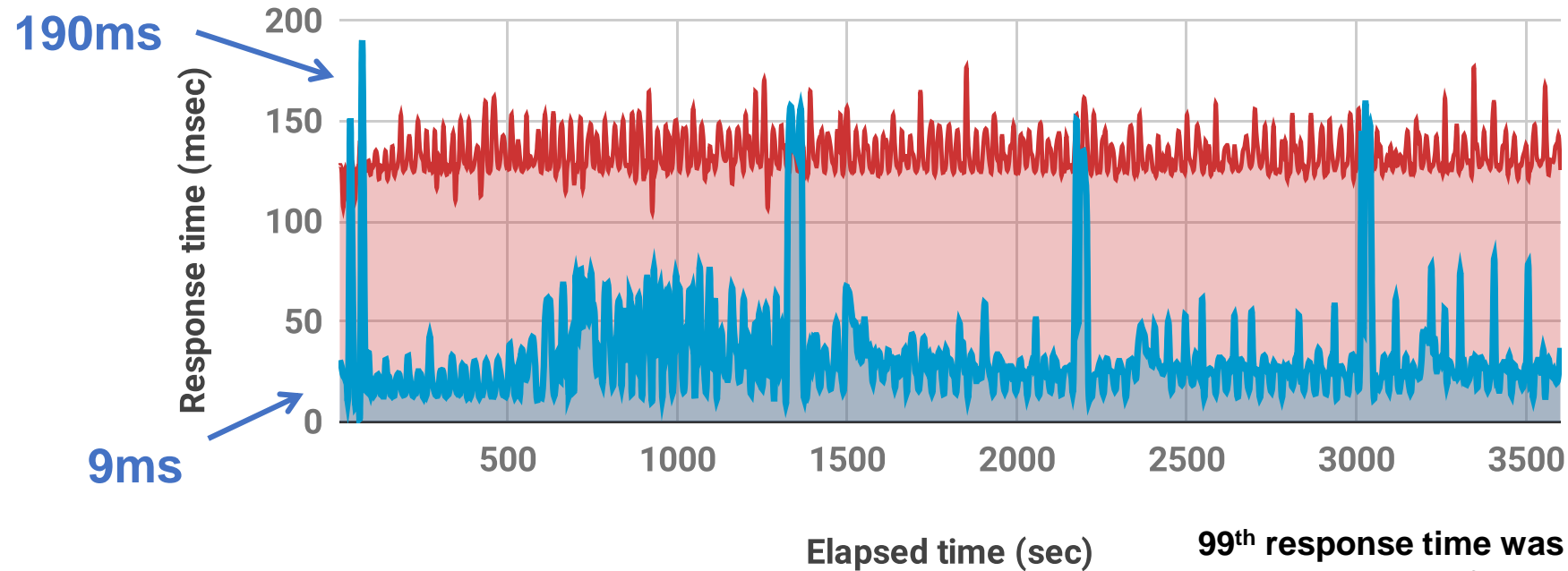
**Done..?**

# Performance Variability

## Checkpoint tuning **makes PostgreSQL unpredictable**

sysbench oltp-write-only p99 response time, 50GiB, 100 clients

■ 1GB + 5min timeout (default) ■ 64GB + 30min timeout



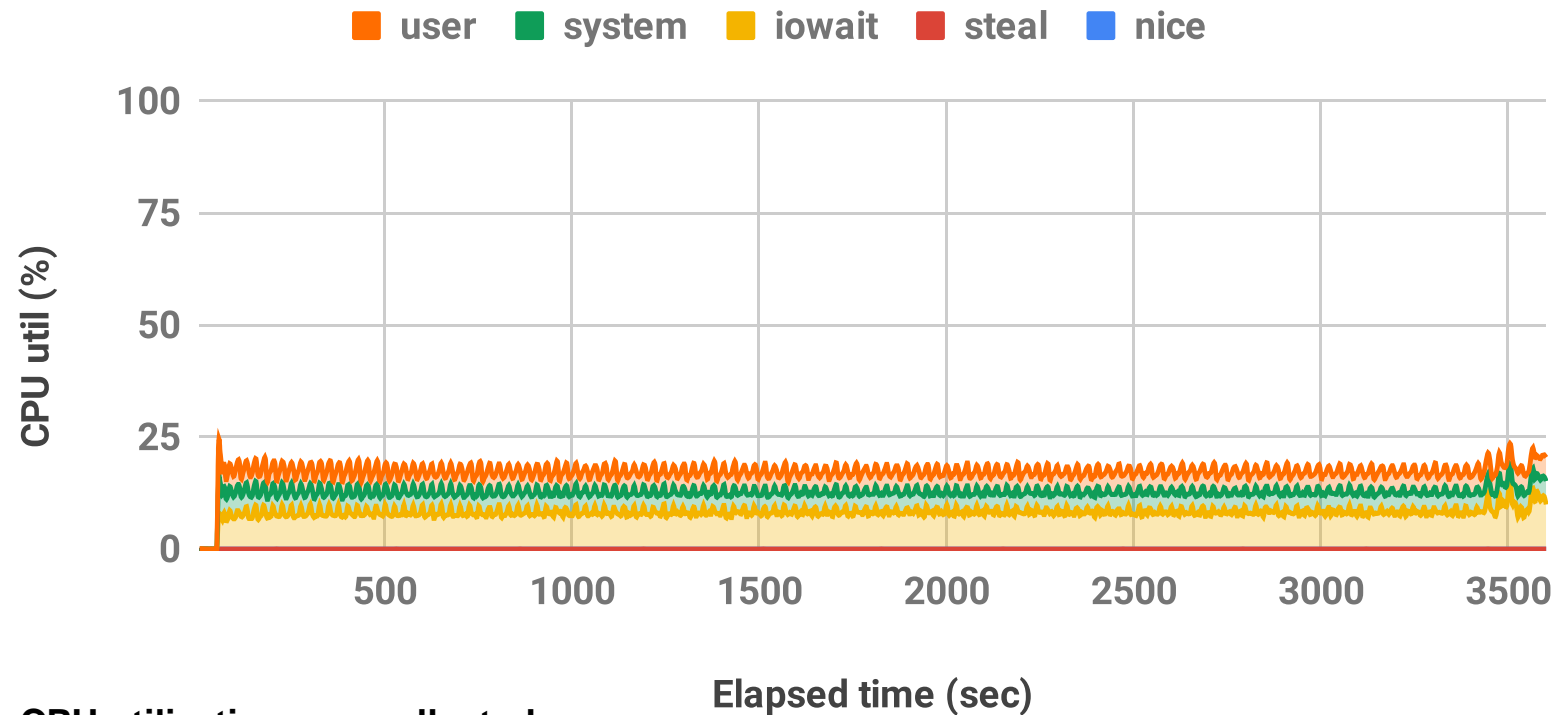
**99<sup>th</sup> response time was collected every 5 seconds for both cases**

**Default is better  
in terms of variability**

# What's the Problem?

**iowait** is the main bottleneck for the frequent checkpoint case

CPU utilization with 1GB + 5min timeout



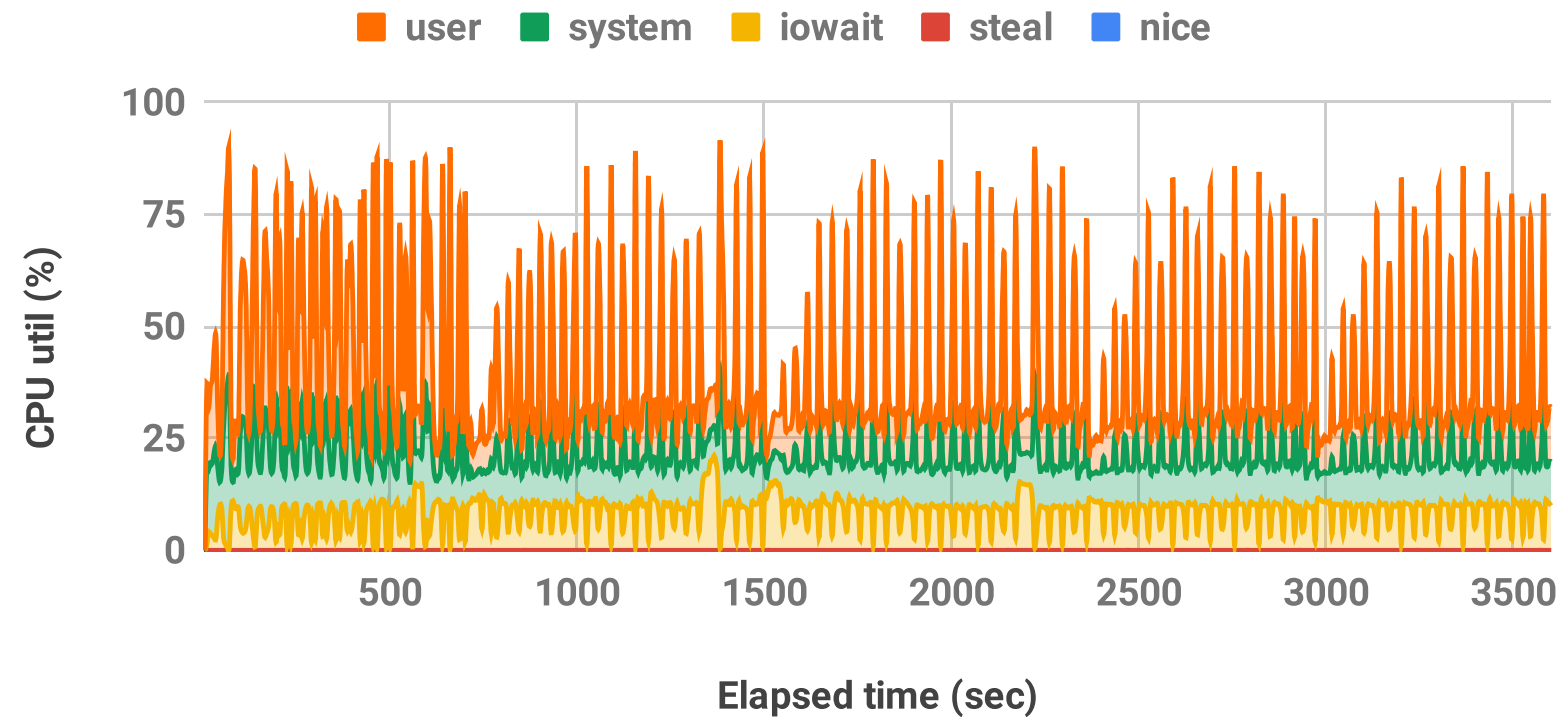
CPU utilization was collected every 5 seconds using sar



# What's the Problem?

**CPU utilization is highly fluctuated even without checkpoint**

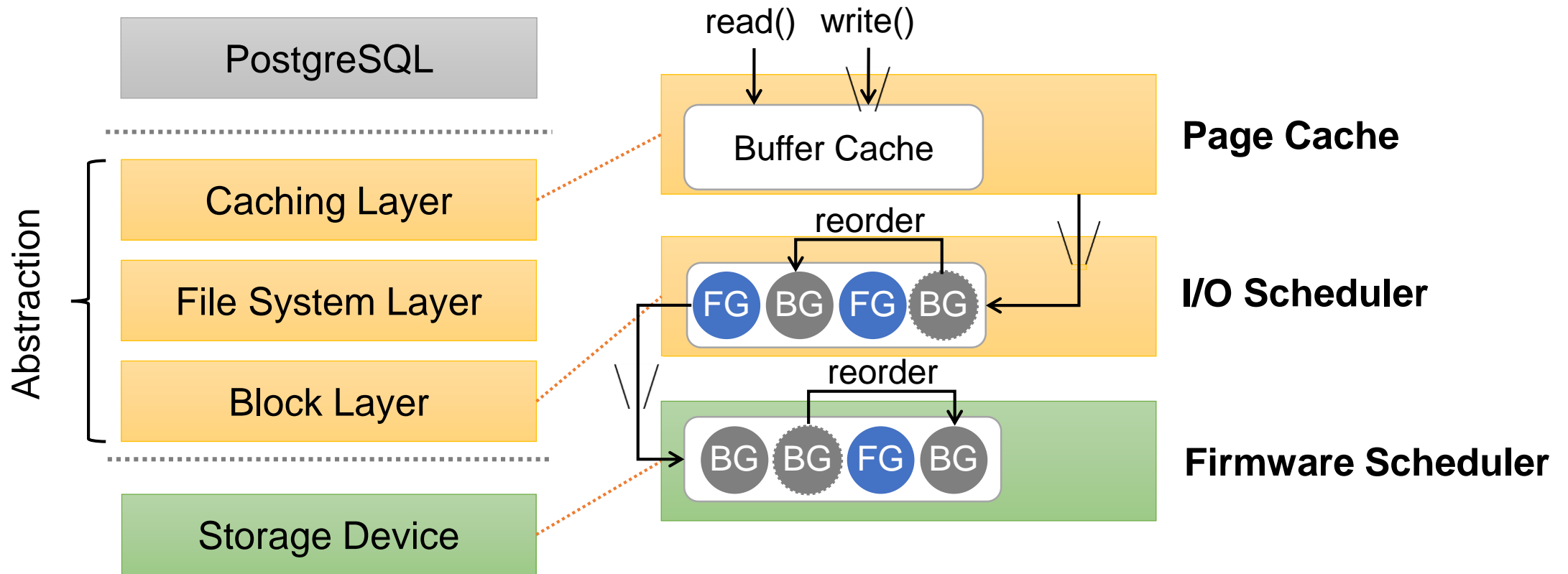
CPU utilization with 64GB + 30min timeout



CPU utilization was collected every 5 seconds using sar

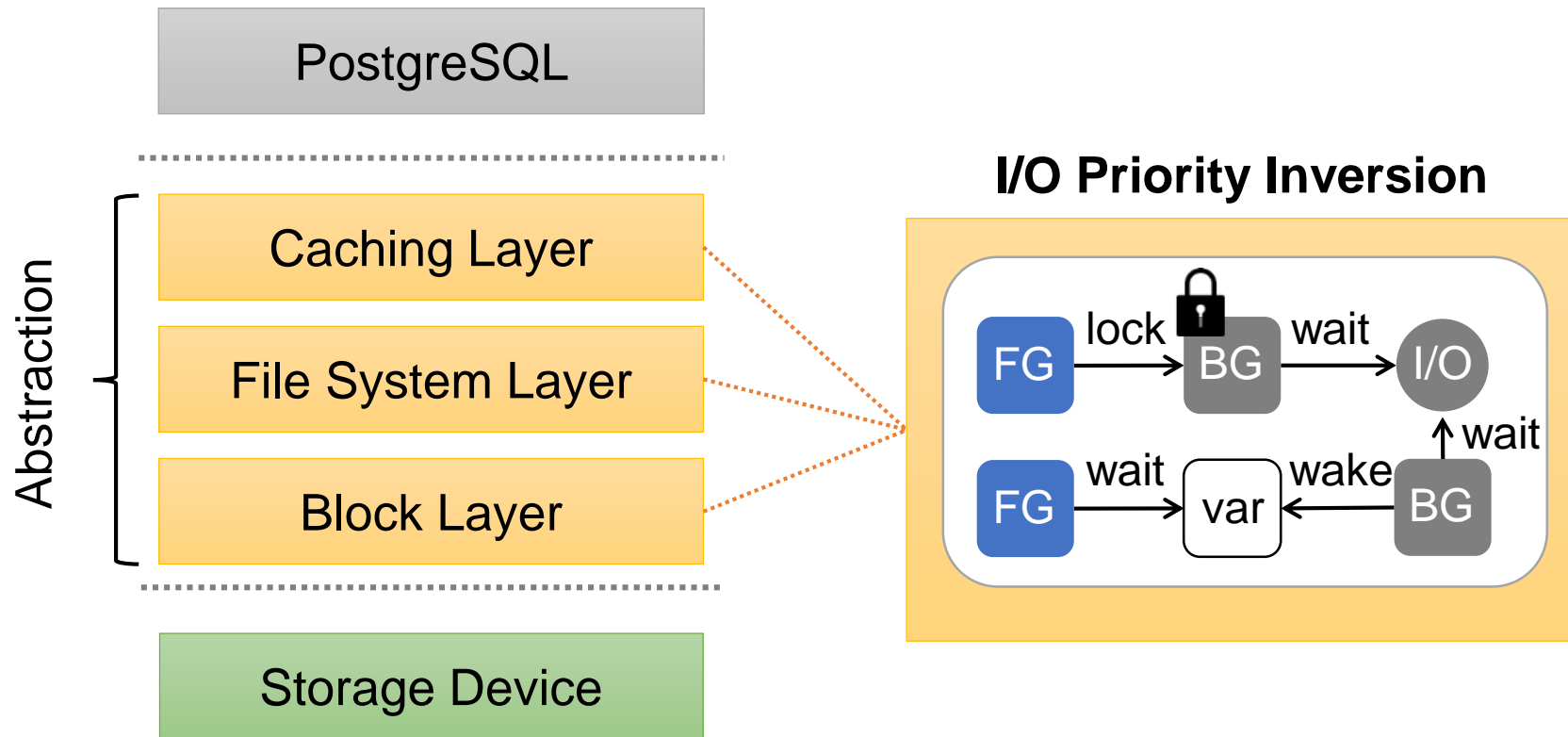
# What's the Problem?

**Background I/Os interfere** foreground I/Os inside Linux



# What's the Problem?

Linux makes backends to **indirectly wait** for background I/Os

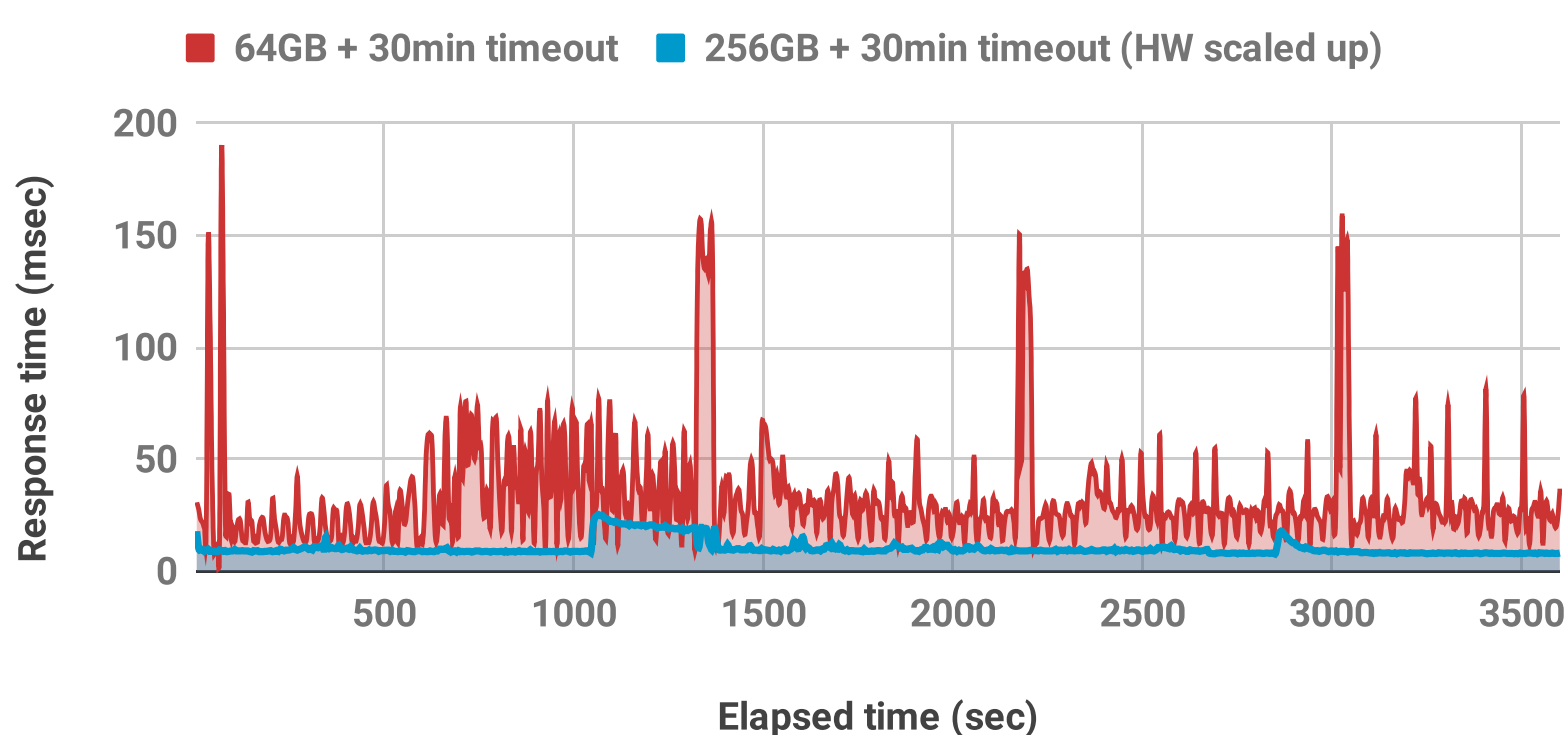


# Hardware Upgrade to the Rescue

## 6x HW scale up mostly resolves the variability

**31.37 ms vs 9.56 ms**

sysbench oltp-write-only p99 response time, 50GiB, 100 clients



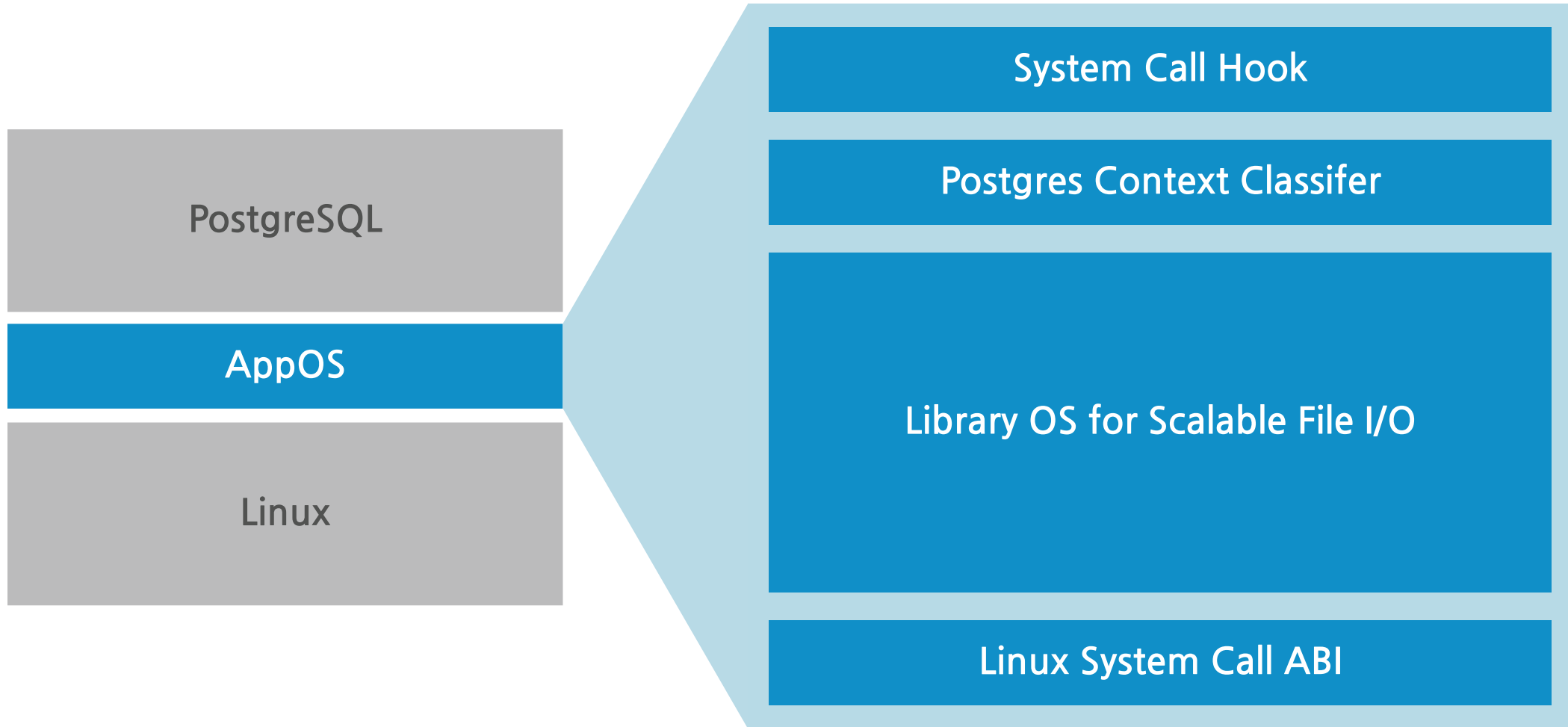
● 6x CPU cores +  
6x DRAM +  
10x storage throughput  
= **6x infra cost**

Whole dataset fit in  
PostgreSQL buffer

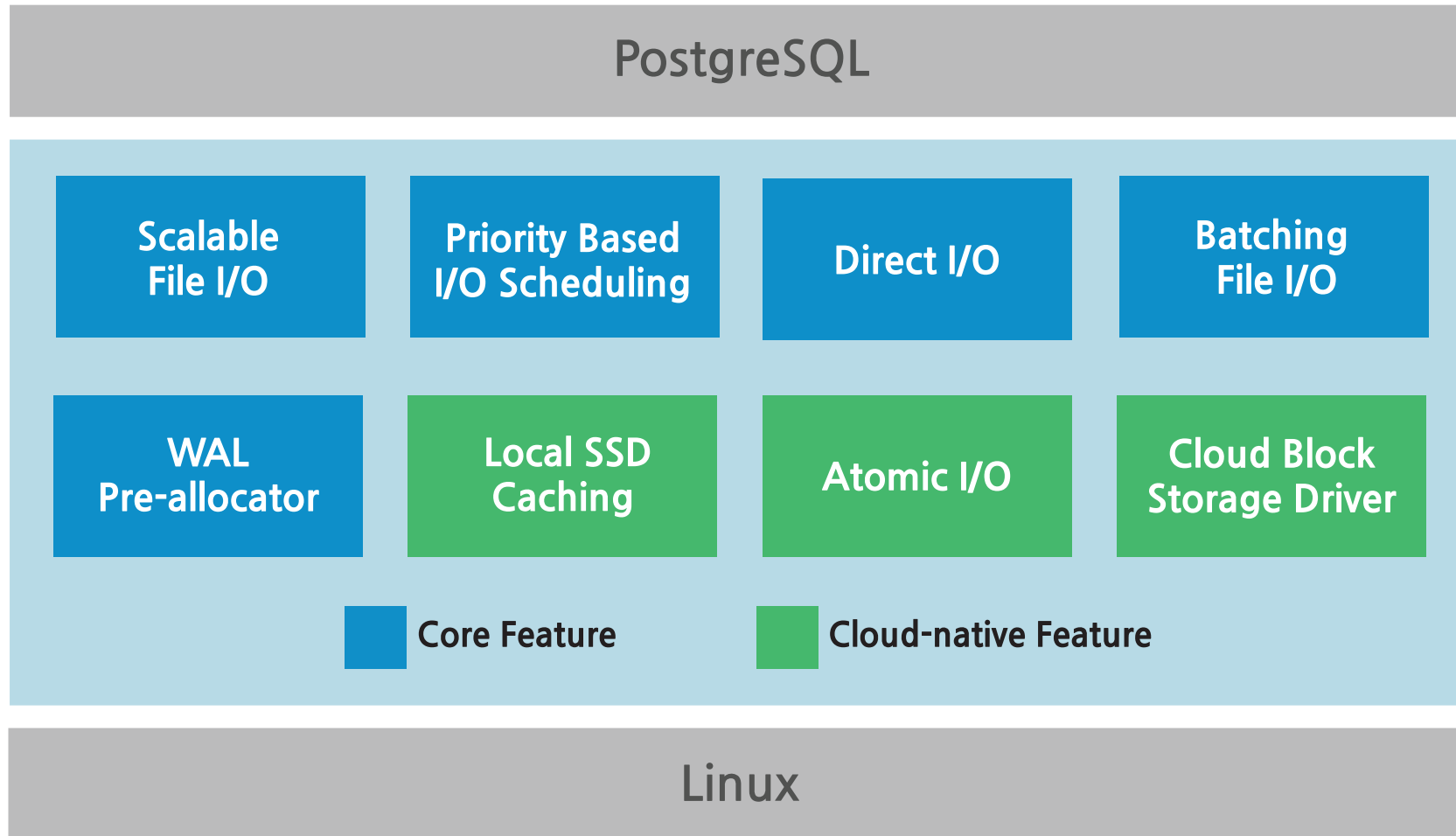
# Software solution..?

**AppOS** is a PostgreSQL extension that provides **specialized file I/O stack**

# AppOS Extension



# AppOS Extension





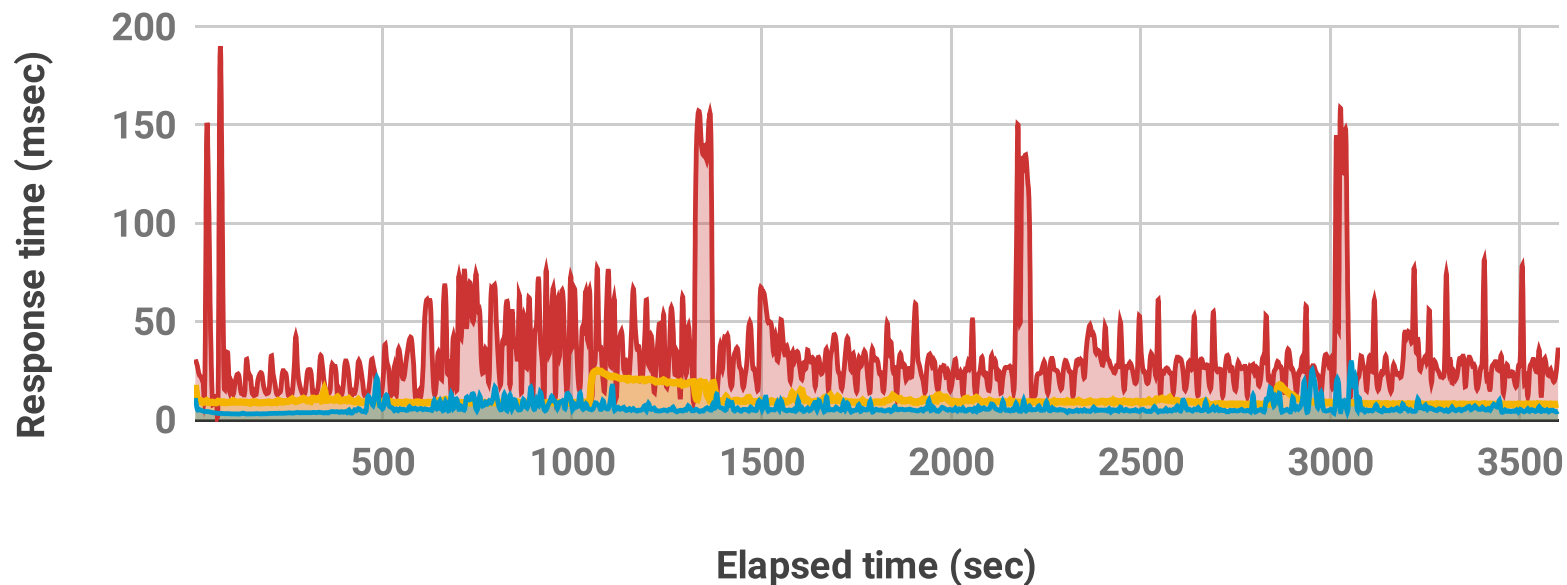
# Performance Variability with AppOS

**AppOS extension efficiently resolves the variability**

**31.37 ms vs 9.56 ms vs 5.77 ms**

sysbench oltp-write-only p99 response time, 50GiB, 100 clients

■ 64GB + 30min timeout   ■ 256GB + 30min timeout (HW scaled up)  
■ 64GB + 30min timeout + AppOS extension

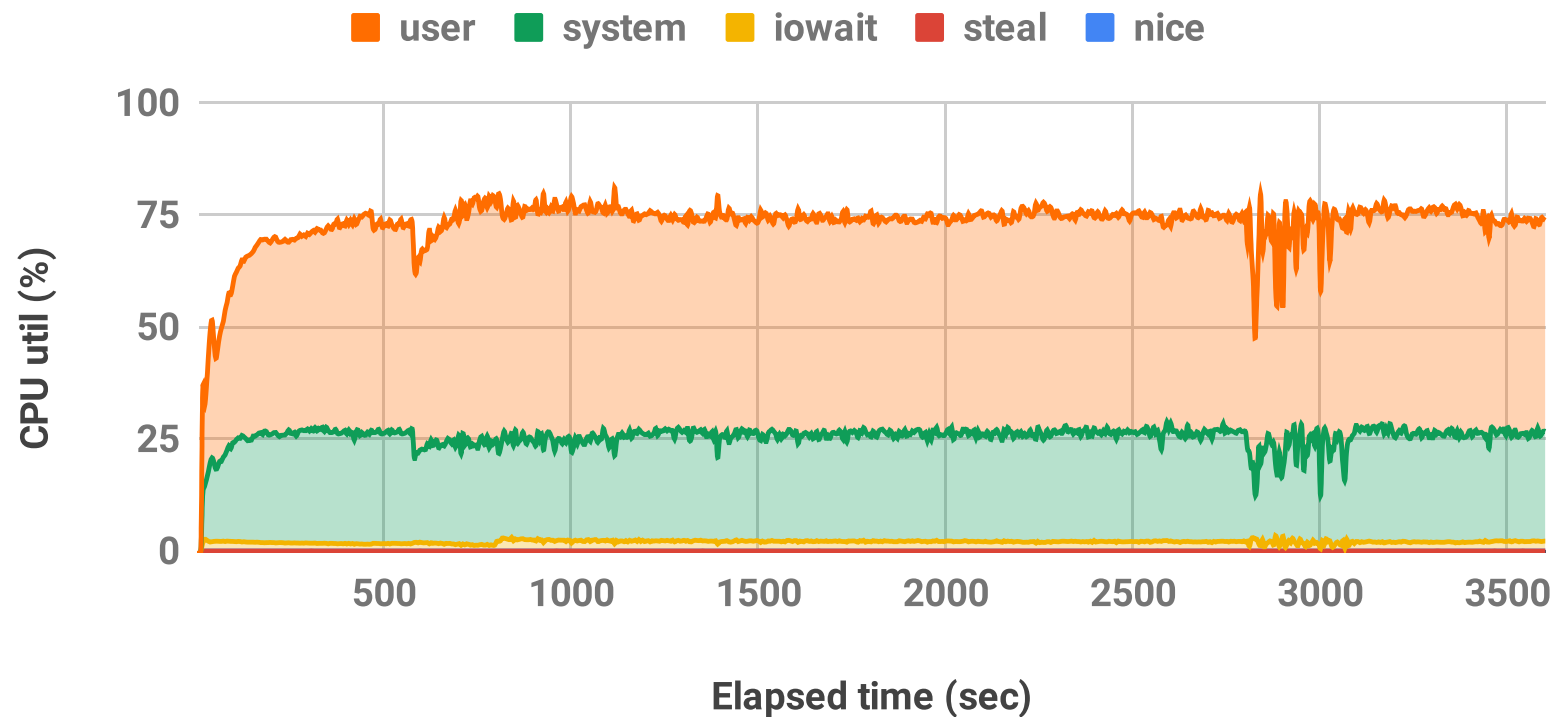


**Stable response  
without extra HW**

# Performance Variability with AppOS

## AppOS extension efficiently resolves the variability

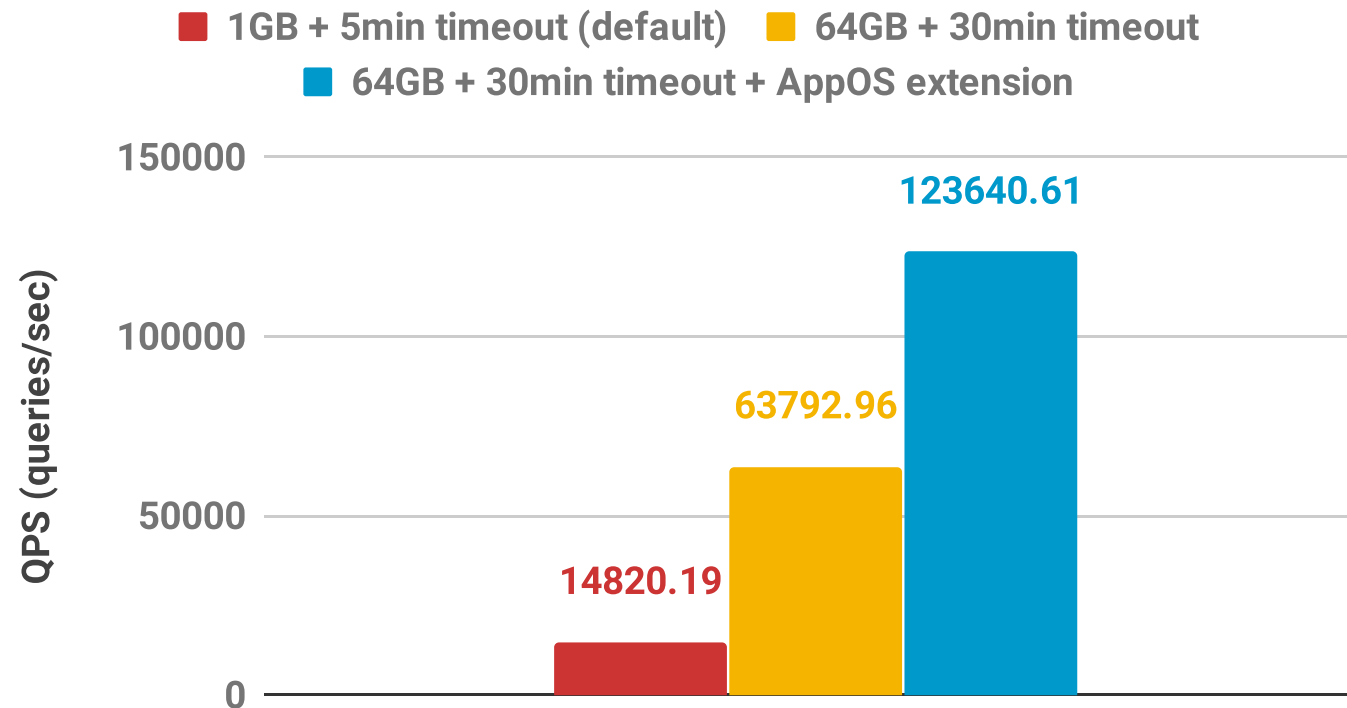
CPU utilization with 64GB + 30min timeout + AppOS extension



# Performance Variability with AppOS

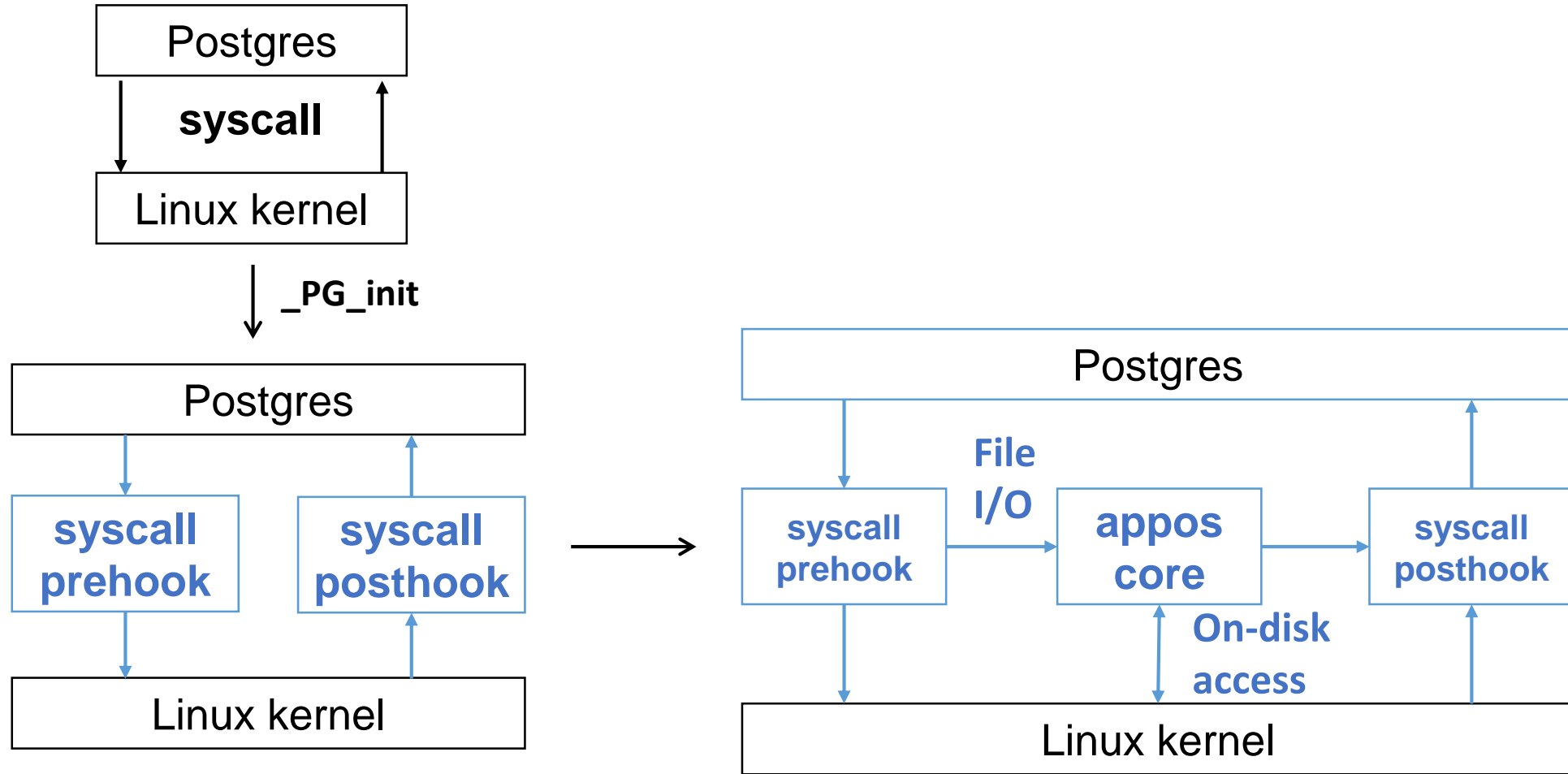
## AppOS extension improves average throughput as well

sysbench oltp-write-only throughput, 50 GiB, 100 clients

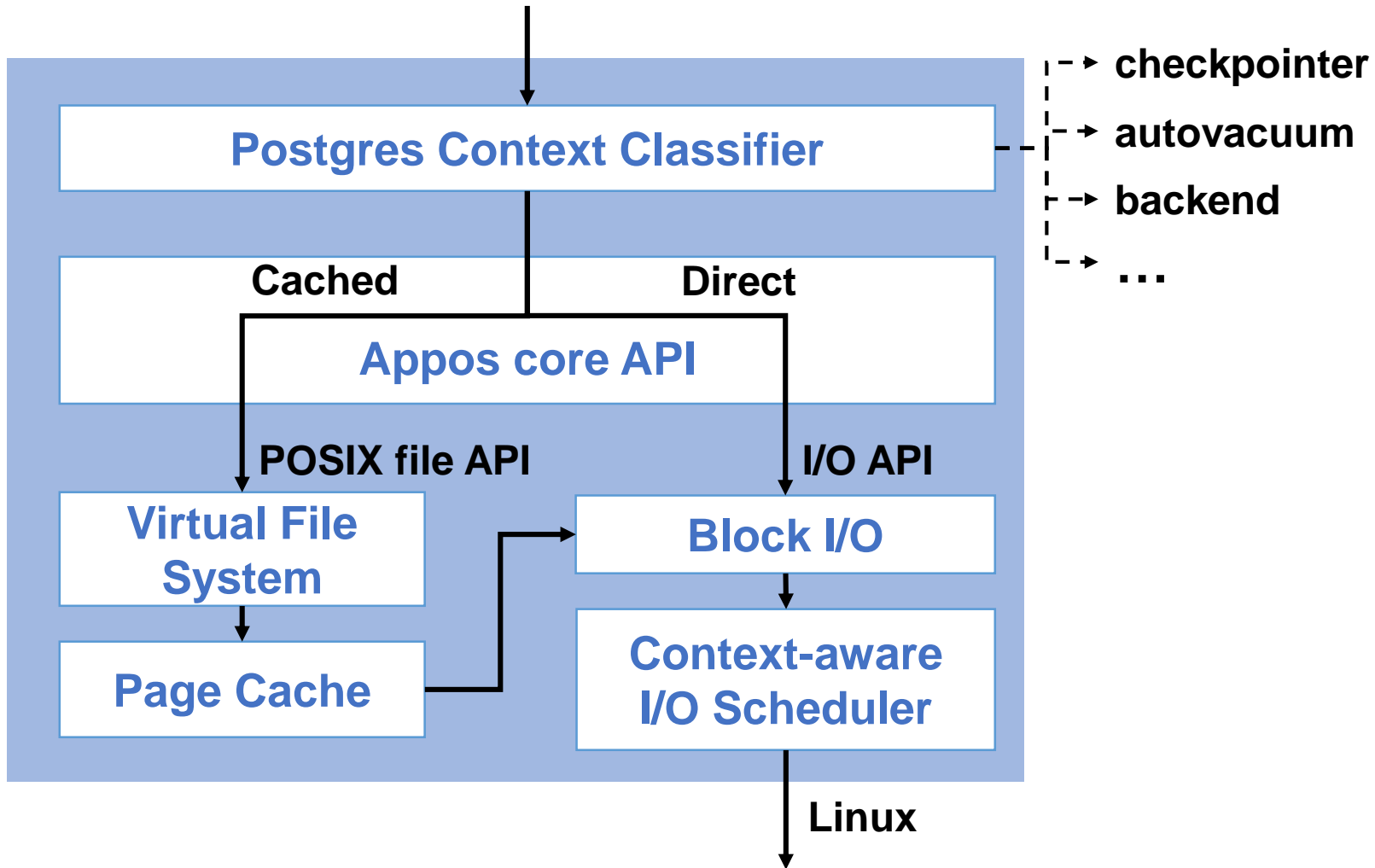


**2x throughput  
with AppOS extension**

# AppOS Internals

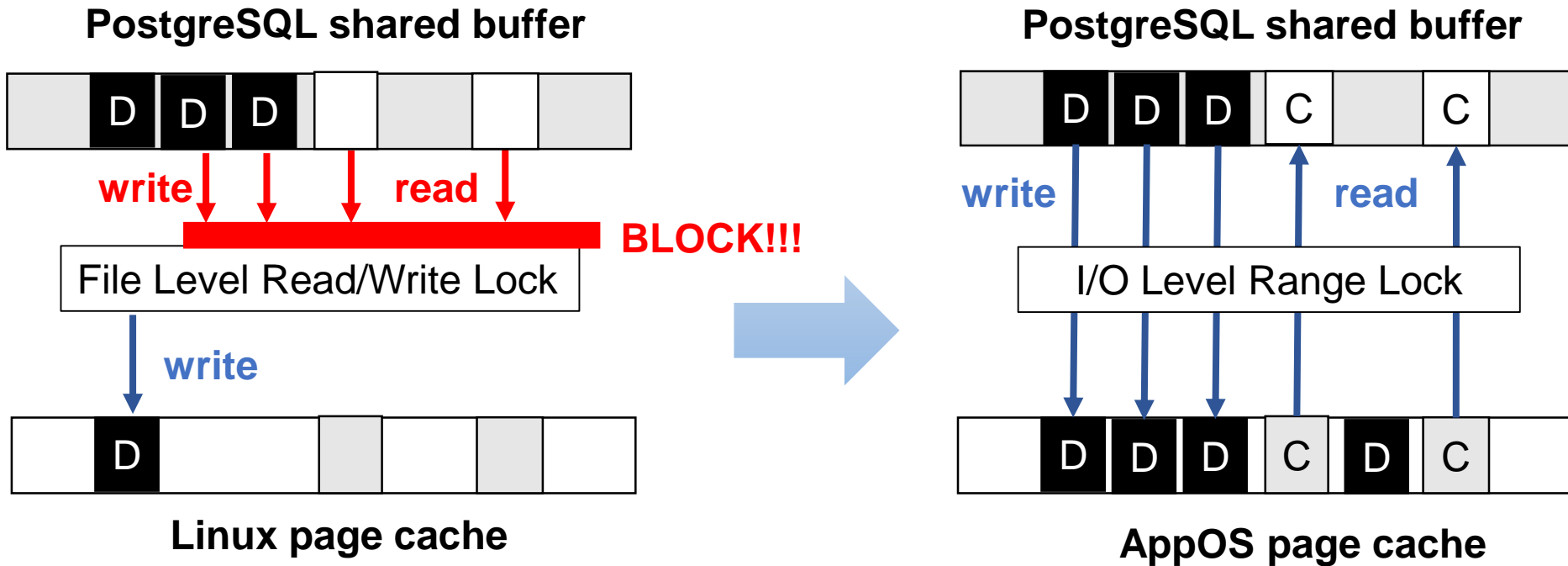


# AppOS Internals



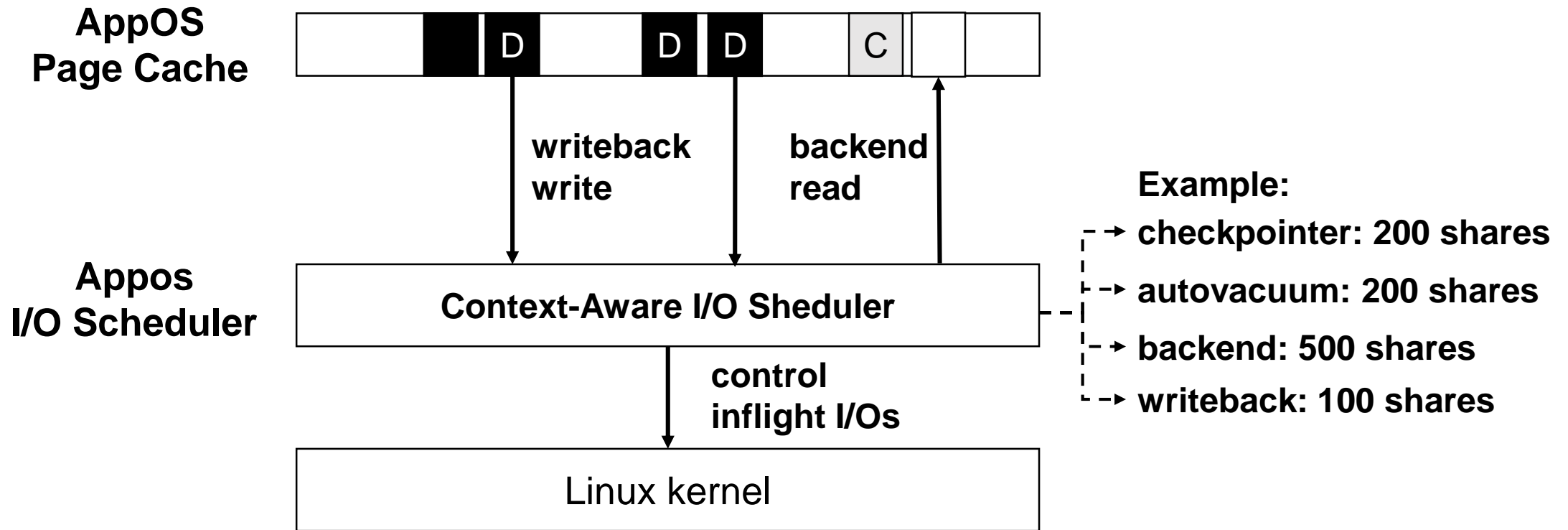
# AppOS Internals

AppOS uses **range lock** instead of big file lock



# AppOS Internals

**AppOS schedules I/Os based on context and congestion**



## Use Cases (1)

**AppOS makes PostgreSQL more predictable**

**Real-time SLA      Autovacuum      Replication lag**

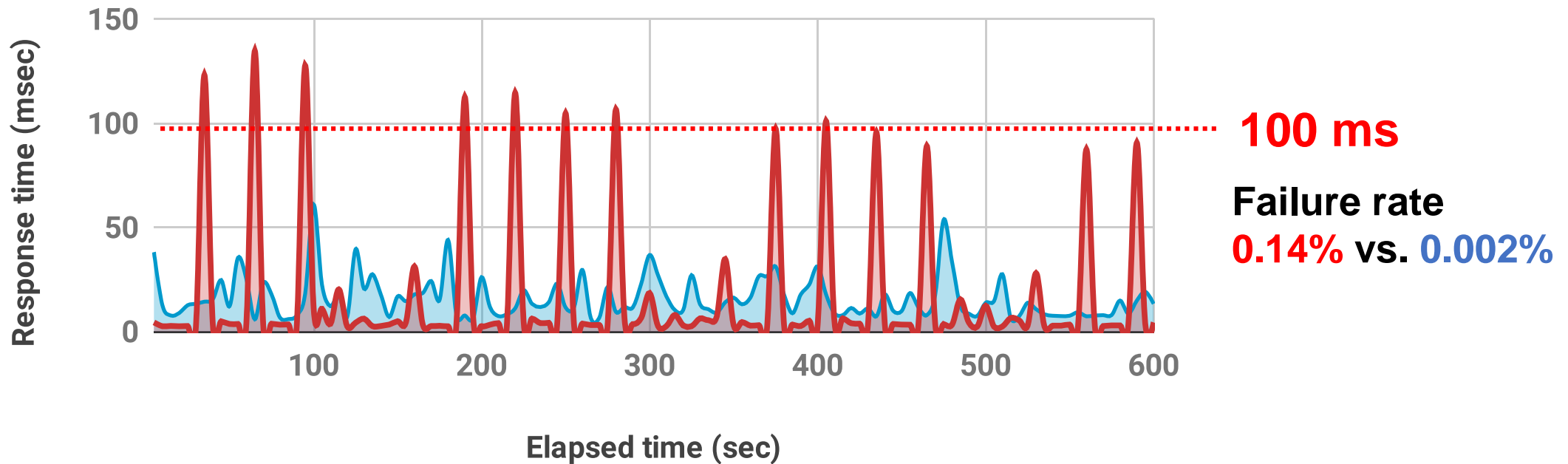


# Use Cases (1)

## Real-time SLA

sysbench oltp-insert p99 response time, 50GiB

- 4 core + 15GB mem + 300GB SSD + AppOS
- 64 core + 240GB mem + 2TB SSD

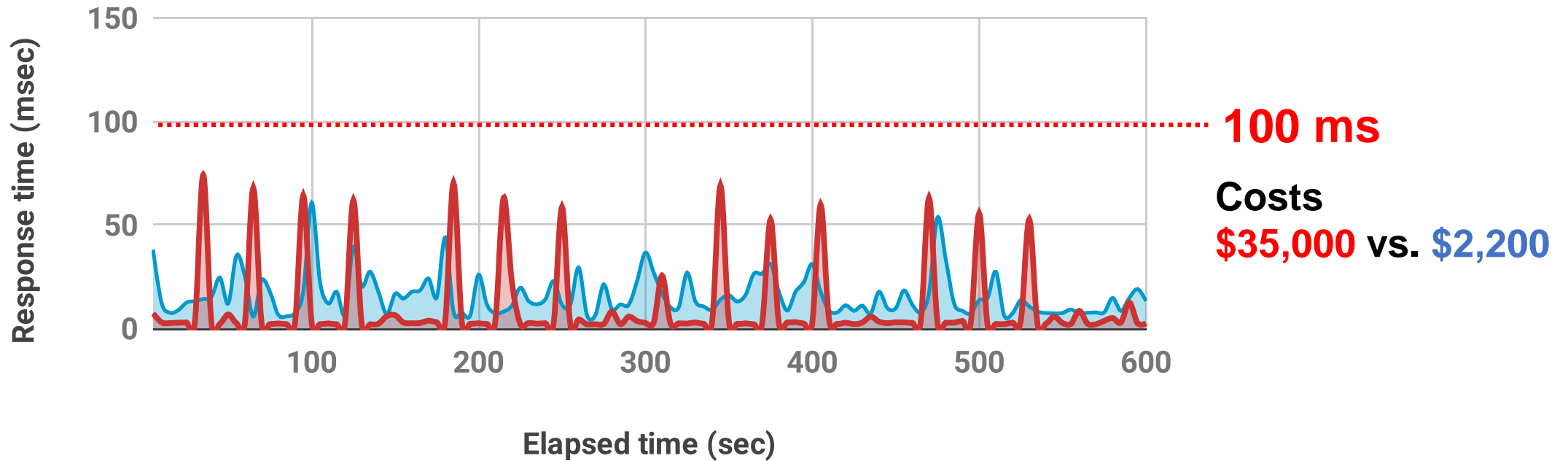


# Use Cases (1)

## Real-time SLA

sysbench oltp-insert p99 response time, 50GiB

- 4 core + 15GB mem + 300GB SSD + AppOS
- 96 core + 360GB mem + 2TB SSD



# Use Cases (1)

## Autovacuum



**Vacuum**

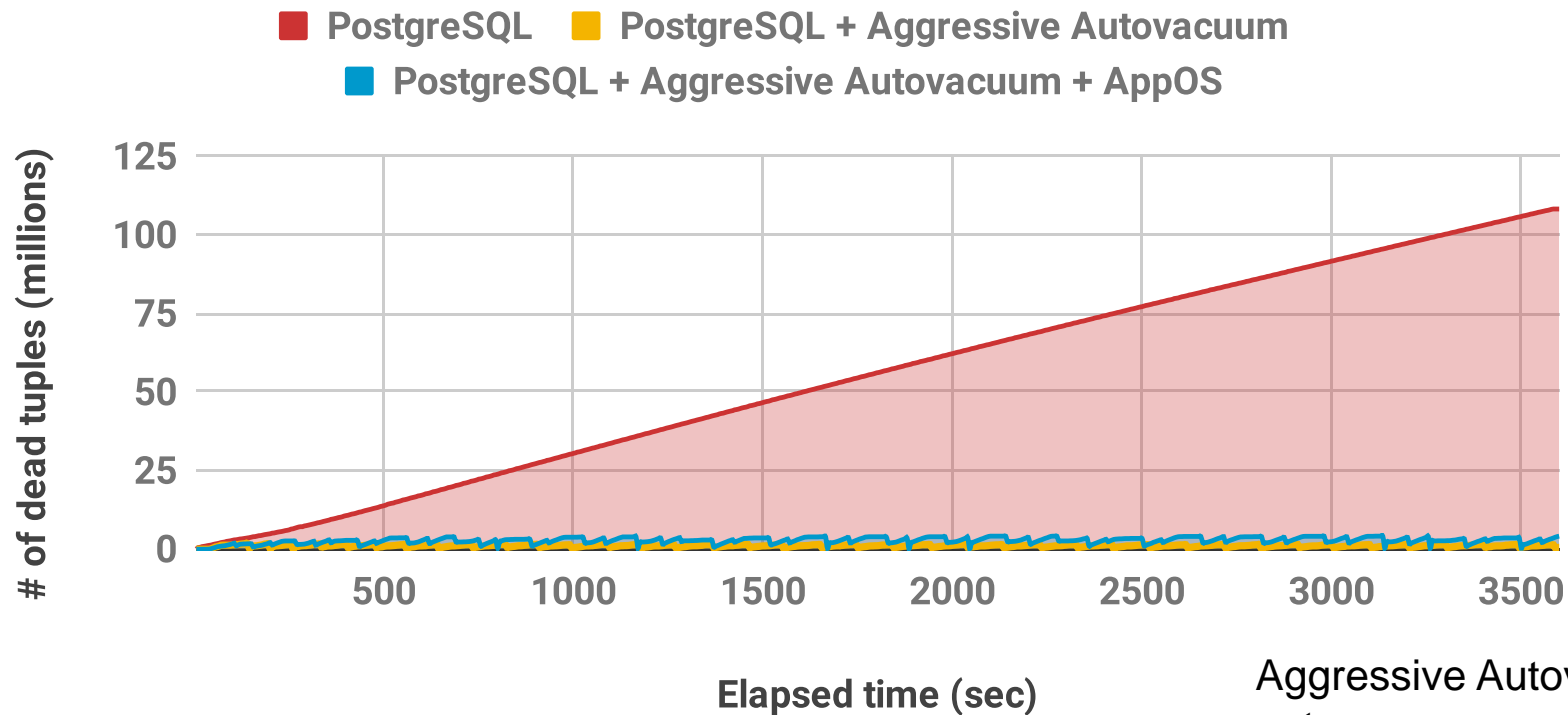


**Analyze**

# Use Cases (1)

## Autovacuum

sysbench oltp-write-only # of dead tuples, 50GiB, 100 clients



Aggressive Autovacuum:  
autovacuum\_vacuum\_cost\_delay=1s  
autovacuum\_vacuum\_cost\_limit=10000

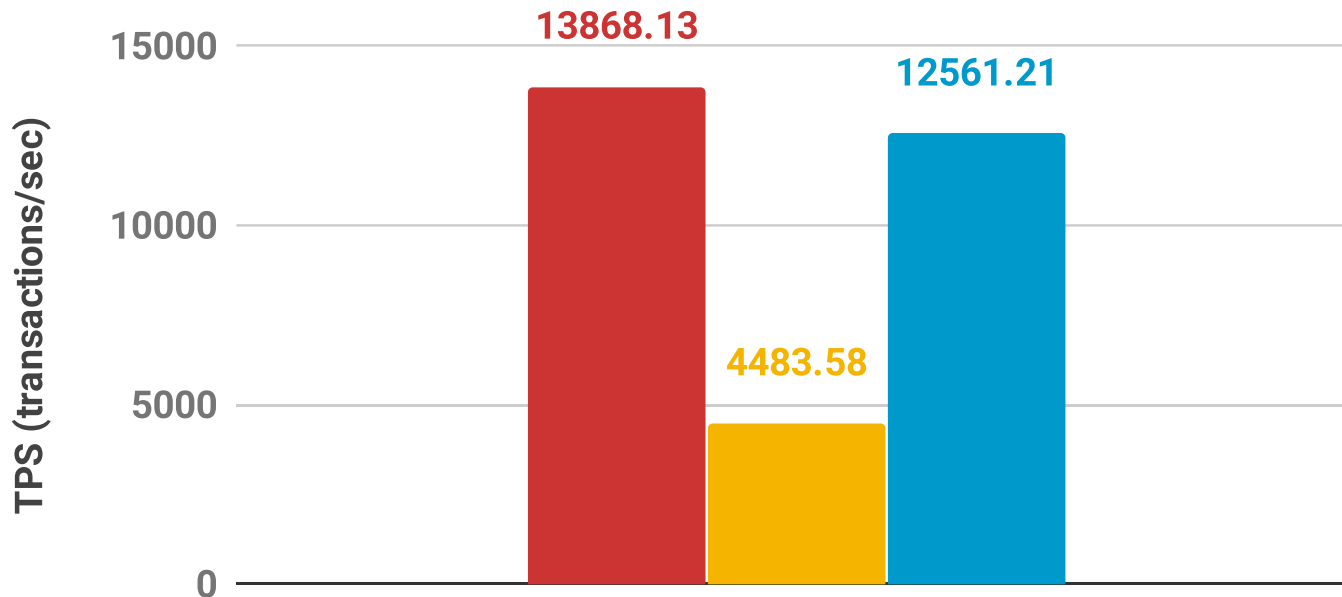
# Use Cases (1)

## Autovacuum

**13,868 TPS vs 4483 TPS vs 12,561 TPS**

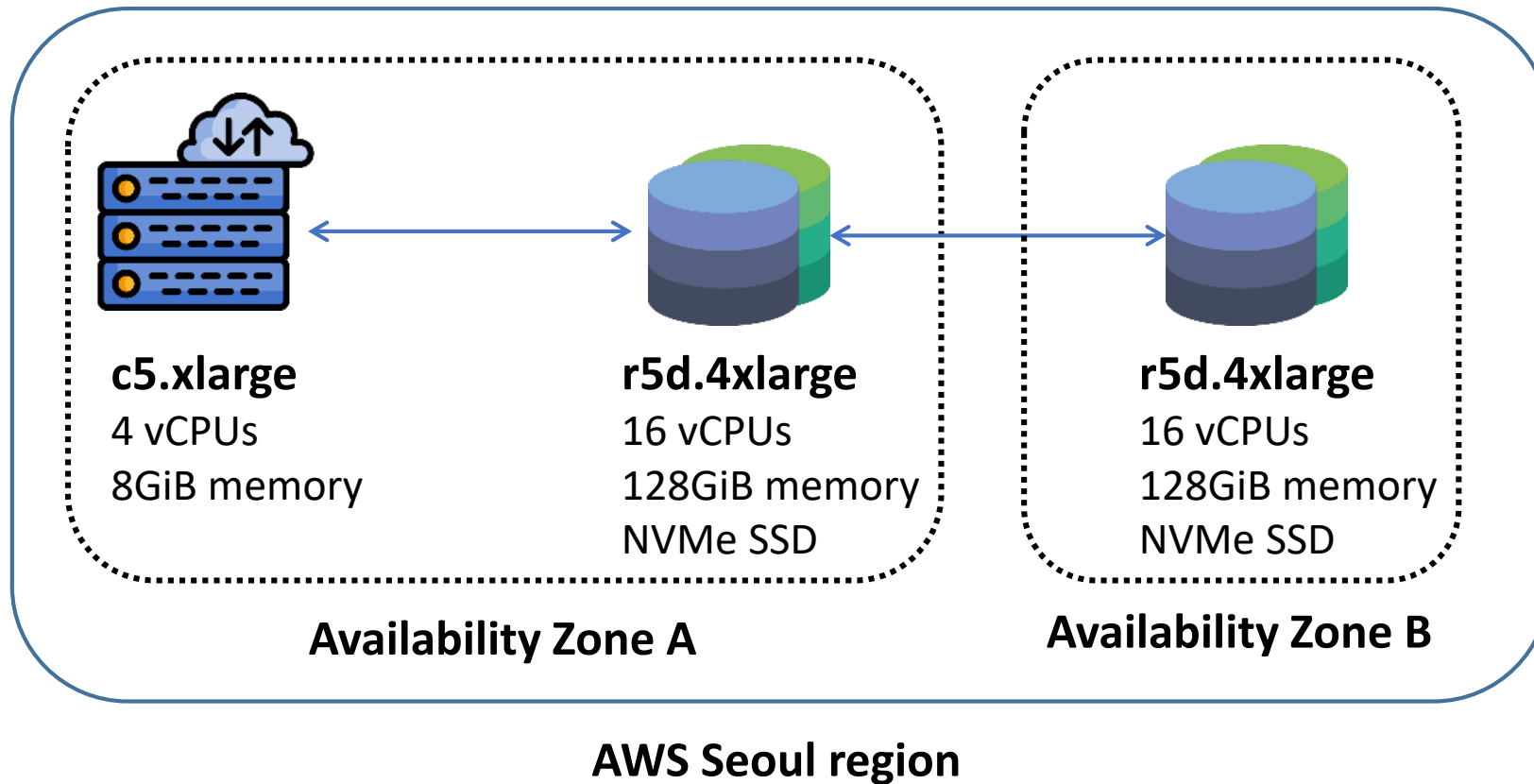
sysbench oltp-write-only throughput, 50 GiB, 100 clients

■ PostgreSQL   ■ PostgreSQL + Aggressive Autovacuum  
■ PostgreSQL + Aggressive Autovacuum + AppOS



# Use Cases (1)

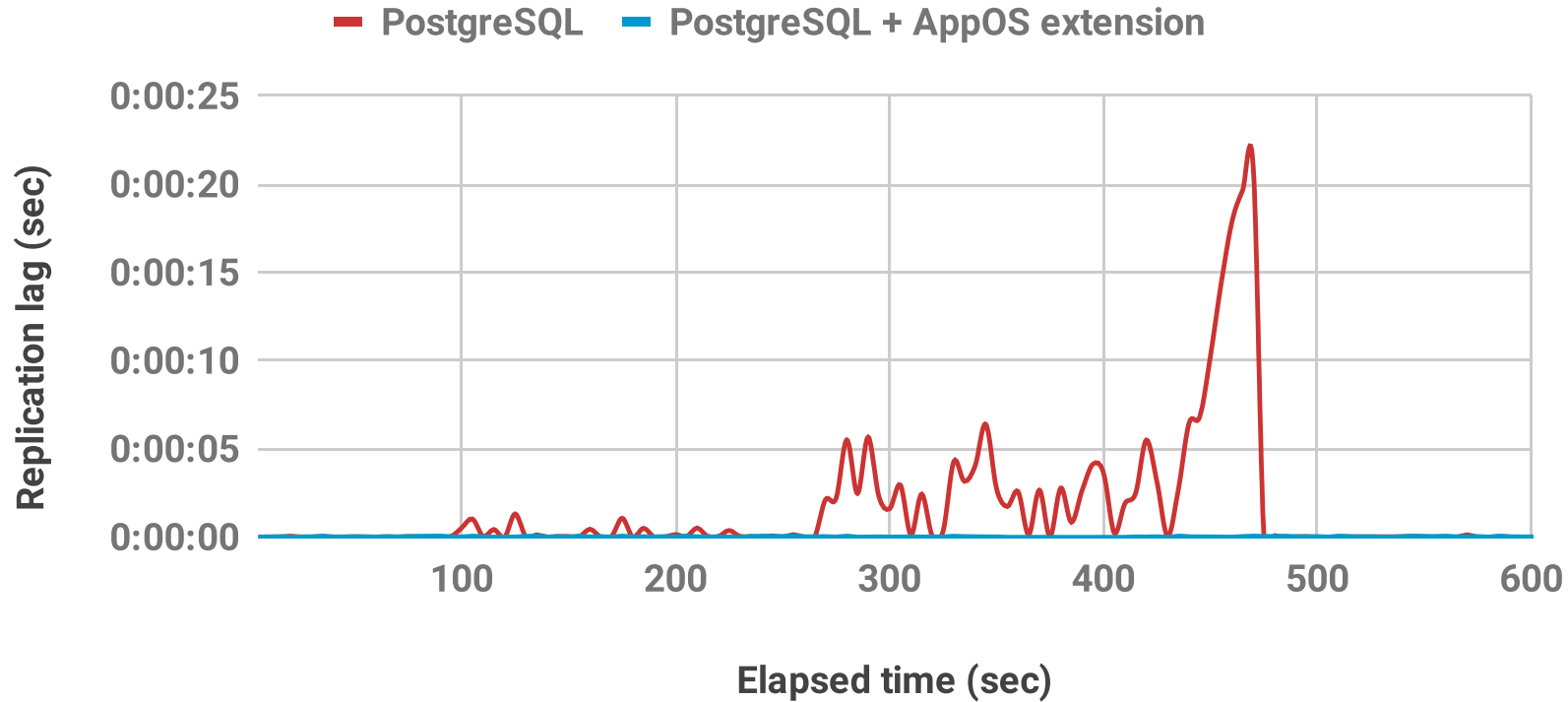
## Replication lag



# Use Cases (1)

## Replication lag

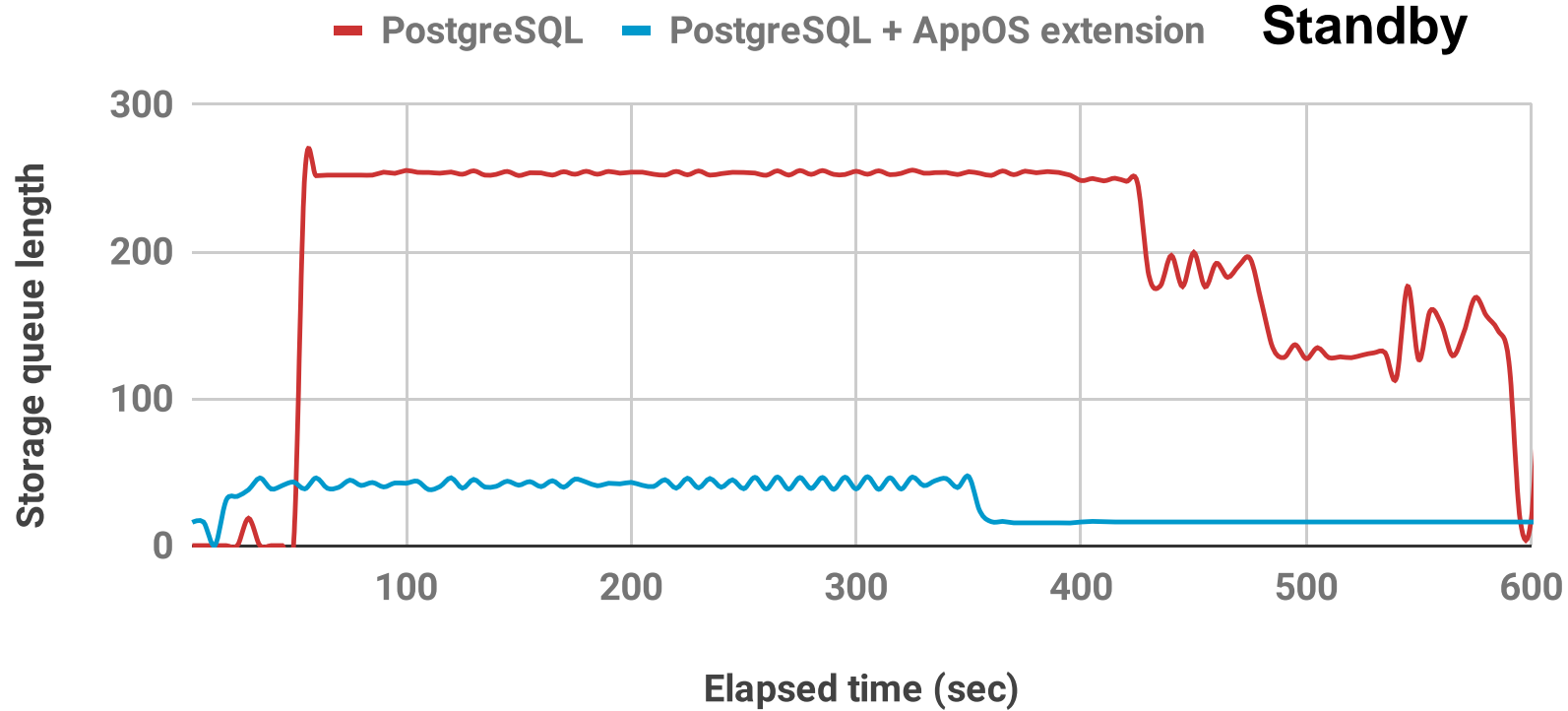
pgbench data loading replication lag, 3k scale



# Use Cases (1)

## Replication lag

pgbench data loading storage queue length, 3k scale





## Use Cases (2)

**AppOS makes PostgreSQL cloud storage-native**

**Cloud block  
storage**

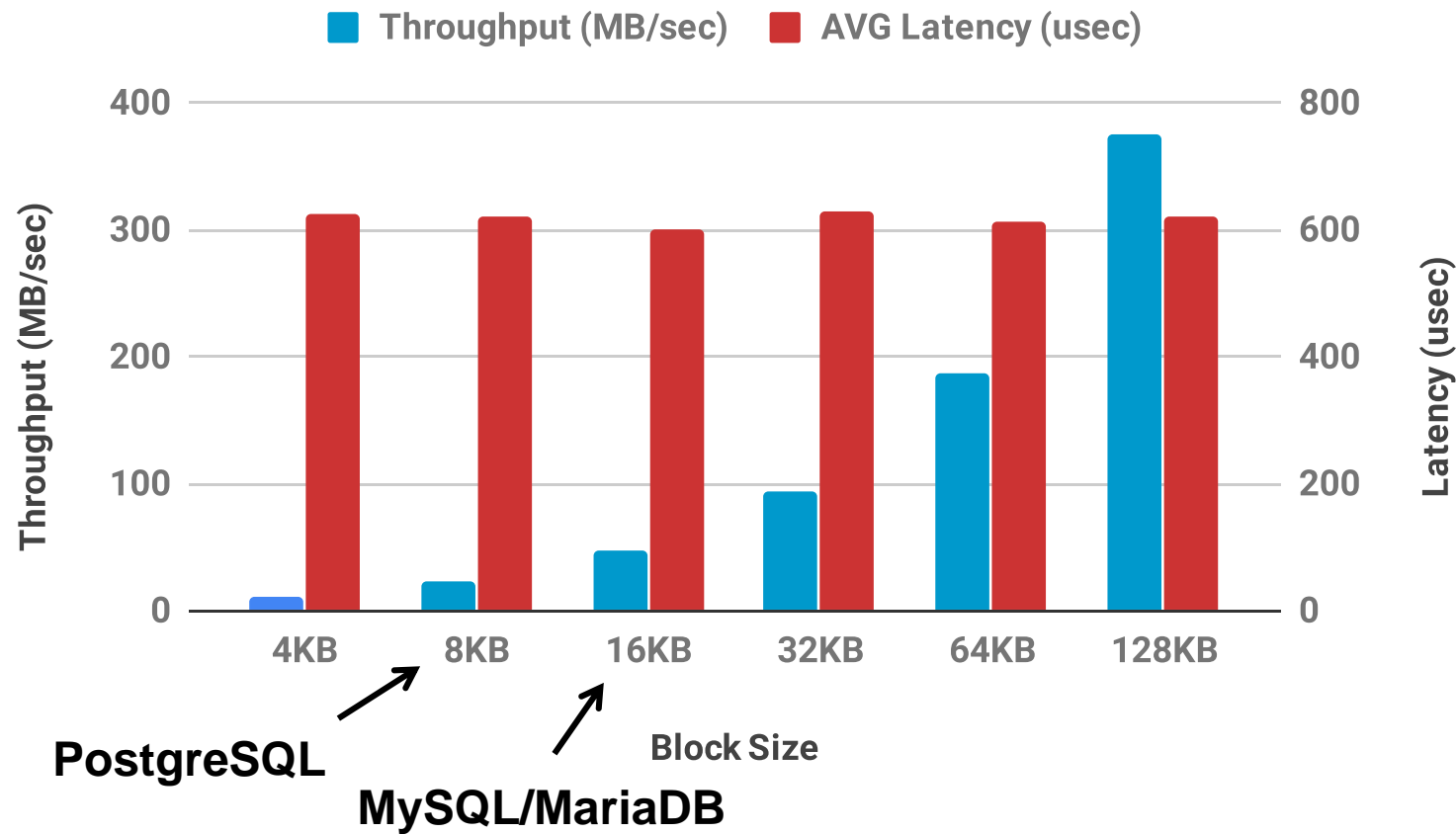
**Atomic write  
support**

**Local SSD**

# Use Cases (2)

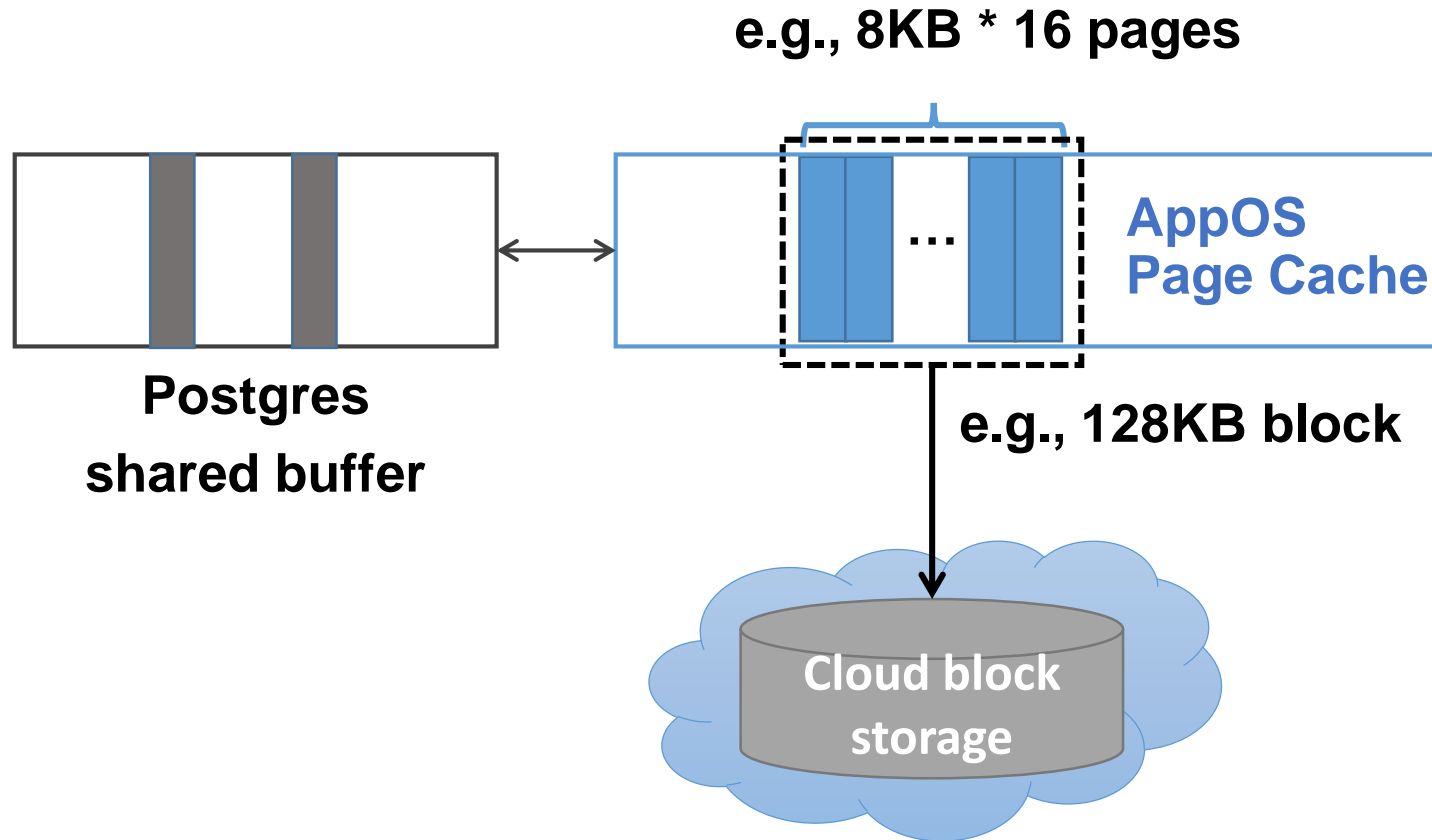
## Cloud block storage

FIO random I/O test with 3000 IOPS AWS EBS



## Use Cases (2)

### Cloud block storage



# Use Cases (2)

## Atomic write support

### Compute Products

Contact sales

- instances
  - Managing your instances
    - Creating and managing groups of instances
  - Networking
  - Deploying containers
  - Scaling your application
  - Monitoring activity
  - Labeling resources
  - Granting access to Compute Engine resources
  - Working with regions and zones
  - Migrating VMs to Compute Engine
  - Advanced VM configurations
  - ▾ Best practices
    - Designing robust systems

### Building a 16 KB atomic write path from database to block device

You can build an end-to-end 16 KB atomic write path from the database to the block device leveraging a 16 KB persistent disk, so you can safely disable the doublewrite feature in MySQL/ InnoDB and achieve a more stable and better performance for a high-write load.

Create and attach a persistent disk through the [Google Cloud Platform Console](#), the [gcloud tool](#), or the [API](#).

1. [Create a 16 KB block size persistent disk](#) and attach it to your VM. The 16 KB persistent disk provides 16 KB write atomicity at the physical block level.

Although optional, it is recommended that you configure your MySQL instance to store data files only to the 16 KB persistent disk. Store log files, especially redo log and binlogs, to a 4 KB persistent disk that is attached to the same VM. This ensures log file writes continue to be high performance because small log writes on 16 KB persistent disk might trigger lots of read-modify-writes, which are slower.

#### Contents

Before you begin

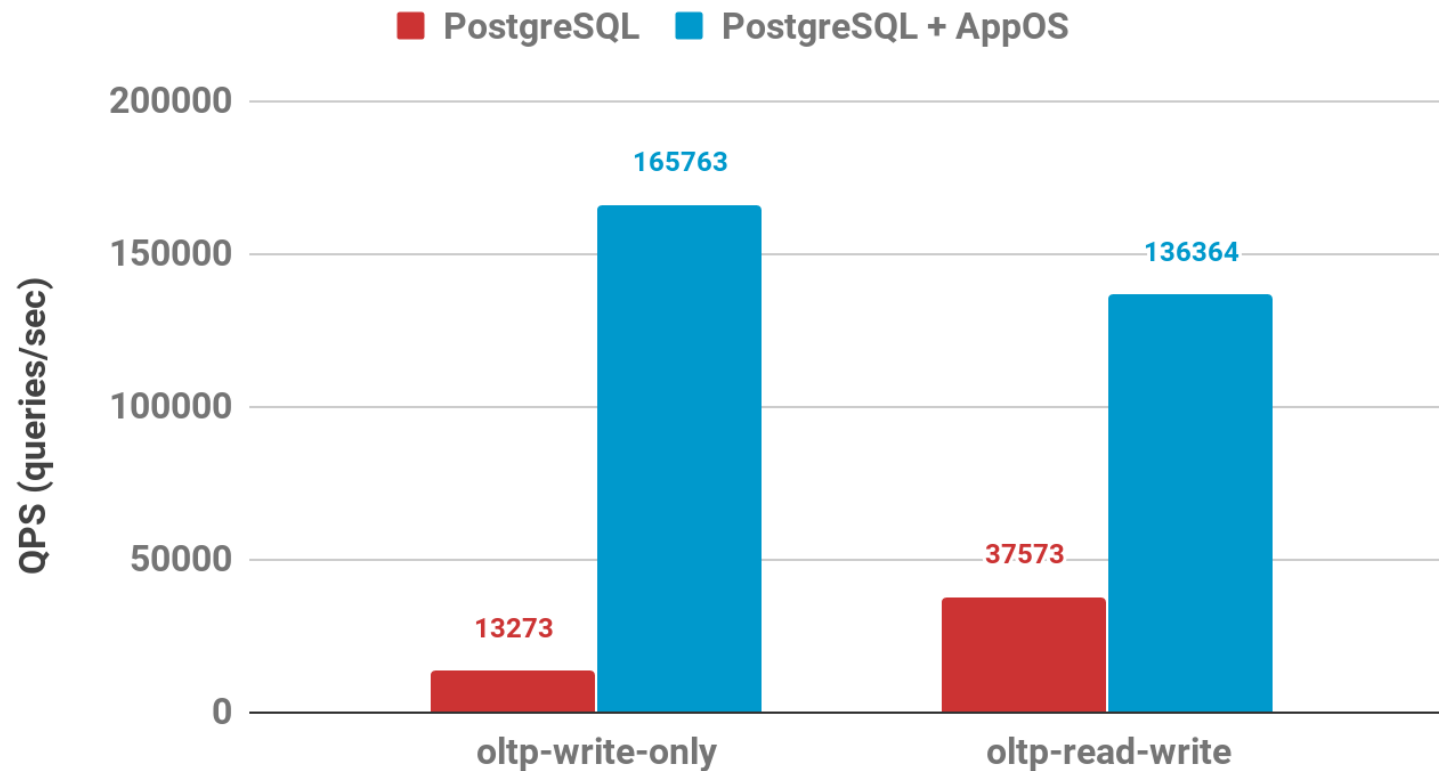
[Building a 16 KB atomic write path from database to block device](#)

What's next

# Use Cases (2)

## Atomic write support

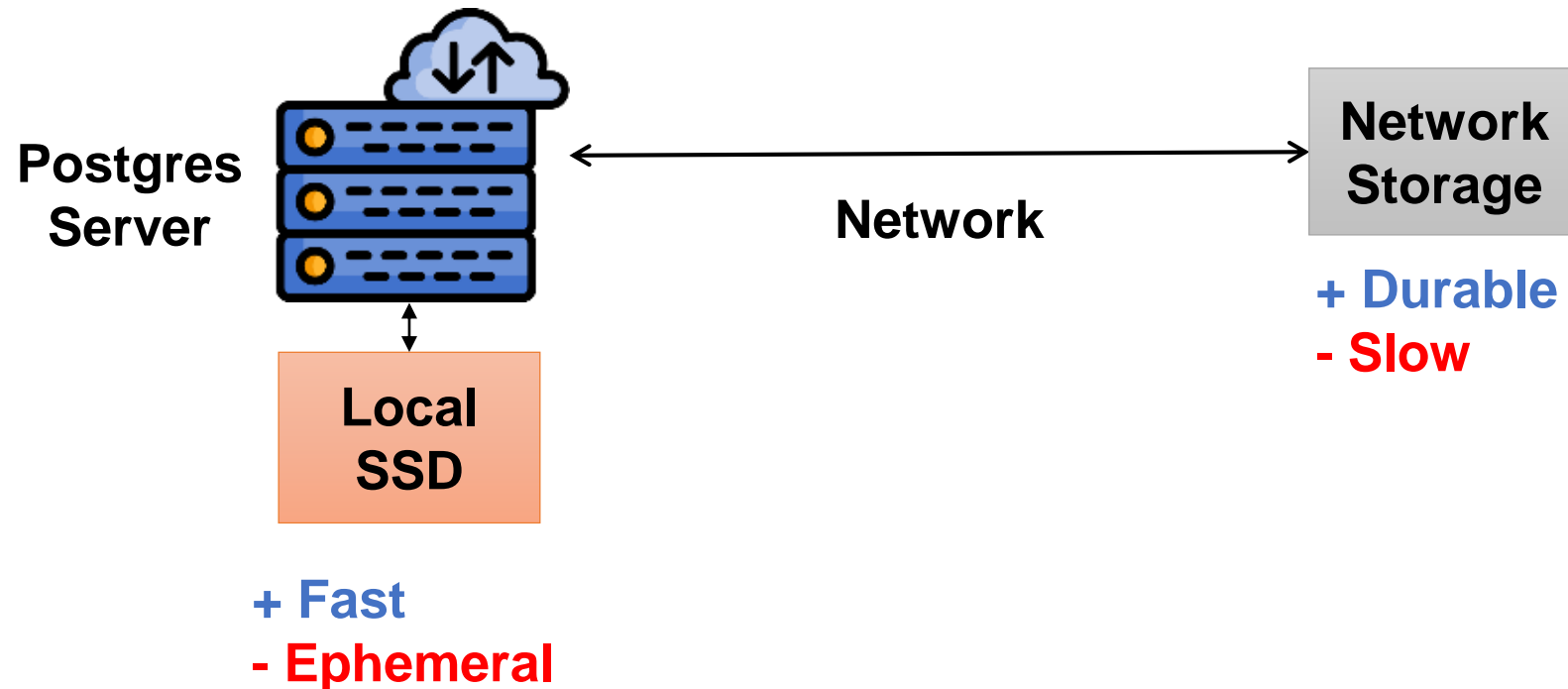
sysbench throughput, 50GiB, 100 clients



**4x ~ 12x**  
with atomic write support  
on Google Cloud

## Use Cases (2)

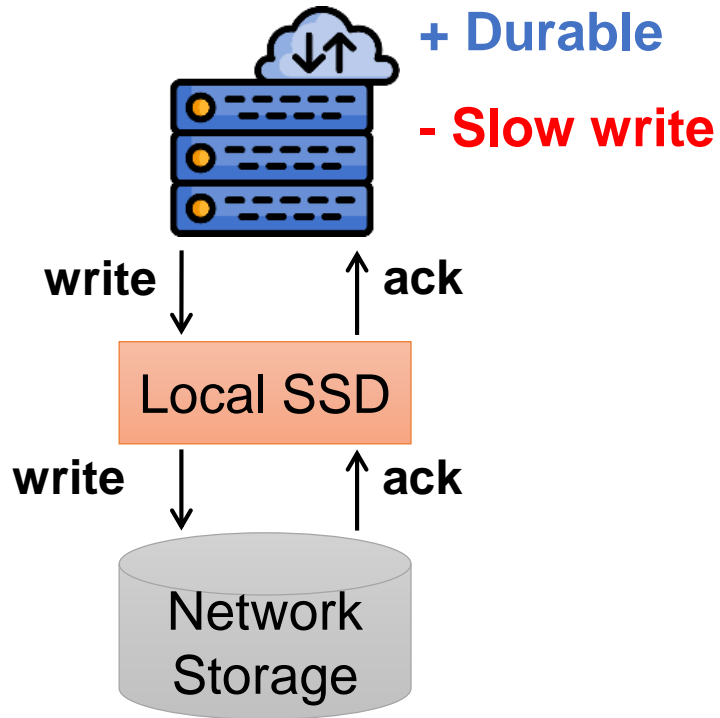
### Local SSD as a memory extension



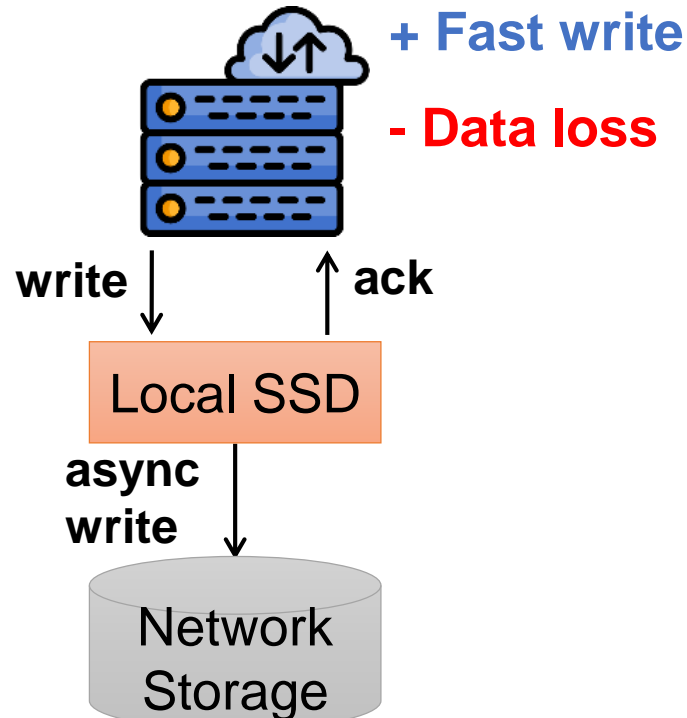
# Use Cases (2)

## Local SSD as a memory extension

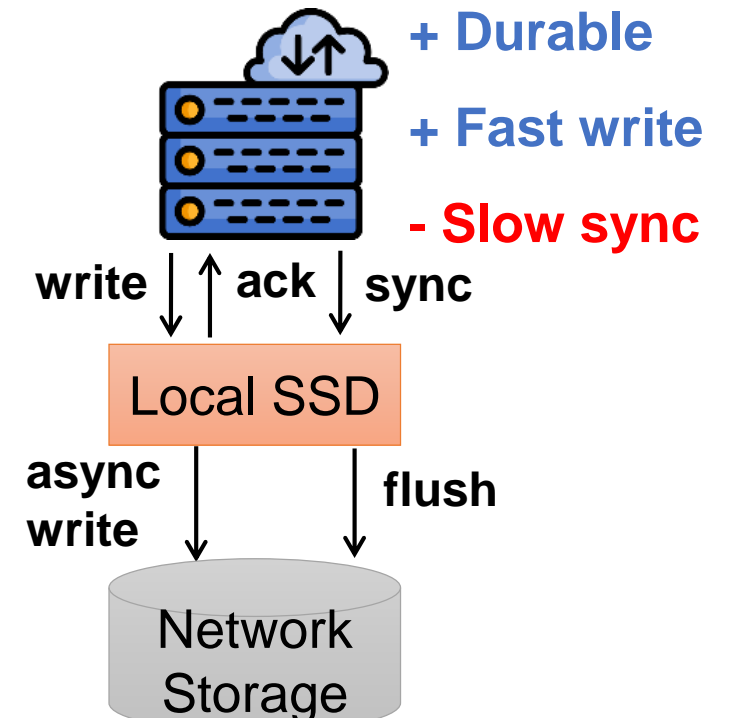
### Write through



### Write back



### Flush on sync



## Use Cases (3)

AppOS can seamlessly work with PostgreSQL-derived DBs

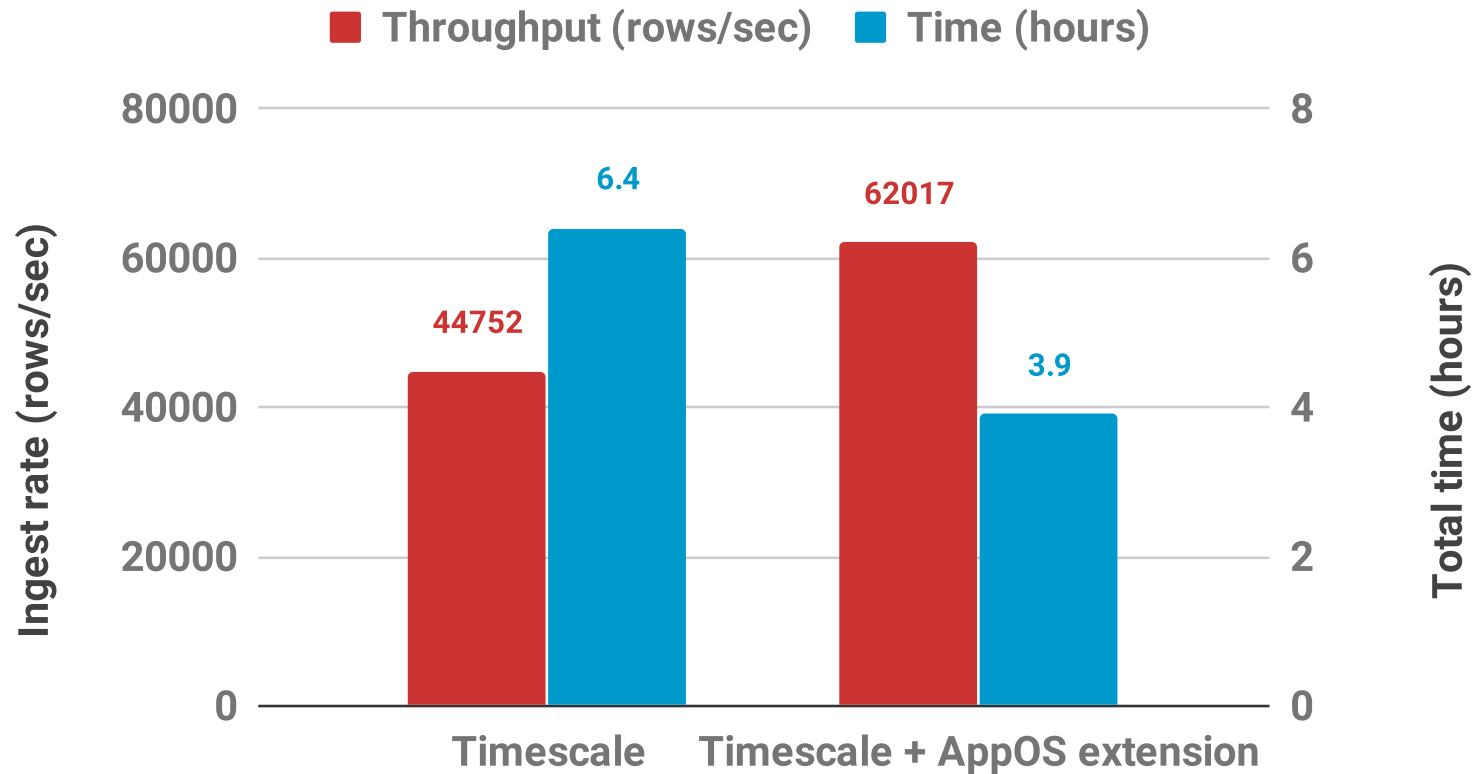




# Use Cases (3)

## TimescaleDB

TSBS Ingest Rate and Time m5.2xlarge with 300 IOPS EBS



<https://blog.timescale.com/blog/timescaledb-vs-6a696248104e>



apposha



<https://apposha.io>  
shawn@apposha.io