# AWS re:Invent

**DECEMBER 2 – 6, 2024 | LAS VEGAS, NV**
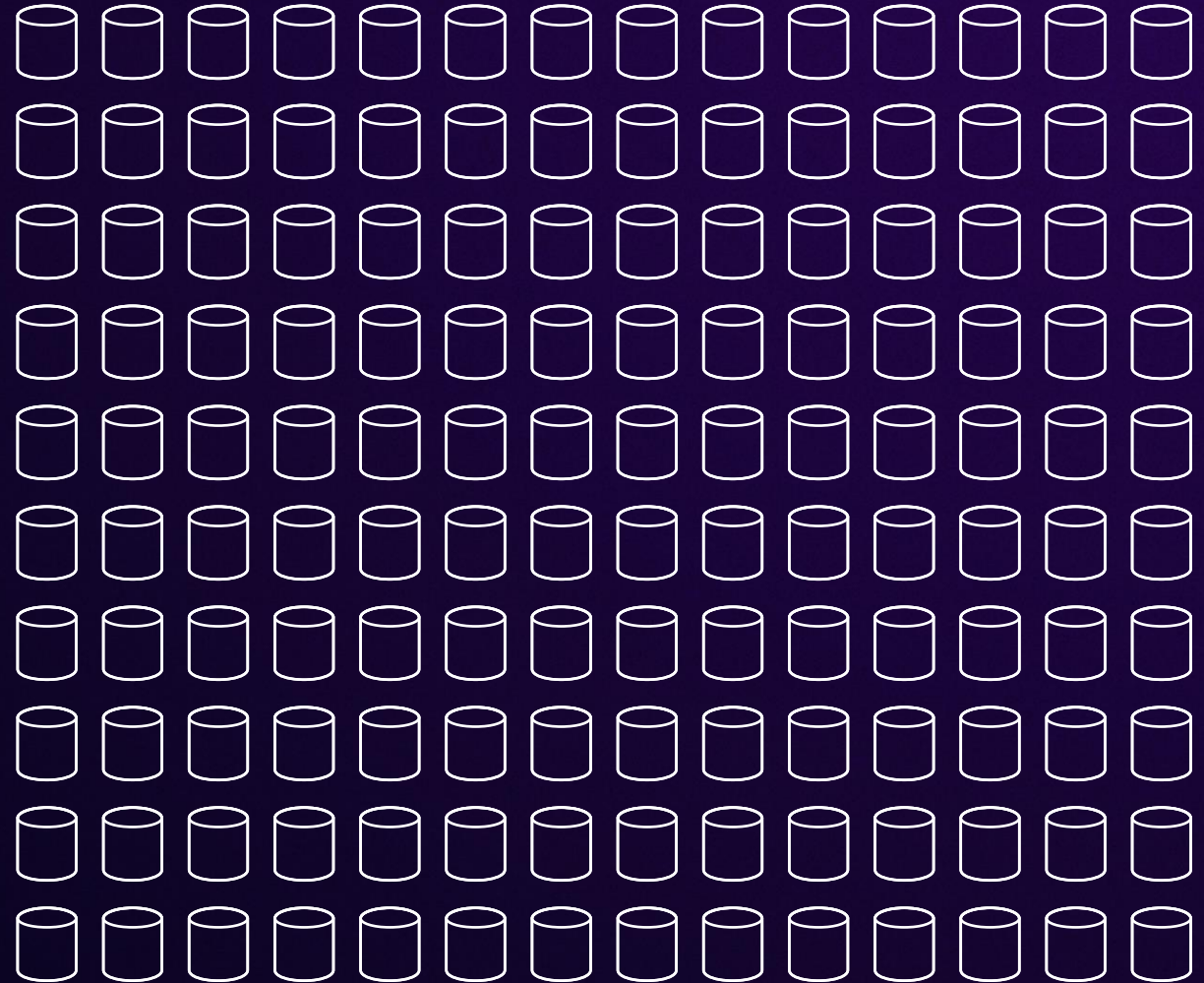
# Using scale to our advantage ... and yours

# Amazon S3 at scale

Over 400 trillion objects

Over a quadrillion requests per year

Over 200 billion events daily

Over 1 PB/s transferred at peak

**01**     Physics of data

**02**     Designing decorrelated systems

**03**     Engineering for failure is engineering for velocity

01 Physics of data

02 Designing decorrelated systems

03 Engineering for failure is engineering for velocity
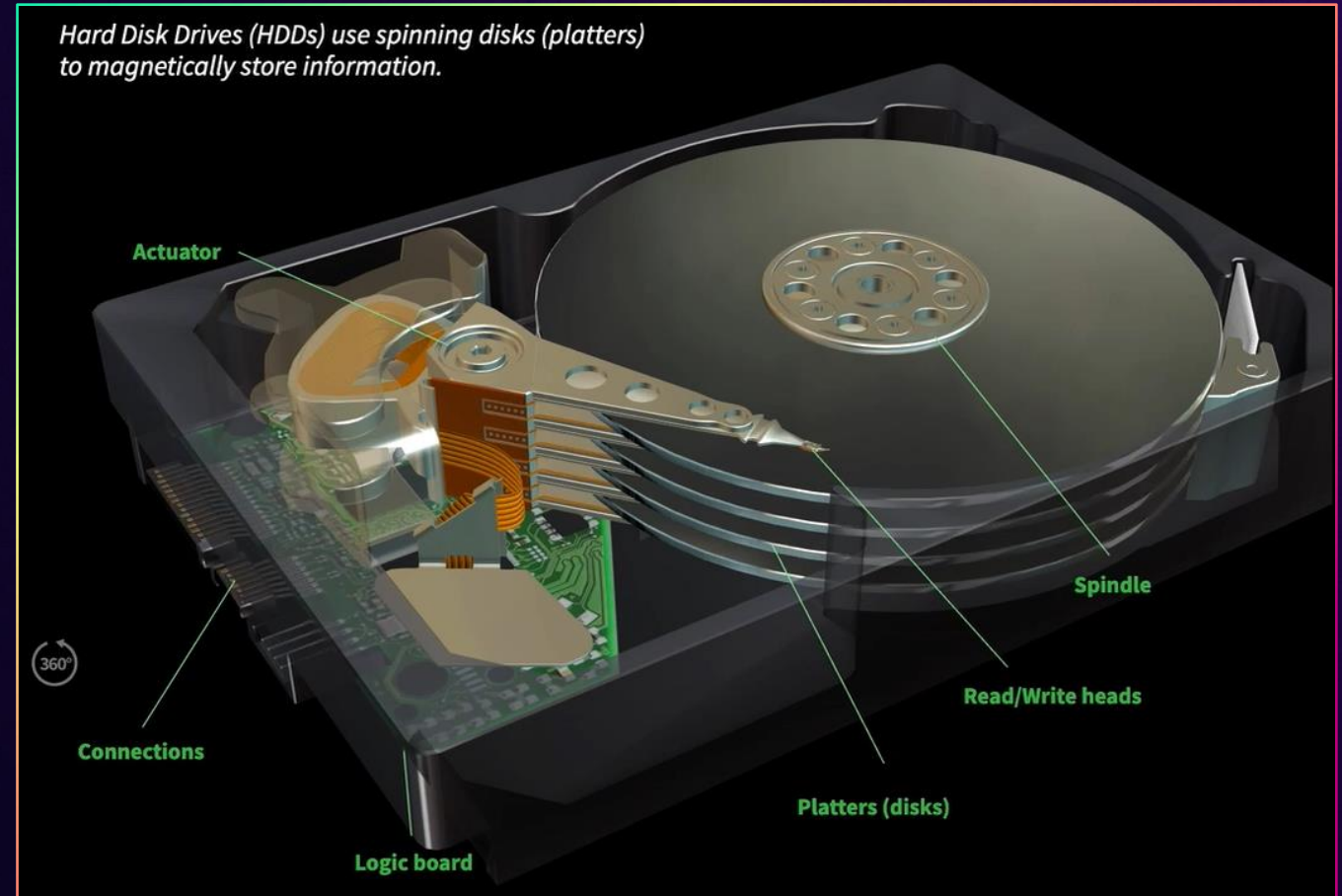
# The hardware, HDDs

**Rotational latency
Seek time**
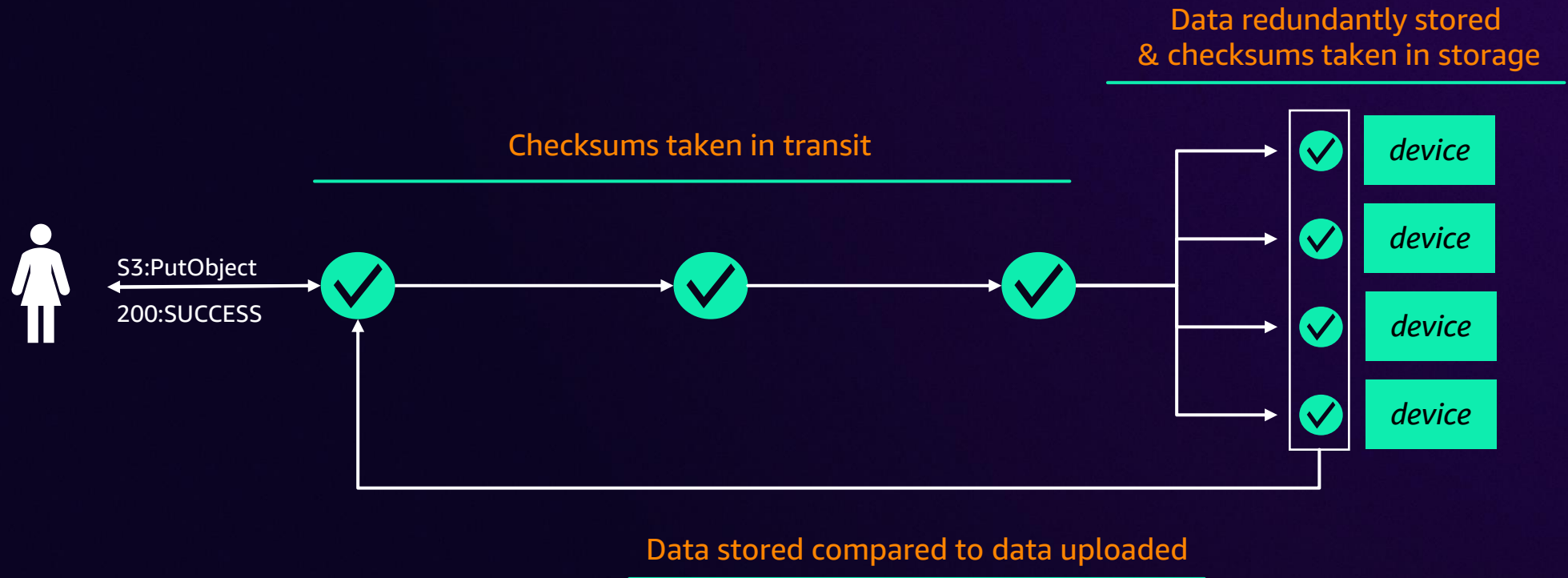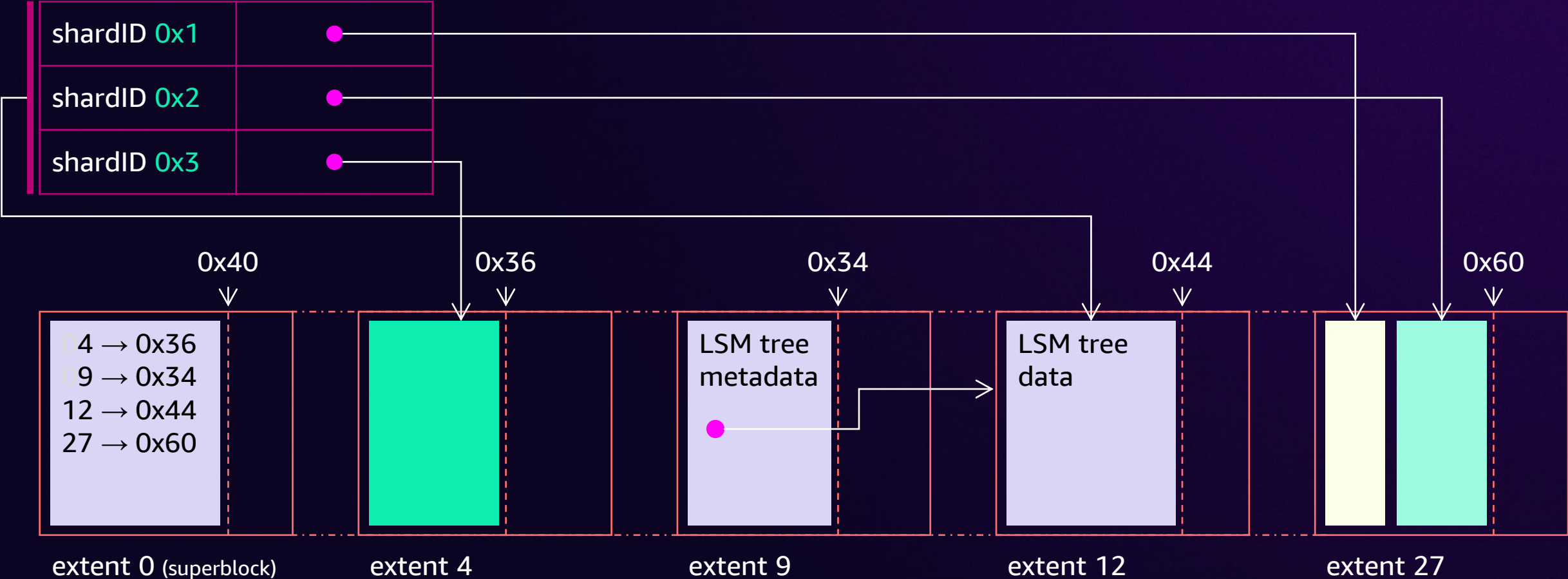


Hard Disk Drives (HDDs) use spinning disks (platters)
to magnetically store information.

Actuator

Spindle

360°

Read/Write heads

Connections

Platters (disks)

Logic board

*Image by Animagraffs*
*https://animagraffs.com/hard-disk-drive/*

# The software, Replication



Data redundantly stored
& checksums taken in storage

Checksums taken in transit

device

device

device

device

S3:PutObject

200:SUCCESS

Data stored compared to data uploaded

# The software, ShardStore



LSM tree

| shardID 0x1 | ● |
|---|---|
| shardID 0x2 | ● |
| shardID 0x3 | ● |

0x40

0x36

0x34

0x44

0x60

4 → 0x36
9 → 0x34
12 → 0x44
27 → 0x60

LSM tree metadata
●

LSM tree data

extent 0 (superblock)　　extent 4　　extent 9　　extent 12　　extent 27

# Individual workloads

- Individual workloads are bursty

- Provisioning for peak

# Individual workloads

- Individual workloads are bursty

- Provisioning for peak

- Customer Example: FINRA

# The physics of data

- Example:
  - 1 PB of data
  - 1 MB per object
  - Accessed in a single hour

# The physics of data

- Rotation = 4 ms, average

- Seek = 4 ms, average

- 0.5 MB Transfer = 2 ms, average

- Total: 10 ms per read  (100 reads/second)

- 50 MB per second at 0.5 MB/read
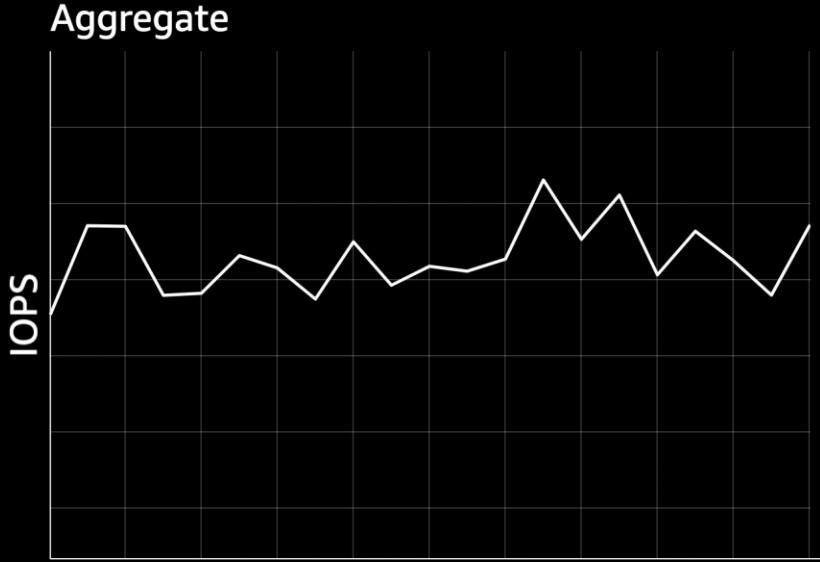
# The physics of data
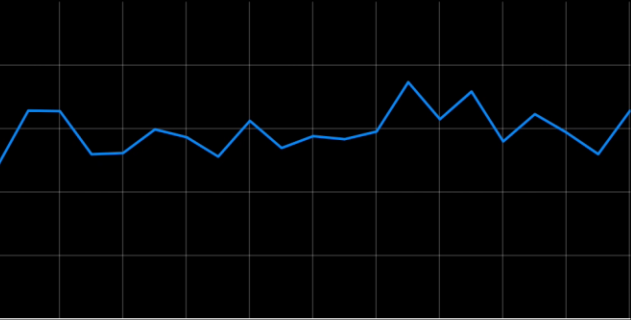
- 1 PB of data

- 275 GB per second

# The physics of data

- 1 PB of data

  - 50 drives at 20 TB per drive

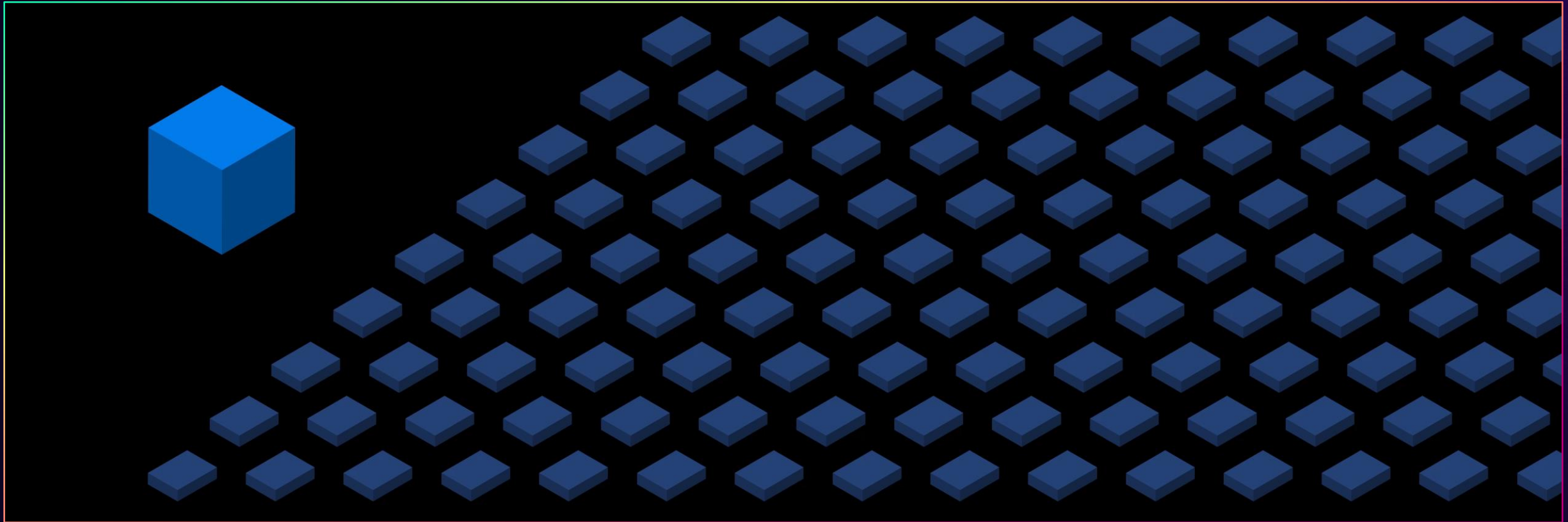- 275 GB per second

  - 5500 drives at 50 MB per second

# The physics of data

- 1 PB of data
  - 50 drives at 20 TB per drive
- 275 GB per second
  - 5500 drives at 50 MB per second
- 100x difference to support bursts!

# Effect of aggregating decorrelated workloads on net system load



Aggregate

IOPS

# Spread shards across a large number of of diverse disks

# Thermodynamics: Balancing the aggregates

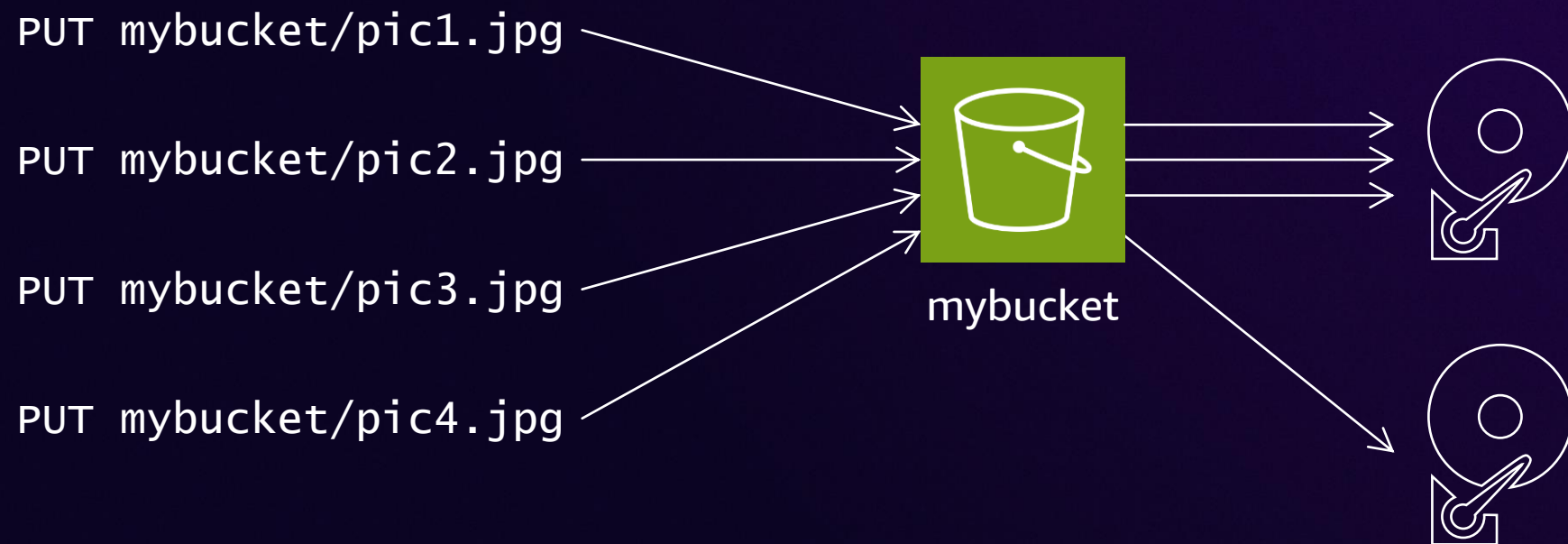# Thermodynamics: Balancing the aggregates

# In-practice: Expanding capacity
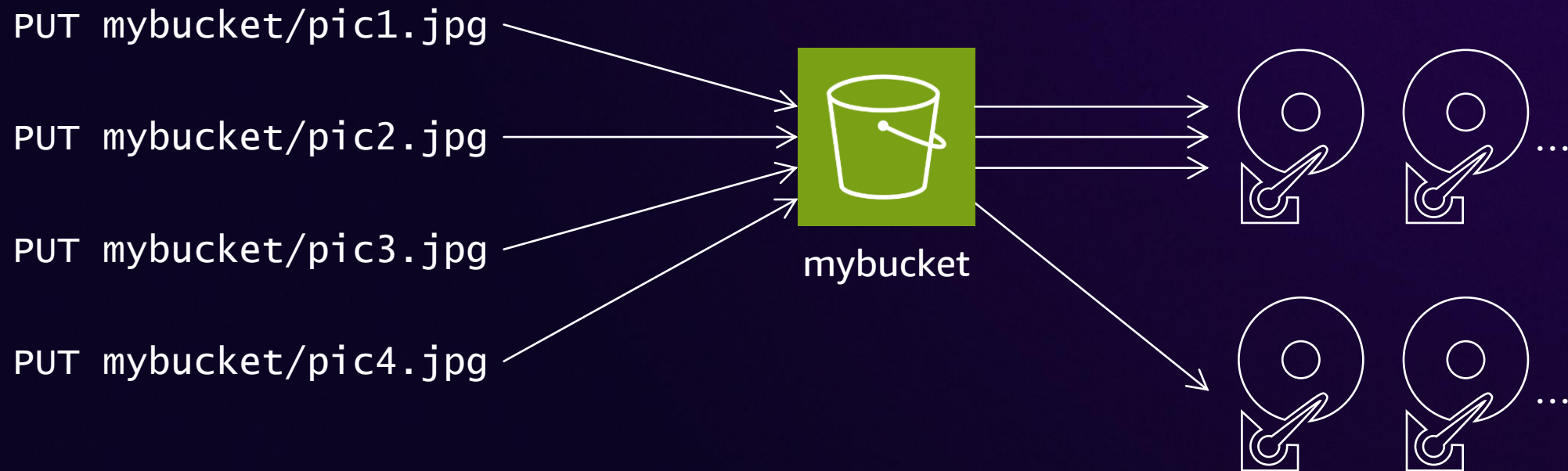
# In-practice: Expanding capacity
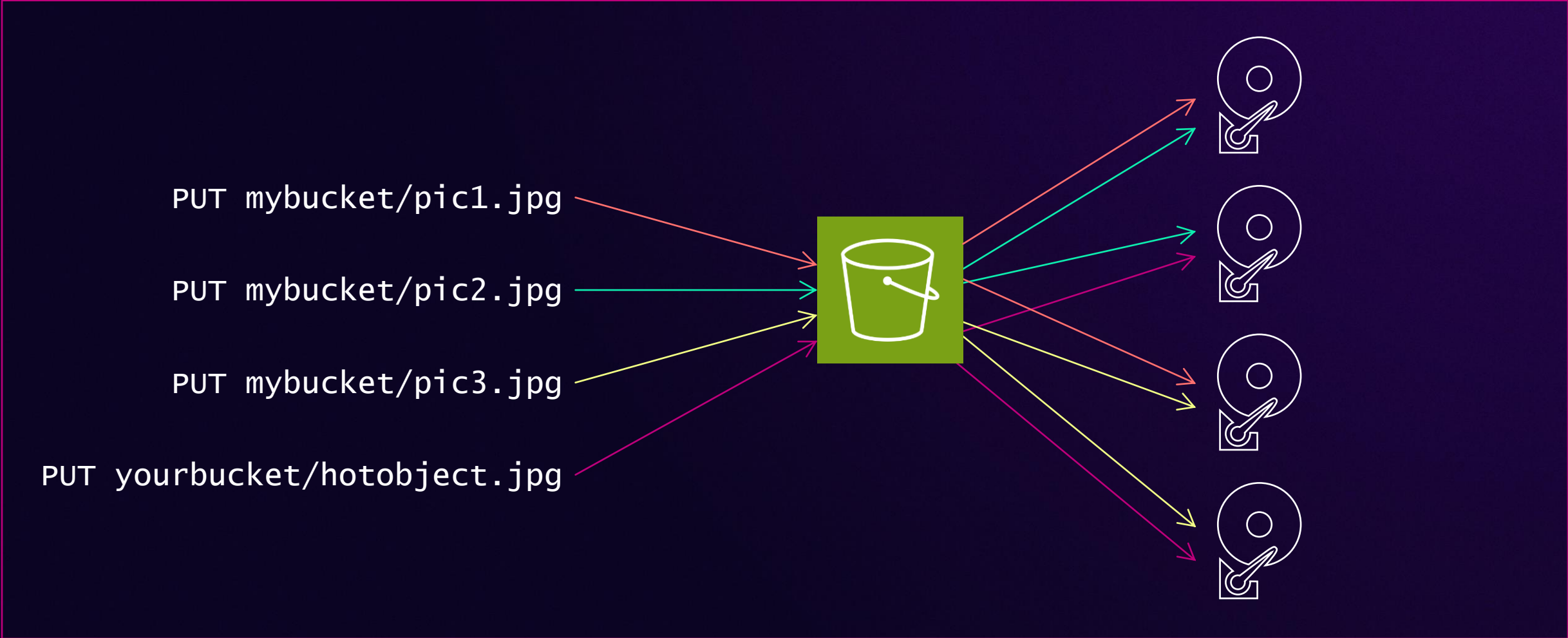
# In-practice: Expanding capacity



EXISTING RACKS

NEW RACK

# Assigning buckets to their storage

PUT mybucket/pic1.jpg

PUT mybucket/pic2.jpg

**mybucket**

PUT mybucket/pic3.jpg

PUT mybucket/pic4.jpg

# Assigning buckets to their storage



PUT mybucket/pic1.jpg

PUT mybucket/pic2.jpg

PUT mybucket/pic3.jpg

PUT mybucket/pic4.jpg

mybucket

# Shuffle sharding



PUT mybucket/pic1.jpg
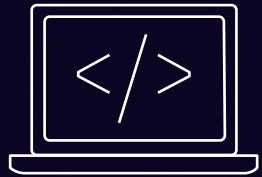
PUT mybucket/pic2.jpg

PUT mybucket/pic3.jpg

PUT yourbucket/hotobject.jpg

**"Elastic" means any S3 customer should be able to use every drive in our fleet on demand – so long as they don't interfere with each other**

# And it's not just drives that shuffle shard

mybucket.s3.amazonaws.com

```
Non-authoritative answer:
mybucket.s3.amazonaws.com
canonical name = s3-us-west-2-
w.amazonaws.com.

Name:s3-us-west-2-w.amazonaws.com
Address:52.92.195.9
Name:s3-us-west-2-w.amazonaws.com
Address:52.92.131.57
Name:s3-us-west-2-w.amazonaws.com
Address:52.218.236.243
...
```
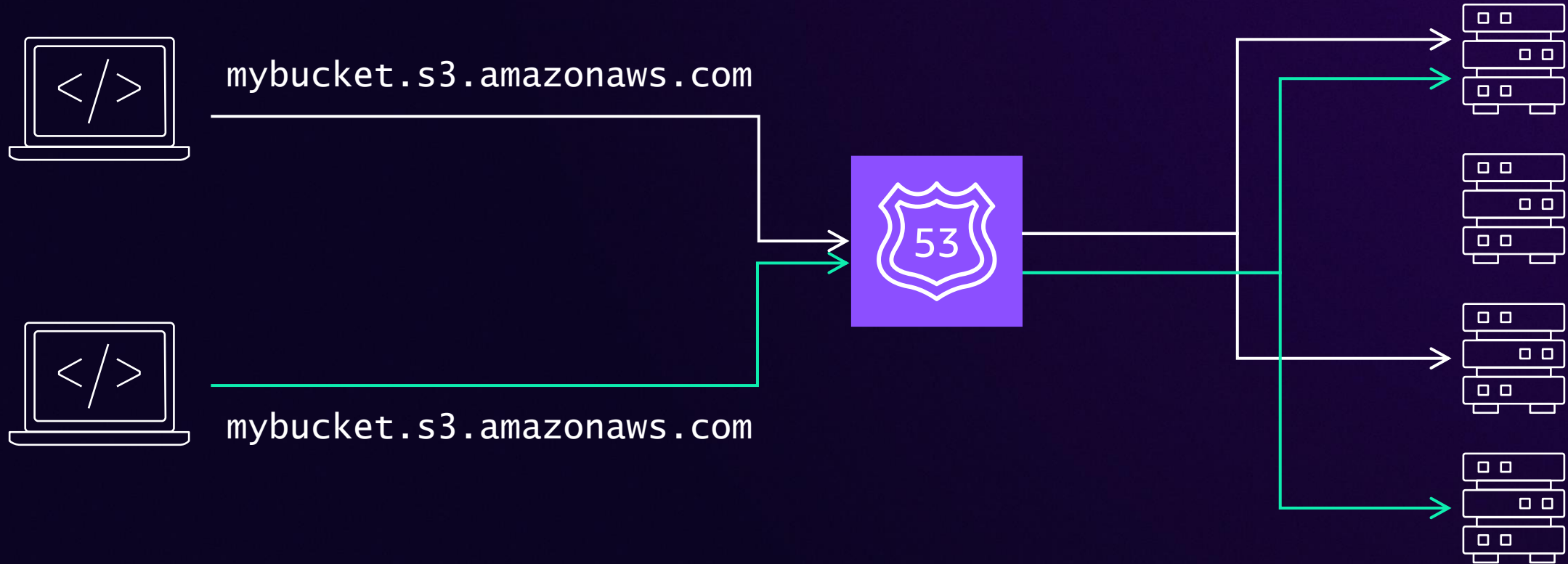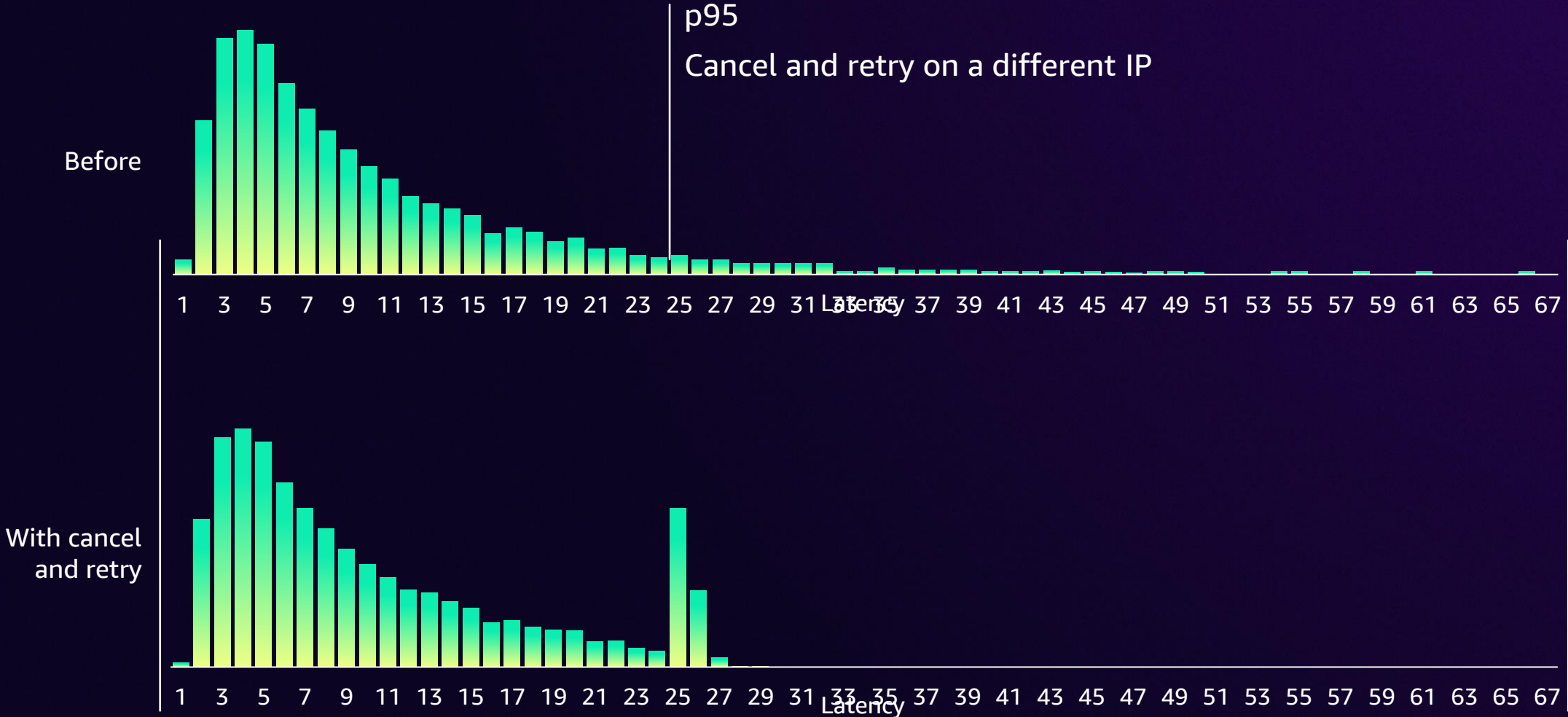
# And it's not just drives that shuffle shard

# Shuffle sharding for fault tolerance



`mybucket.s3.amazonaws.com`

# Shuffle sharding in AWS Common Runtime (CRT)

p95

Cancel and retry on a different IP

Before

With cancel and retry

Latency

Latency

# Placing data for shuffle sharding
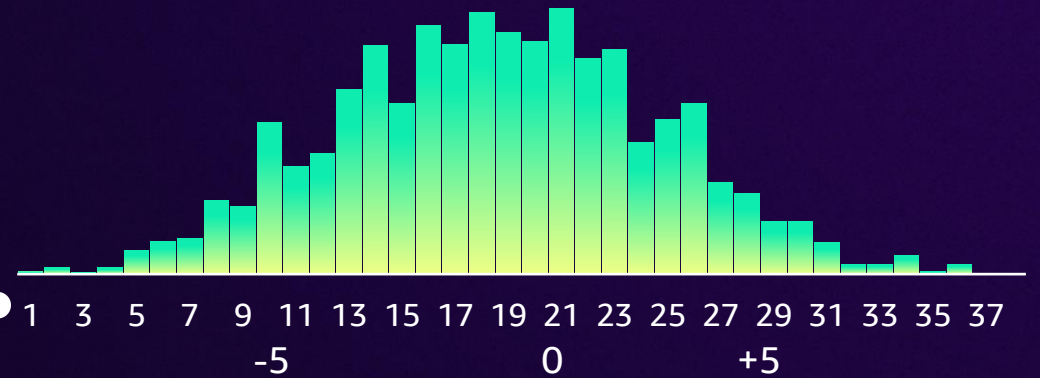
PUT `mybucket/pic1.jpg` ⟶

mybucket

# Placing data for shuffle sharding
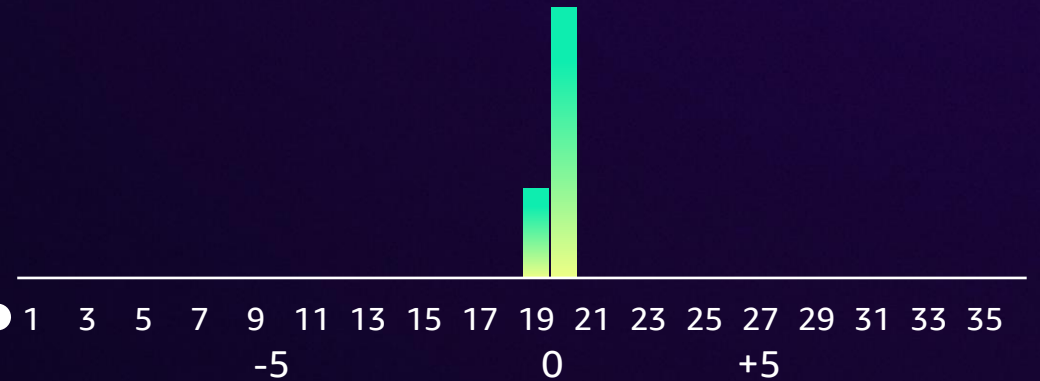


Choose a random drive

PUT mybucket/pic1.jpg → mybucket

"The power of two random choices"
Look at just two random drives and choose the better one

# Engineering for decorrelation

- **Shuffle sharding** decorrelates customers and workloads

- It also enables **scale** and **fault tolerance**

- Use **two random choices** to balance shards



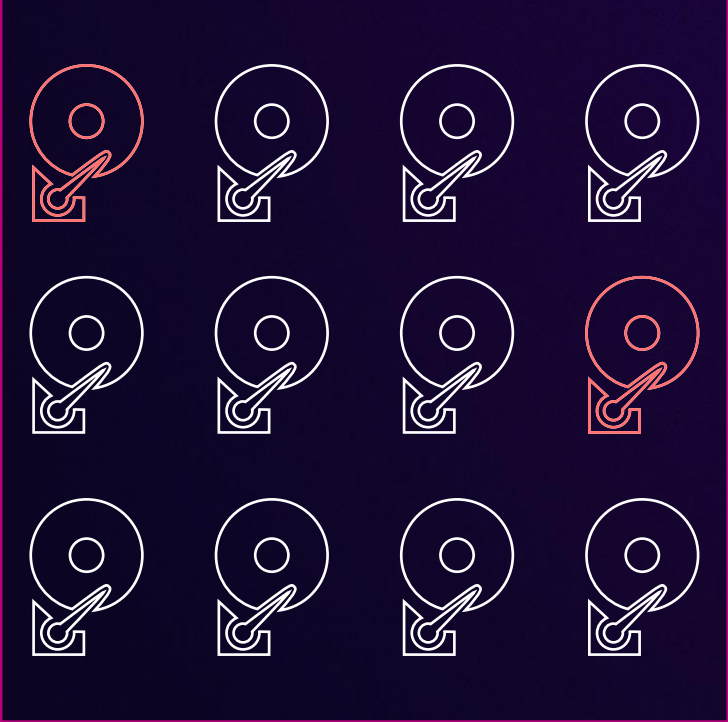**Workload isolation using shuffle sharding**
Amazon Builder's Library
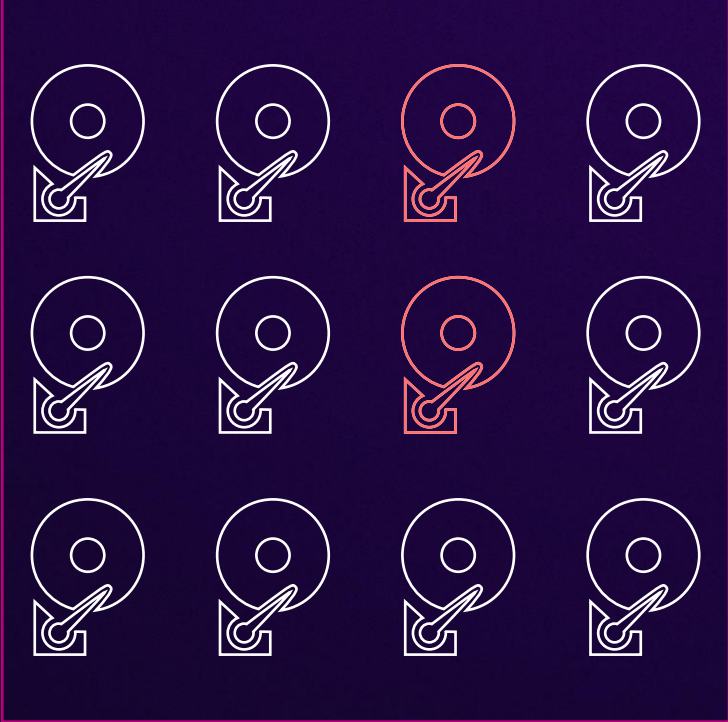
# At scale, failures are a fact of life
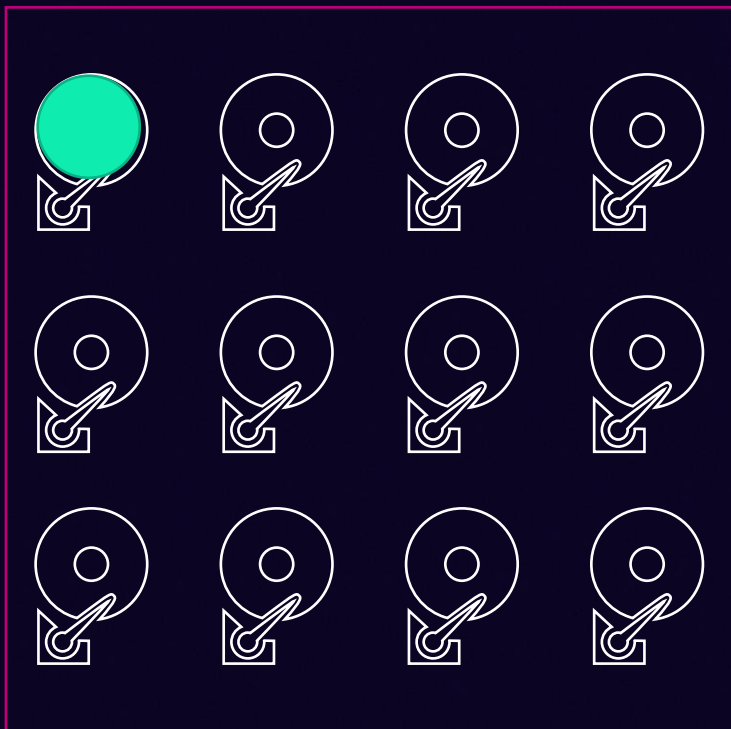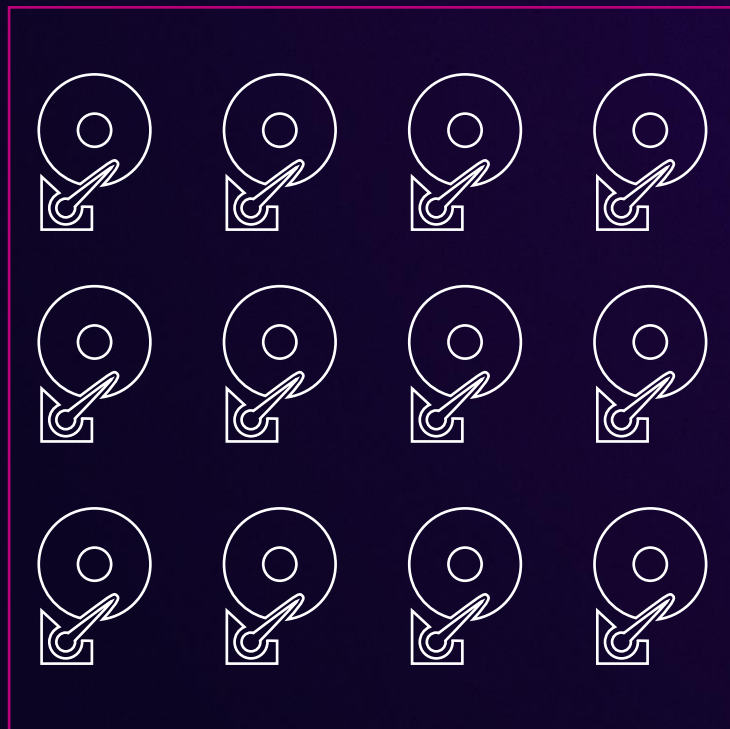


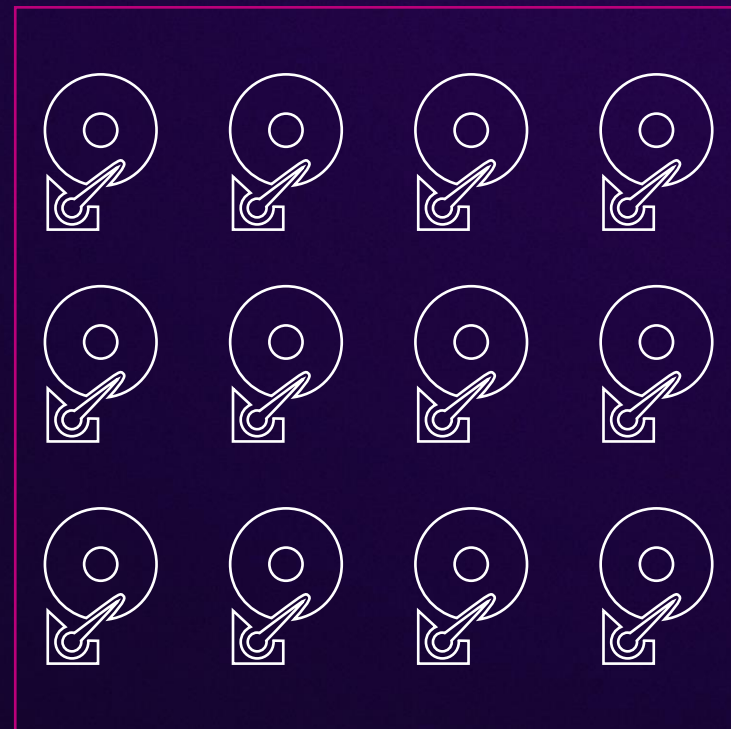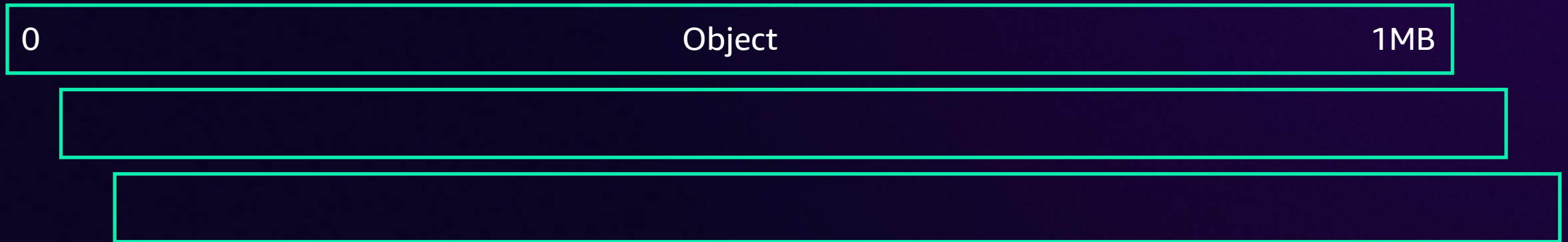AZ1                    AZ2                    AZ3

# At scale, failures are a fact of life



AZ1

AZ2

AZ3

# Erasure coding for fault tolerance

**Replication** is a simple way to tolerate faults – both individual drives and entire Availability Zones – but comes with high overhead
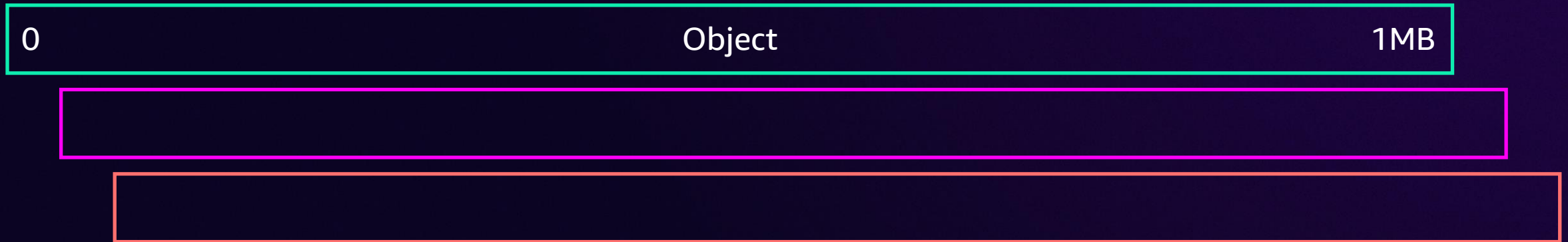
| 0 | Object | 1MB |

# Erasure coding for fault tolerance

**Replication** is a simple way to tolerate faults – both individual drives and entire Availability Zones – but comes with high overhead

| 0 | Object | 1MB |
|---|--------|-----|

# Erasure coding for fault tolerance

**Erasure coding** uses math to split the object into **shards** and create extra **parity** shards; the object can be rebuilt from **any K of the shards**

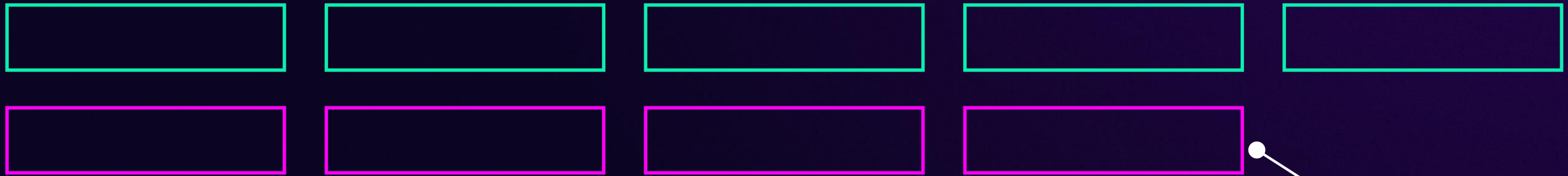| 0 | Object | 1MB |
|---|:---:|---:|

# Erasure coding for fault tolerance

**Erasure coding** uses math to split the object into **shards** and create extra **parity** shards; the object can be rebuilt from **any K of the shards**

**Any 5 of these shards**
is enough to rebuild the object,
with only 1.8× overhead
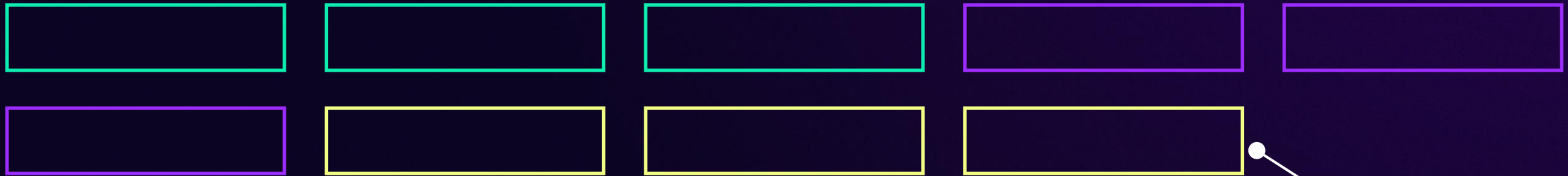
# Erasure coding for fault tolerance

**Erasure coding** uses math to split the object into **shards** and create extra **parity** shards; the object can be rebuilt from **any K of the shards**

**Any 5 of these shards**
is enough to rebuild the object,
with only 1.8× overhead

# Fault tolerance also gives velocity

Erasure coding allows us to deploy new software or hardware safely by **exposing only a few shards** to the changes

# Fault tolerance also gives velocity

Erasure coding allows us to deploy new software or hardware safely by **exposing only a few shards** to the changes

New software, or new drive type, or …

# Fault tolerance and shuffle sharding

The **birthday paradox** tells us that shuffle sharding will expose **some** object to too much risk; we need to layer in constraints

# Fault tolerance and shuffle sharding

The b... ne
objec...

James Bornholt
Amazon Web Services
& The University of Texas at Austin

Rajeev Joshi
Amazon Web Services

Vytautas Astrauskas
ETH Zurich

Brendan Cully
Amazon Web Services

Bernhard Kragl
Amazon Web Services

Seth Markle
Amazon Web Services

Kyle Sauri
Amazon Web Services

Drew Schleit
Amazon Web Services

Grant Slatton
Amazon Web Services

Serdar Tasiran
Amazon Web Services

Jacob Van Geffen
University of Washington

Andrew Warfield
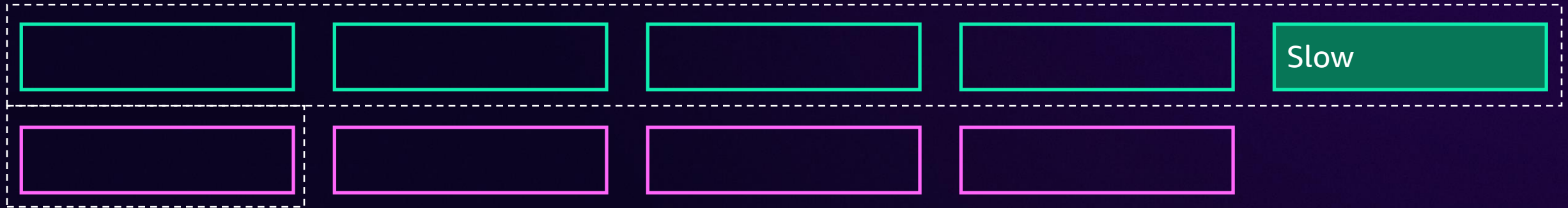Amazon Web Services

## Abstract

This paper reports our experience applying lightweight formal methods to validate the correctness of ShardStore, a new key-value storage node implementation for the Amazon S3 cloud object storage service. By "lightweight formal methods" we mean a pragmatic approach to verifying the correctness

1. Introduction

aws

# Improved performance through fault tolerance

Additional shards and shuffle sharding allow us to hedge against tail latency by **overreading**
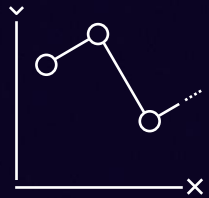
# Fault tolerance and velocity

- Failures are common at scale, and **fault tolerance is more than design**

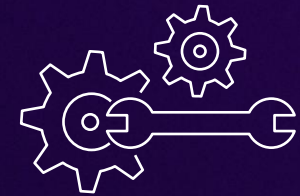- Fault tolerance, used wisely, also **improves developer velocity**



**ShardStore**

**Physics of data**

**Engineering for decorrelation**

**Fault tolerance and velocity**

# Thank you!

**Seth Markle**

smarkle@amazon.com

**James Bornholt**

bornholt@amazon.com

Please complete the session survey in the mobile app