facebook

# Hybrid XFS: Supercharging HDDs with SSDs

**Skanda Shamasunder**

Production Engineer (Storage), Facebook

# Who am I?

# Outline

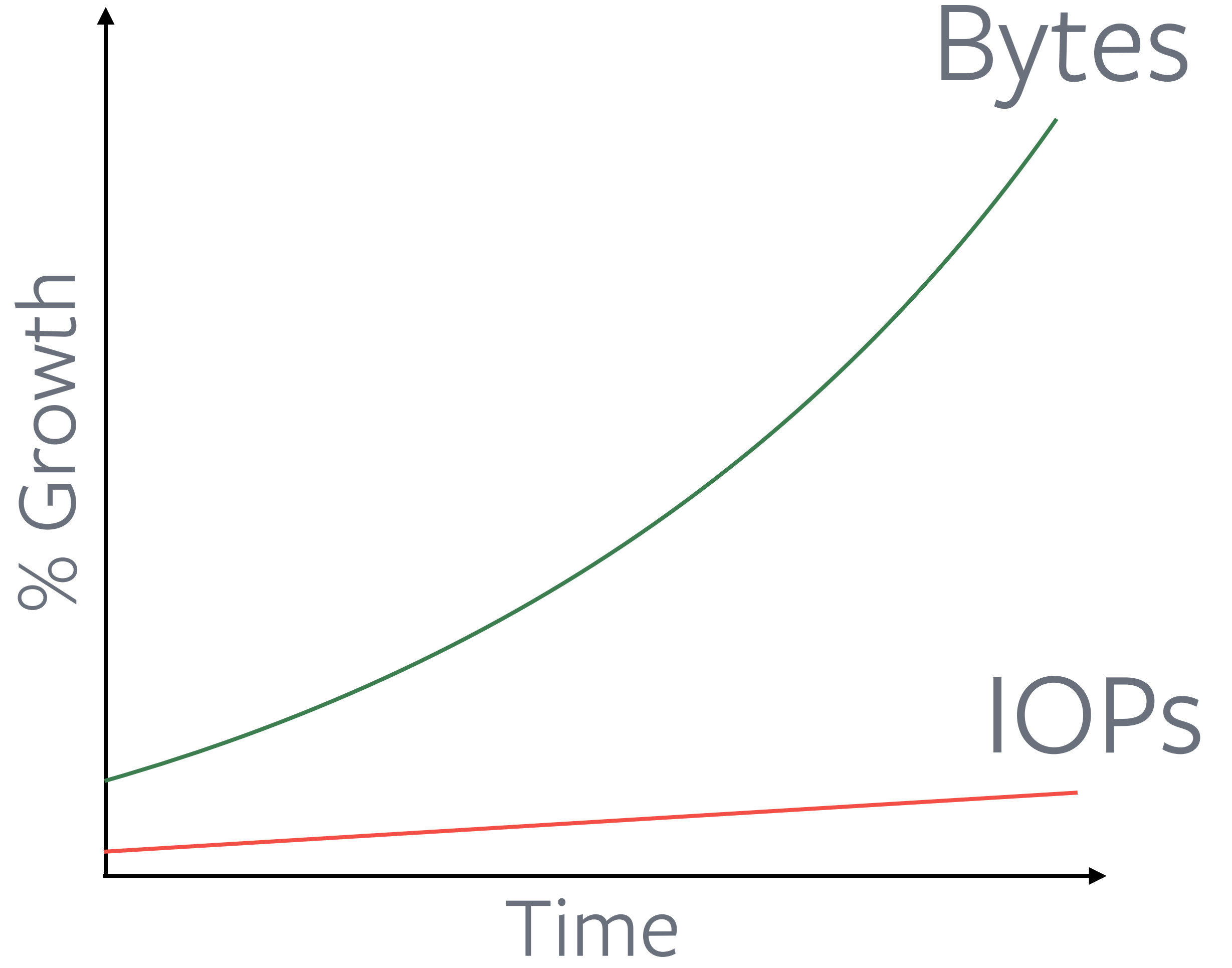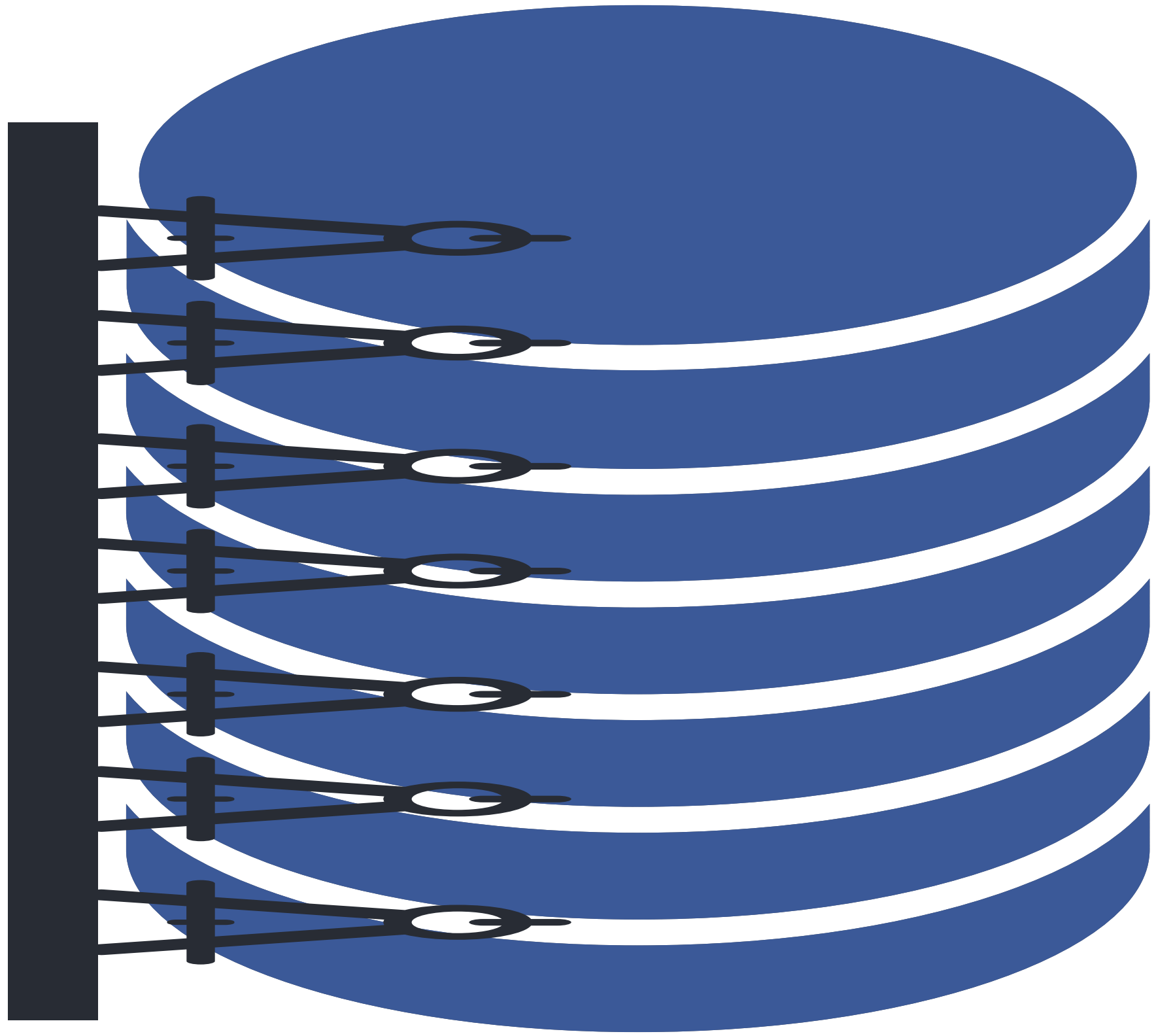# The IO wall
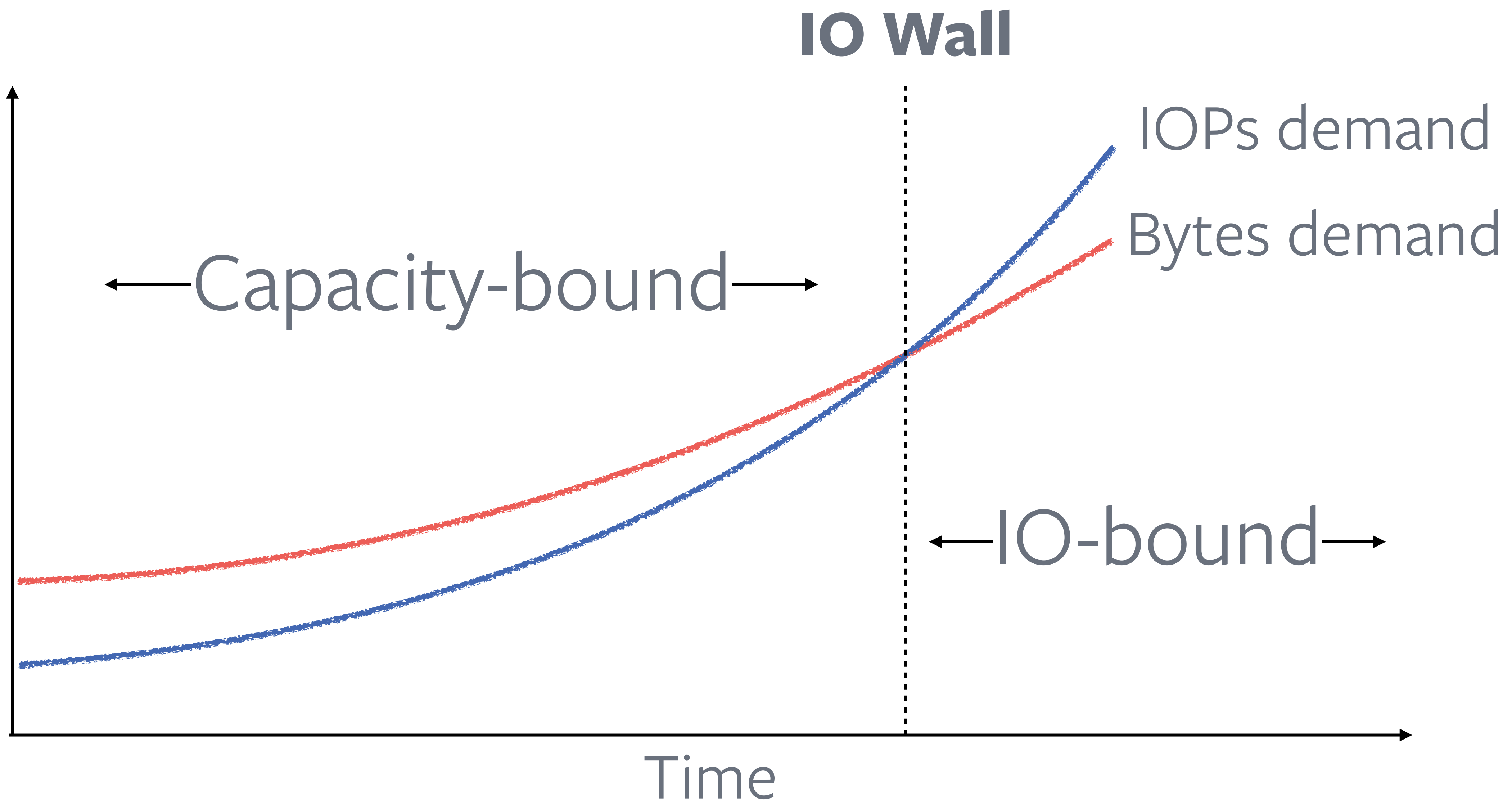
Bytes

IOPs

% Growth

Time

# IO workloads are only getting hotter
Thanks, ML

- AI/ML analytics
- Ephemeral photos
- Live Video
- Offloading Cold Data

**IO Wall**

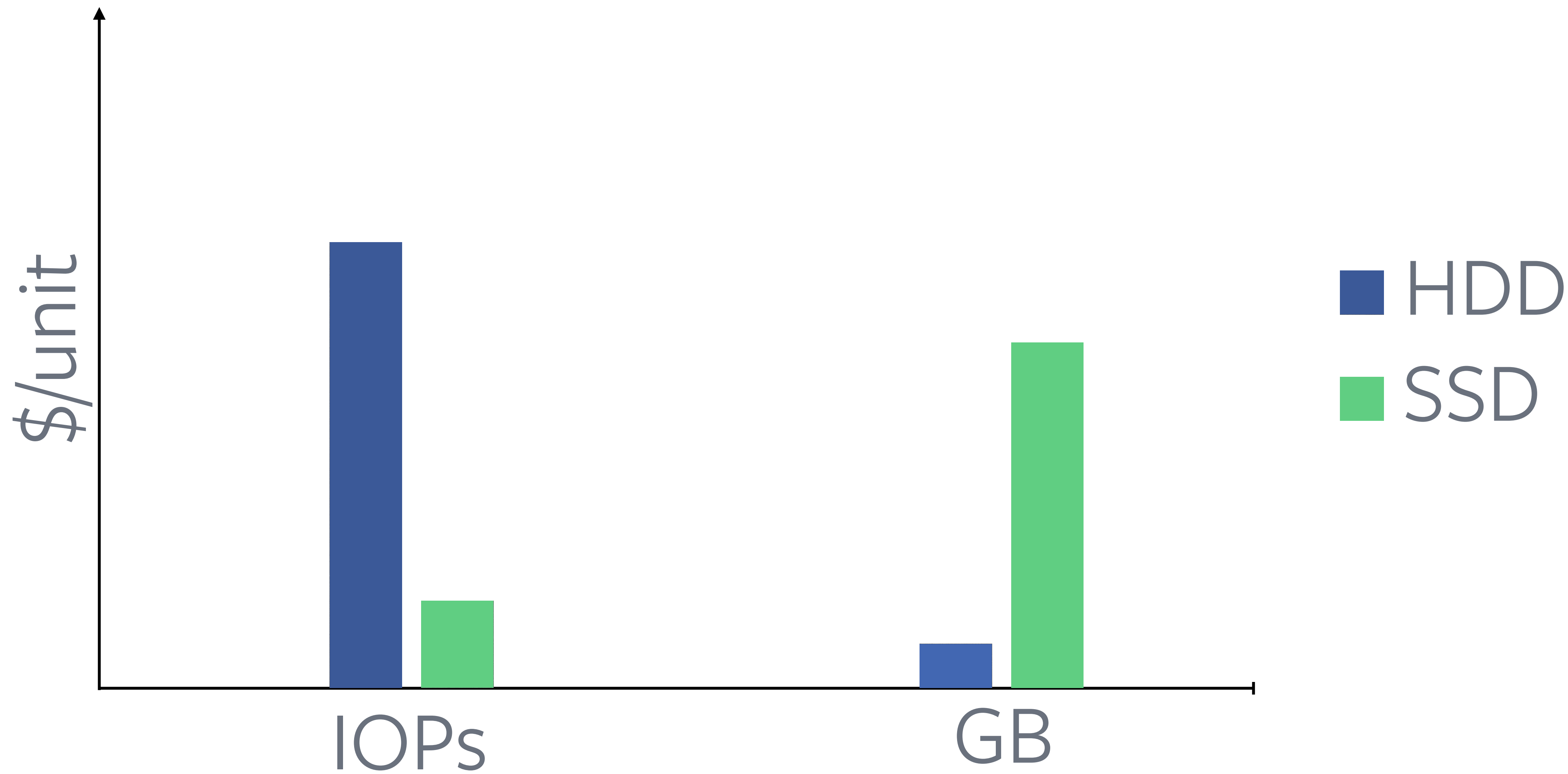IOPs demand

Bytes demand

Capacity-bound

IO-bound

Time

| Presto | Hive | Spark |
| --- | --- | --- |

**Analytics File Formats**

**Caching**

**Distributed File System**

**On-Disk File System (XFS)**

# The Opportunity

```
[rwareing@storage001.dc1 /dev/shm/blktrace] blkparse
      41 N      <---- Un-categorized
    8571 R      <---- Data reads
      14 RM     <---- Metadata reads
    4773 WM     <---- Metadata writes          24%
    6901 WS     <---- (Synchronous) Data writes
```

# Options to reduce metadata IO

- Write less data

- Smaller metadata

- Fewer flushes for writes

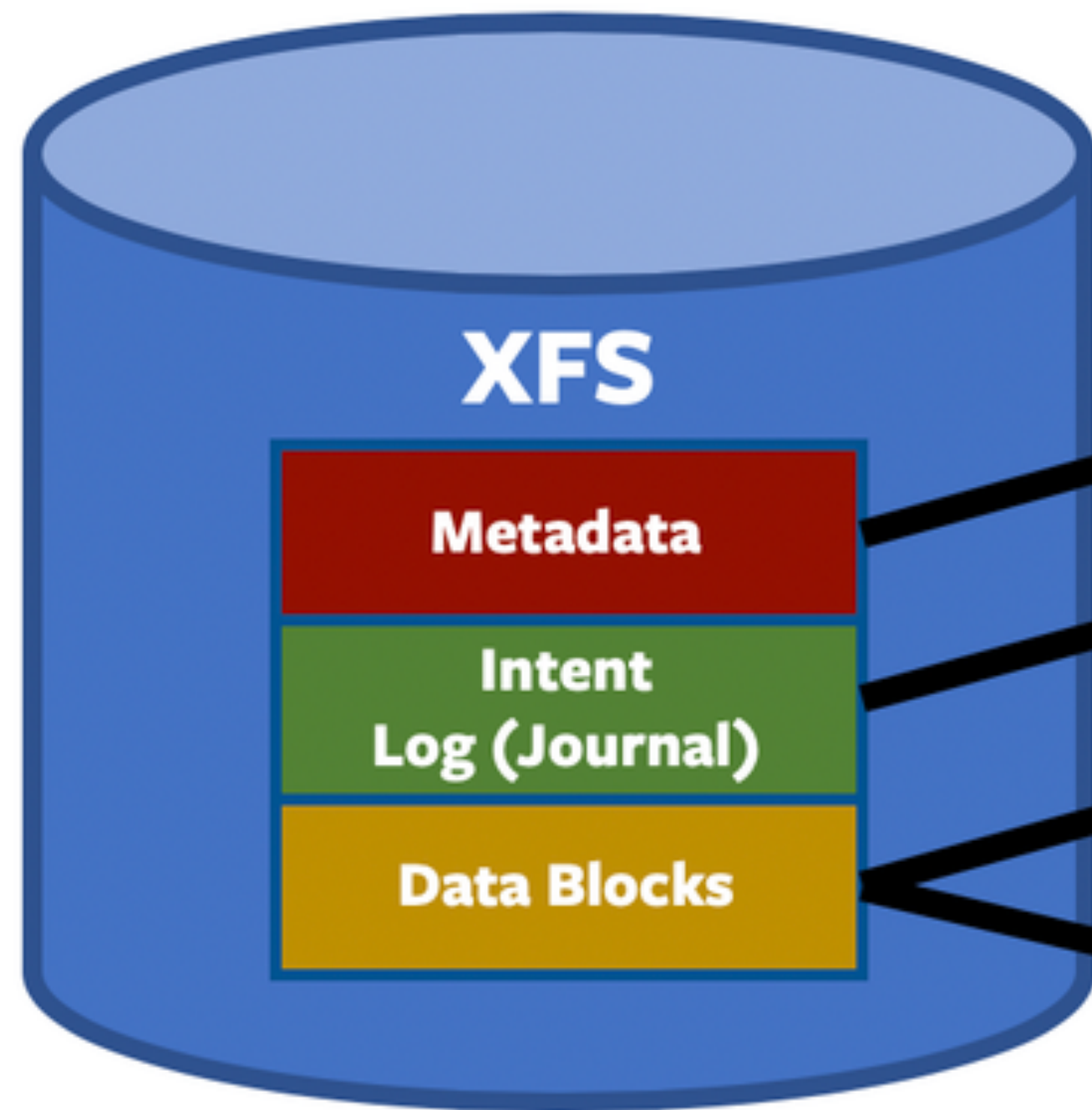- Create a new filesystem
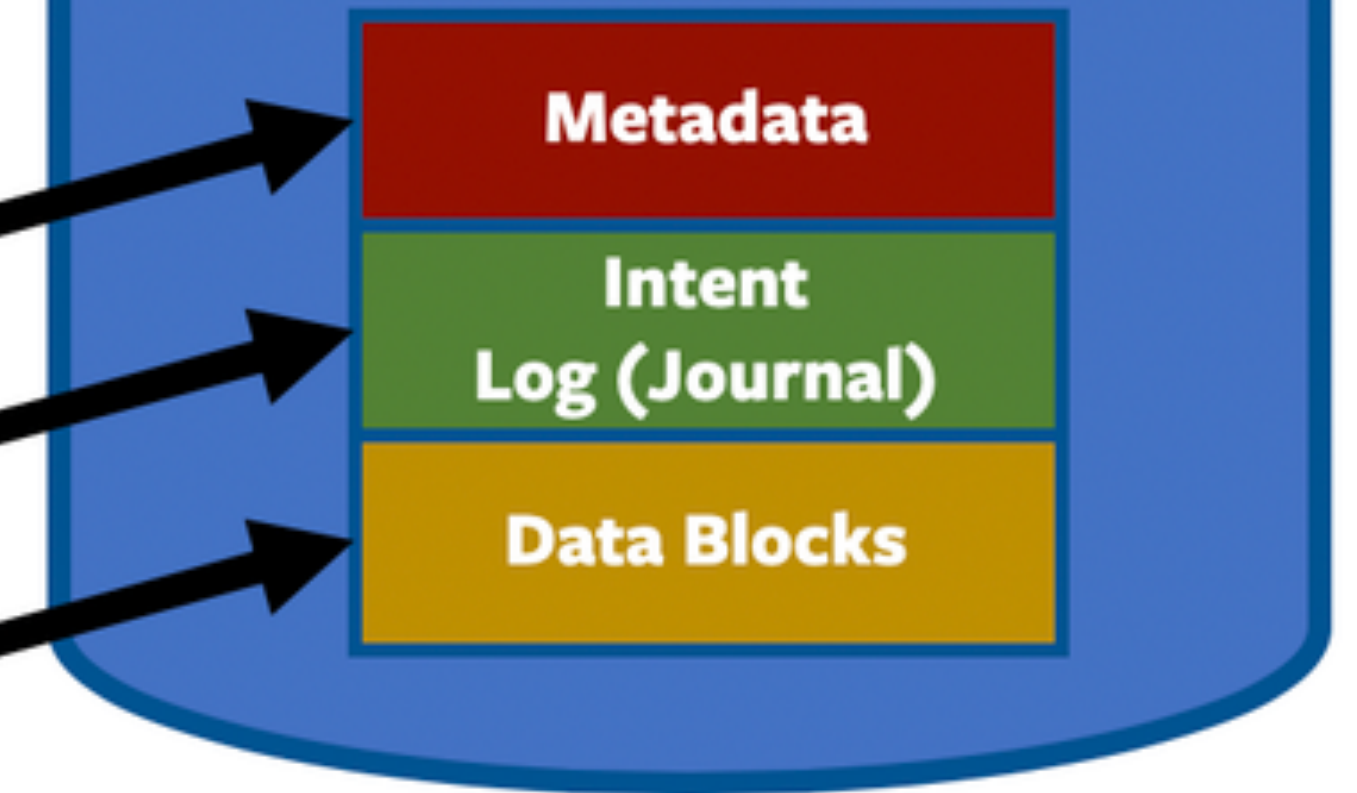
- Put the metadata elsewhere

XFS Realtime mode - a hidden gem

**Traditional XFS**

XFS
- Metadata
- Intent Log (Journal)
- Data Blocks

**XFS w/ Realtime Subvolume**

Data Sub-volume
- Metadata
- Intent Log (Journal)
- Data Blocks

Realtime Sub-volume
- Realtime Data Blocks

image by Richard Wareing, Facebook

# But... we only have one SSD

# Partition the heck out of it

And gave it a name – **Hybrid XFS**

image by Richard Wareing, Facebook

Hybrid XFS vs. Control – Ave. Request Size (FB IO Trace)

1.5MB

1MB

image by Richard Wareing, Facebook

Keep your eye on the Random Writes

Legend: Coalesced Reads · Random Reads · Coalesced Writes · Random Writes

image by Richard Wareing, Facebook

<10 coalesced IOs/sec; more than 50% reduction

Nearly elminated

image by Richard Wareing, Facebook

Hybrid XFS vs. Control – Time Spent in Operation (ms)

15% Reduction

25% Reduction

— Hybrid XFS — Control

image by Richard Wareing, Facebook

Success
Ahead

Image by Gerd Altmann from Pixabay

# The Risks

# 1. SSD Failures
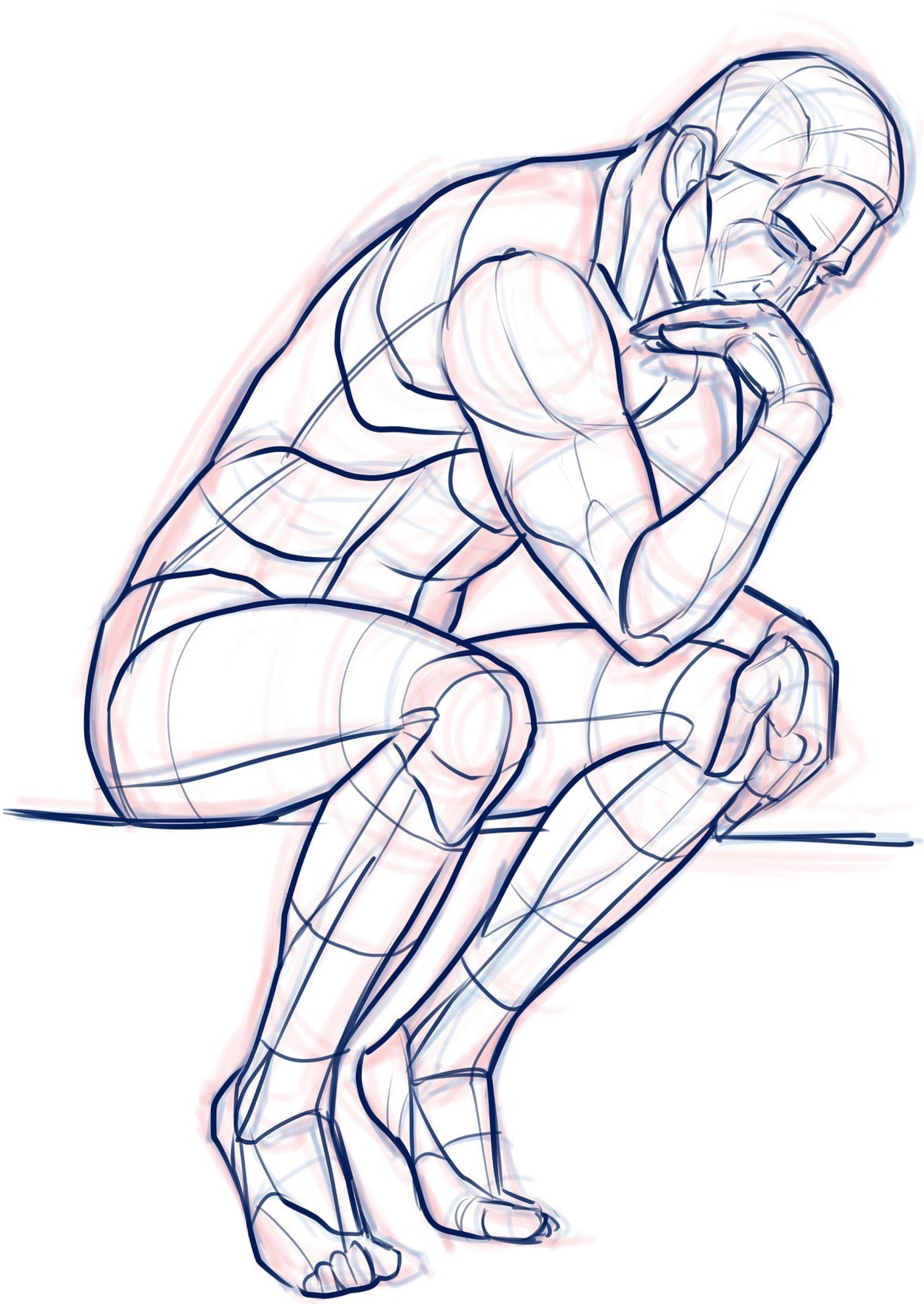
# 2. Endurance

# 3. Hardware changes?

# 4. Operational headache

# Risks.....

Many many risks

- Changes in workload?
- Rollout?
- and...?

# Let the Data guide you

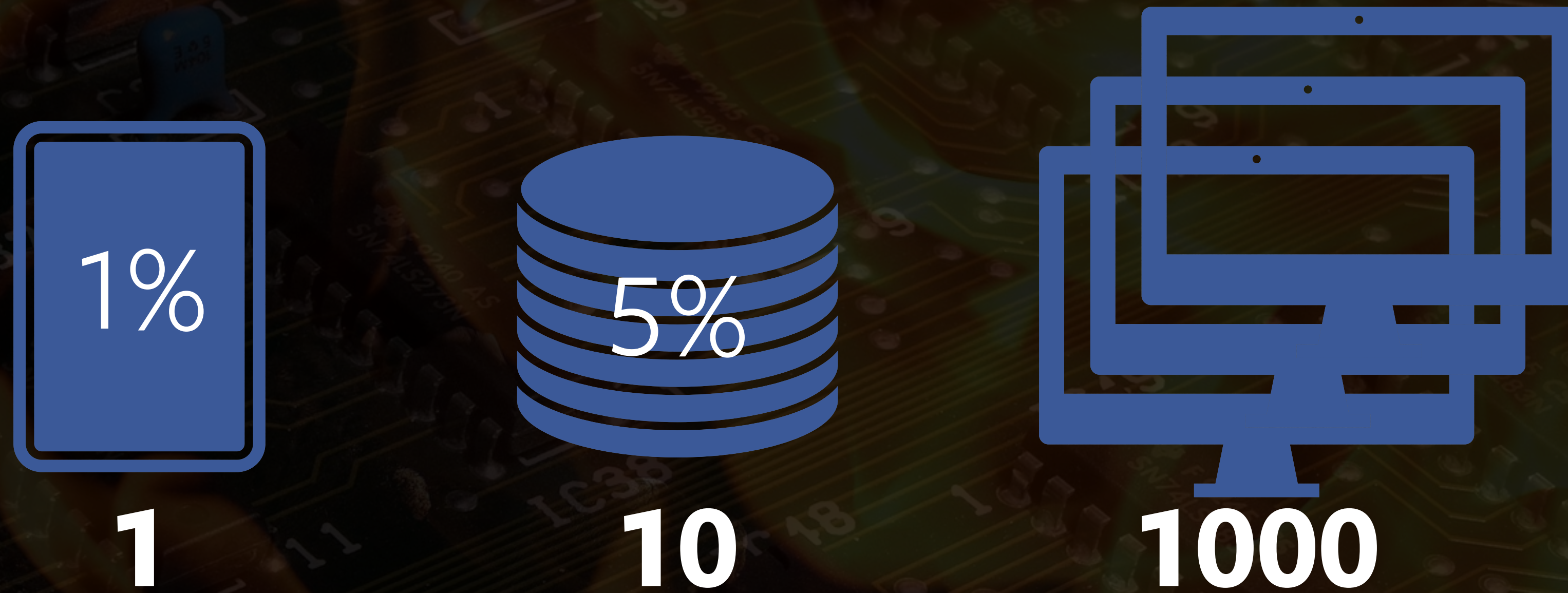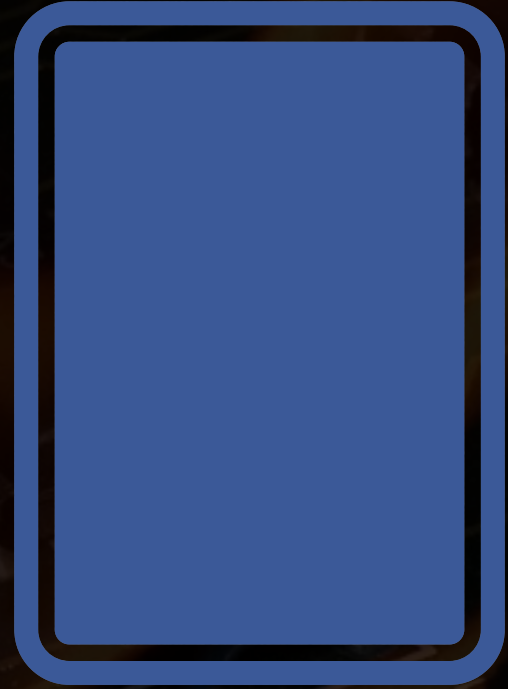Image by skeeze from Pixabay

# The Analysis

# 1. SSD Failures

1%

5%

1

10

1000

# 1. SSD Failures



1000

x 1% = 10 SSDs or
100 HDDs

+

10000

x 5% = 500 HDDs

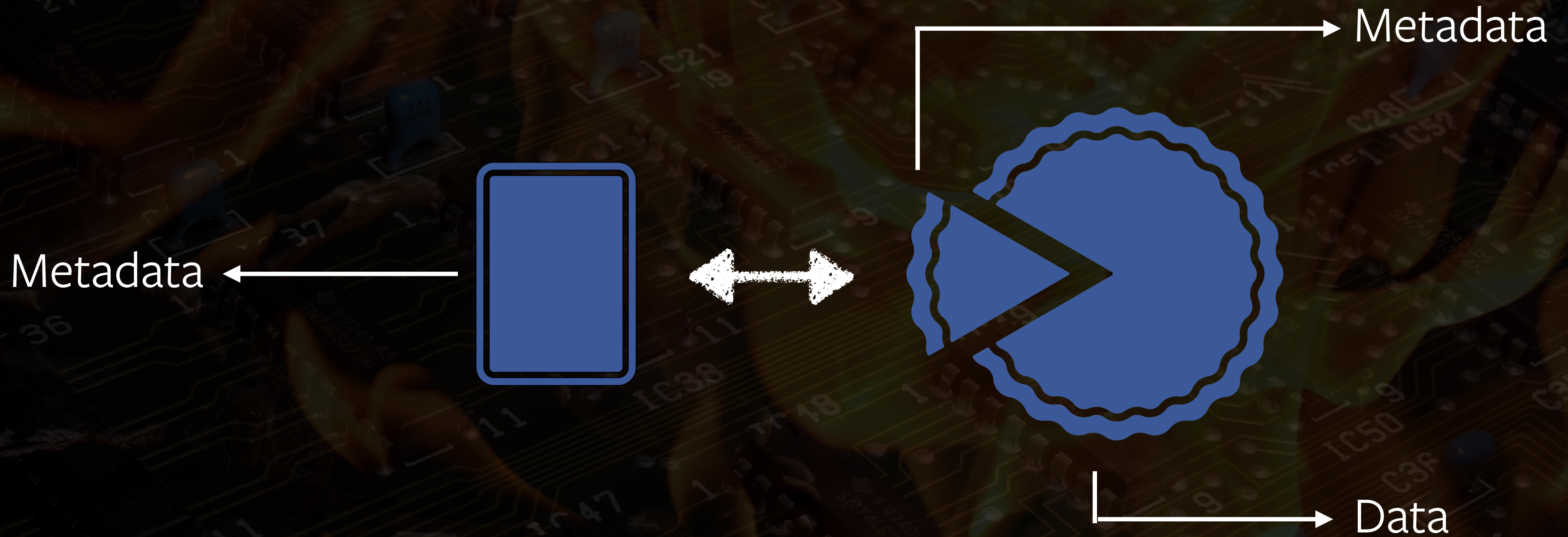= 600 HDDs

1500 ~~HDDs~~

# 1. SSD Failures – en masse!

# 1. SSD Failures

Metadata

Metadata

Data

**Rescue mode**

# 2. Endurance



*DWPD = Drive Writes Per Day

# 3. Will there always be SSDs?

# 4. Operational headache

- Patched XFS statfs call

- Collect stats for both metadata and data

- Make systemd wait for BOTH devices

# Rolling It Out

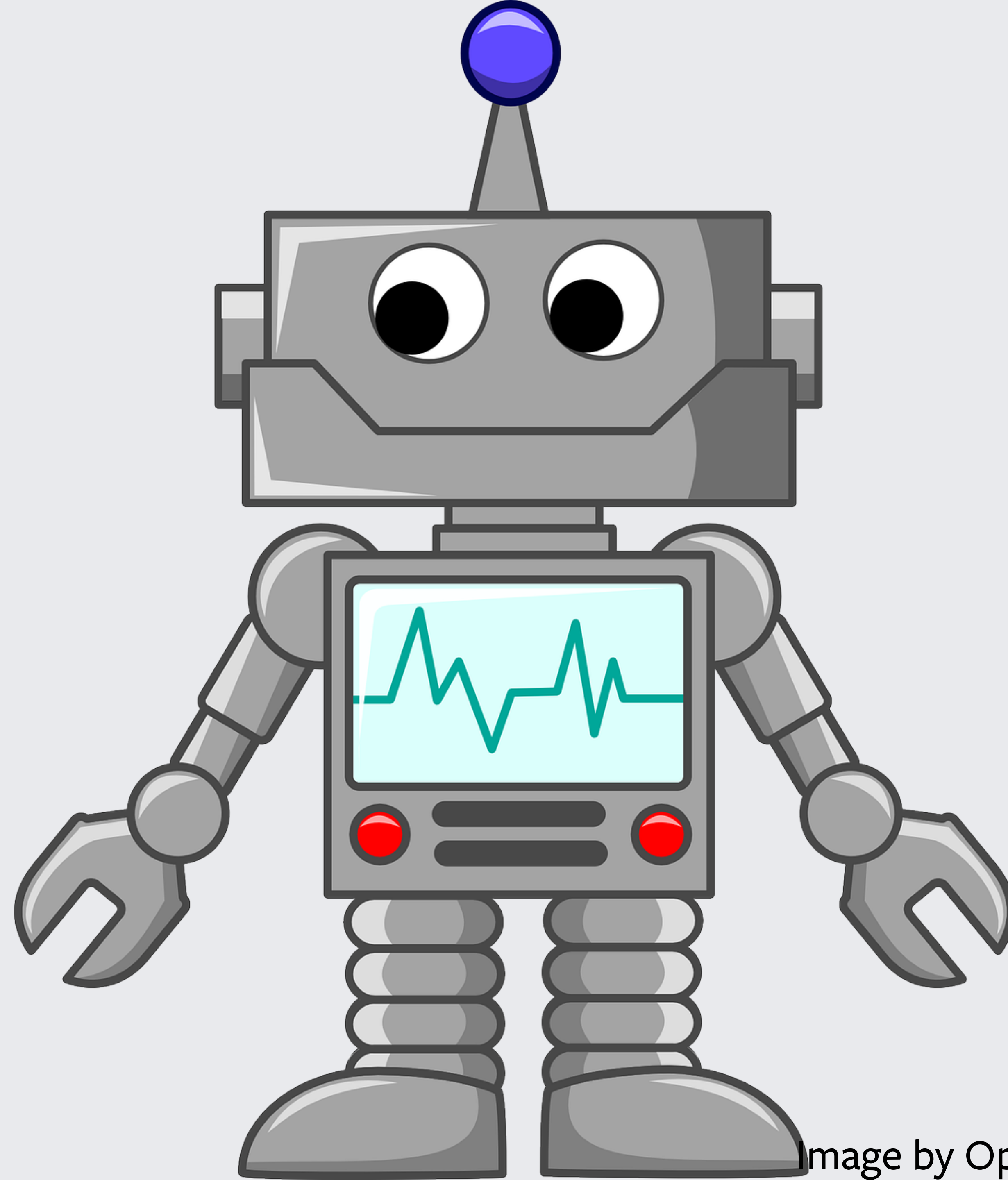# How do you roll out a destructive change to tens of thousands of hosts?

PLEASE BE CAREFUL!

PREVENT WILDFIRES
NEW JERSEY FOREST FIRE SERVICE

# Look out for...

- Impact to:
  - Durability
  - Capacity
  - Performance
- Fallout from failures
- Automation flying blind

# Success!

# Lessons

"Hard problems can have simple solutions"

"Gut feelings can be wrong"

"Data wins arguments"

"Better safe than sorry"

# Thank You!

facebook