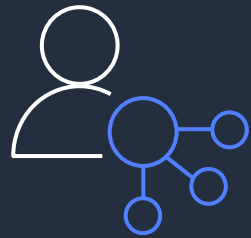# Amazon Elastic File System
## Building Blocks for a Cloud-Native File System

Jacob Strauss, Principal Engineer

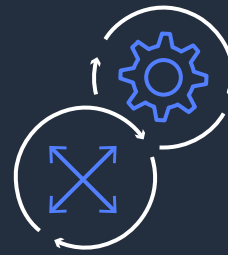Storage Developer Conference
September 2020

# Cloud-Native File System
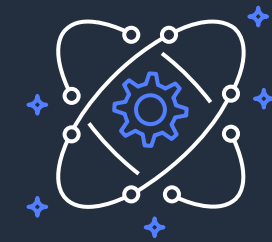
# Cloud-native file system

## Elastic

- Grow & shrink on demand
- No need to provision and manage infrastructure & capacity
- Pay as you go, pay only for what you use
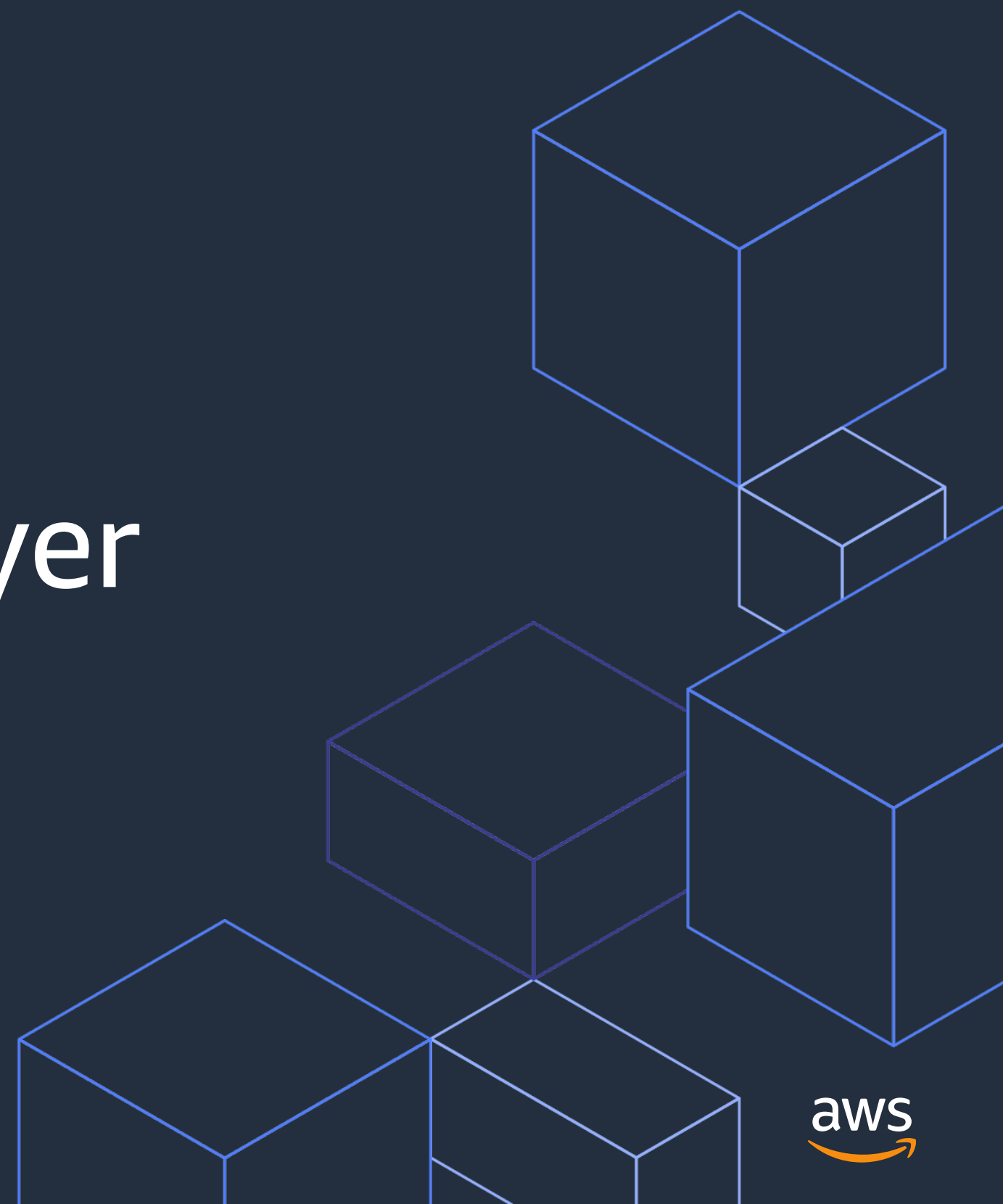- Simple to use, create a file system in seconds

## Scalable

- Grow up to petabytes
- Performance modes for low latencies and maximum I/O
- Throughput that scales with storage
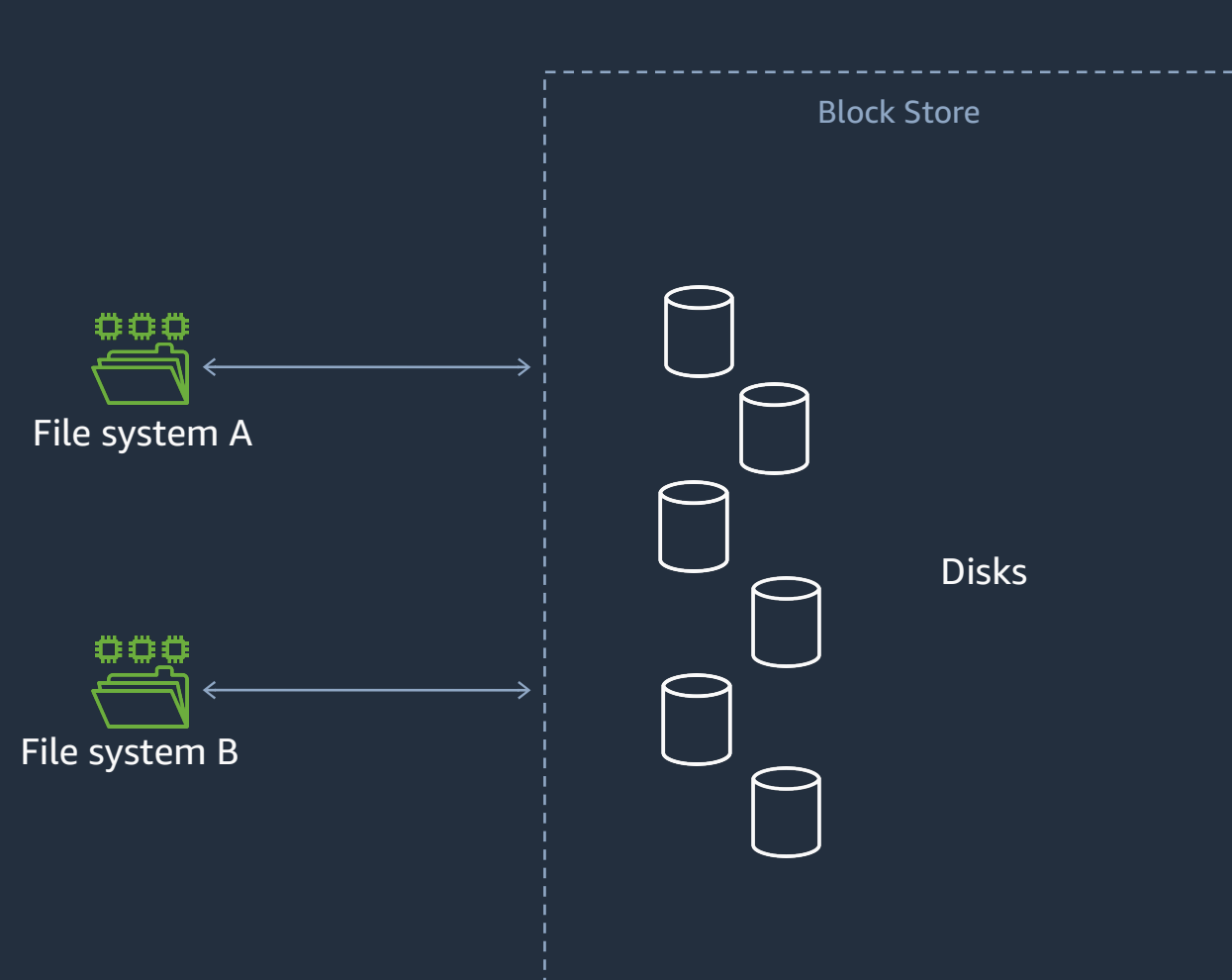- Provisioned throughput available

## Integrated

- Shared access from on-premises, inter region, and cloud-native applications
- Integrated with various AWS computing models
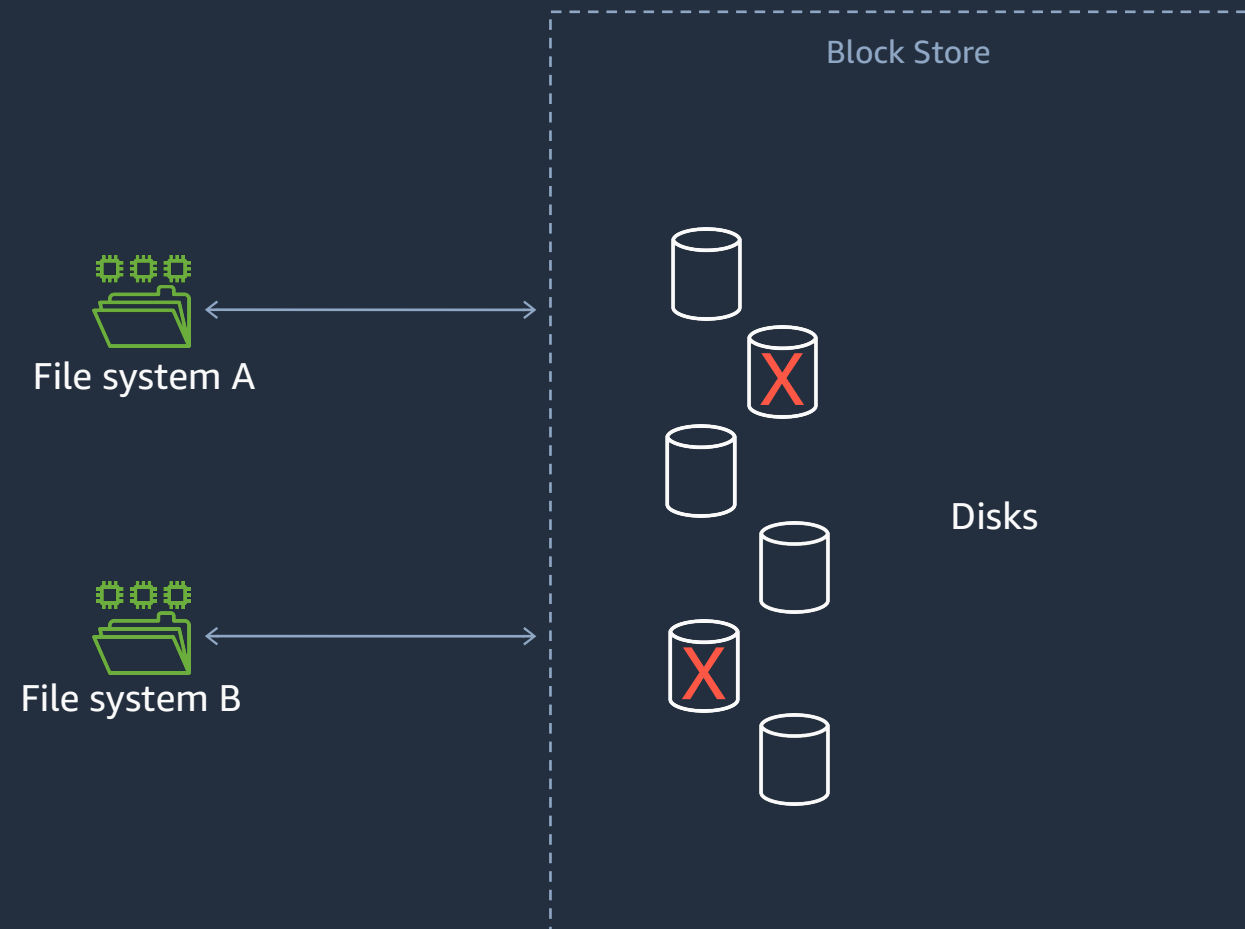- Access concurrently from Amazon EC2, AWS Lambda, and Amazon ECS and EKS containers
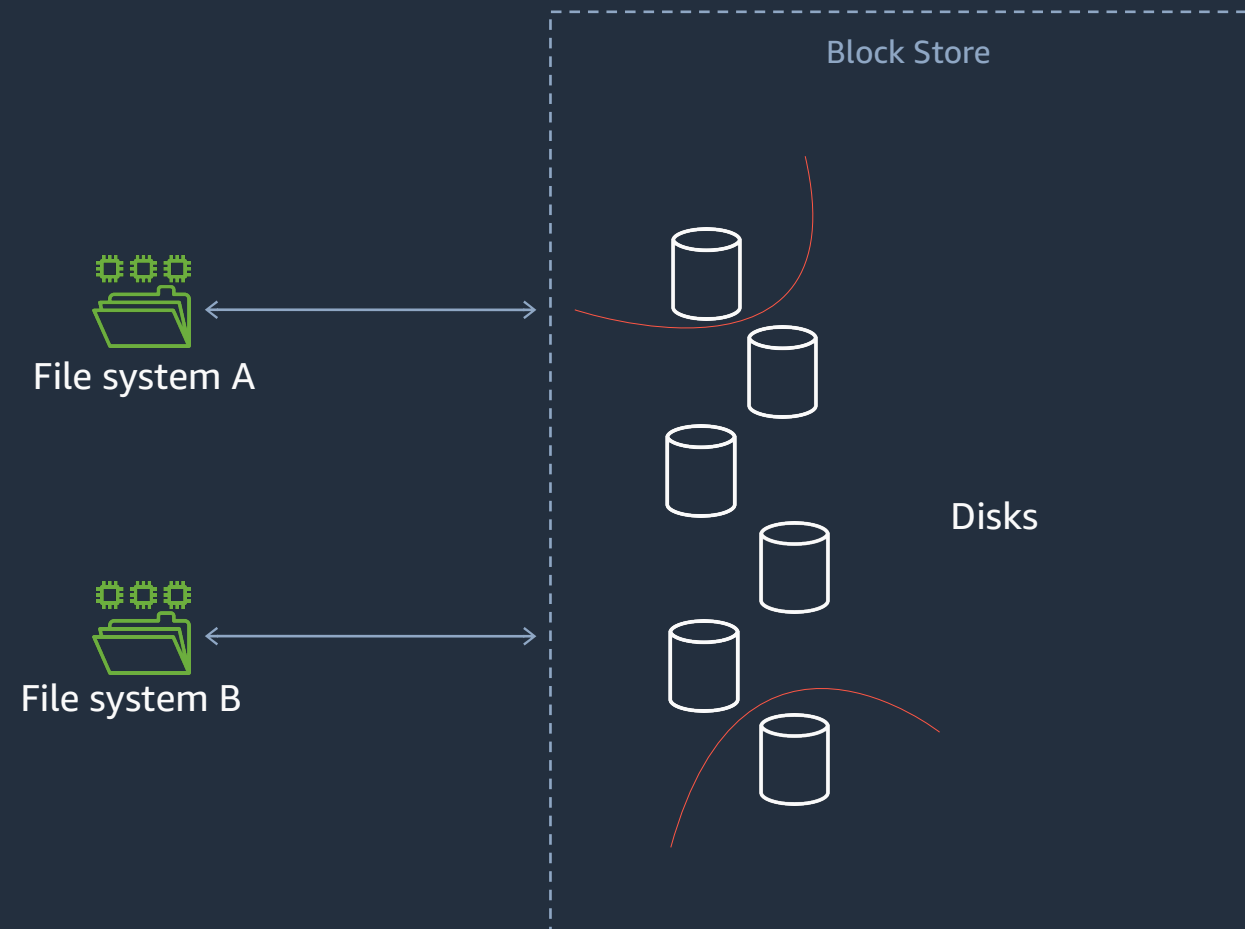
4

aws

# Deep Dive: Block Layer

aws

# Block responsibilities

Block Store

File system A

File system B

Disks

aws

# Block responsibilities: security & durability



Block Store

File system A

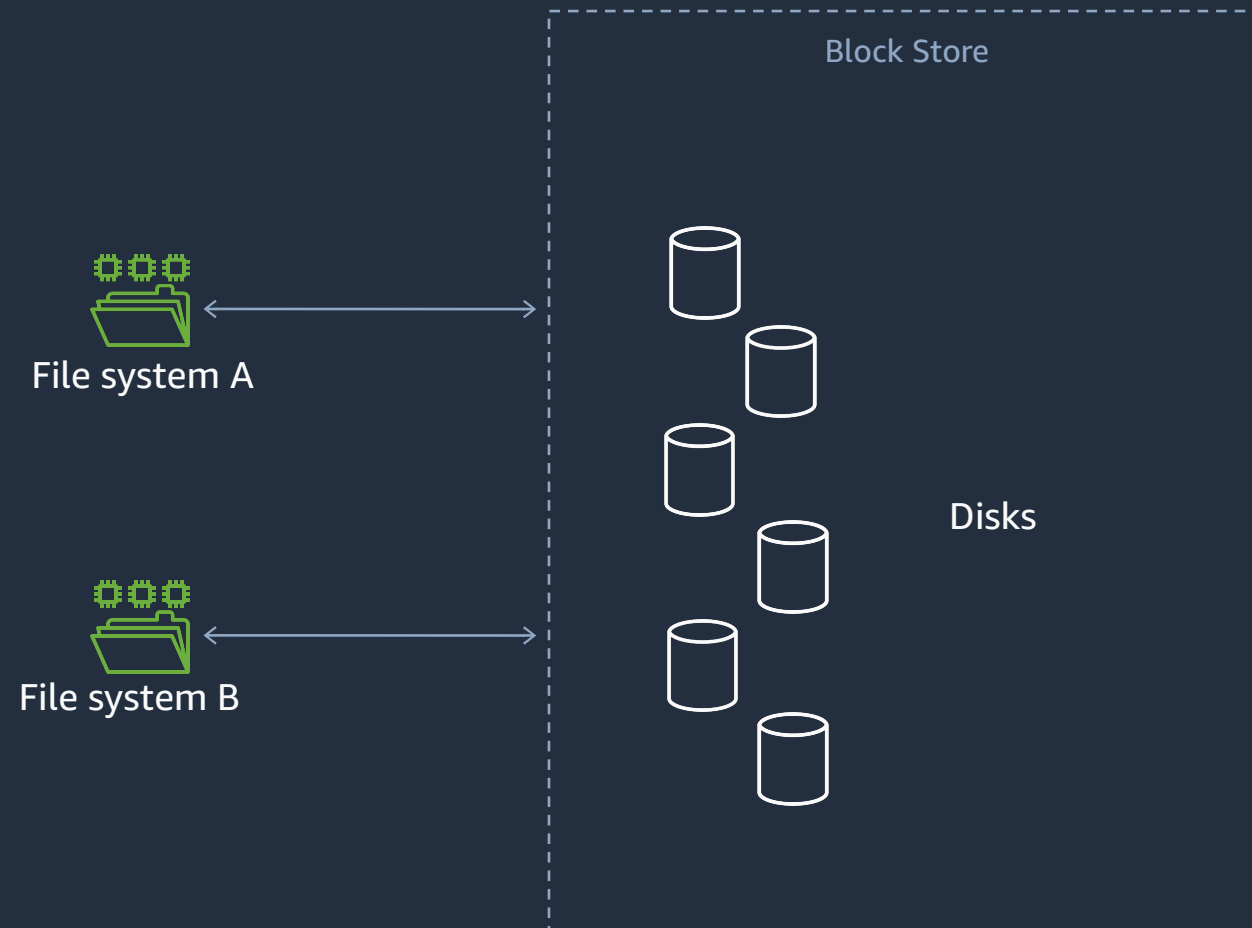File system B
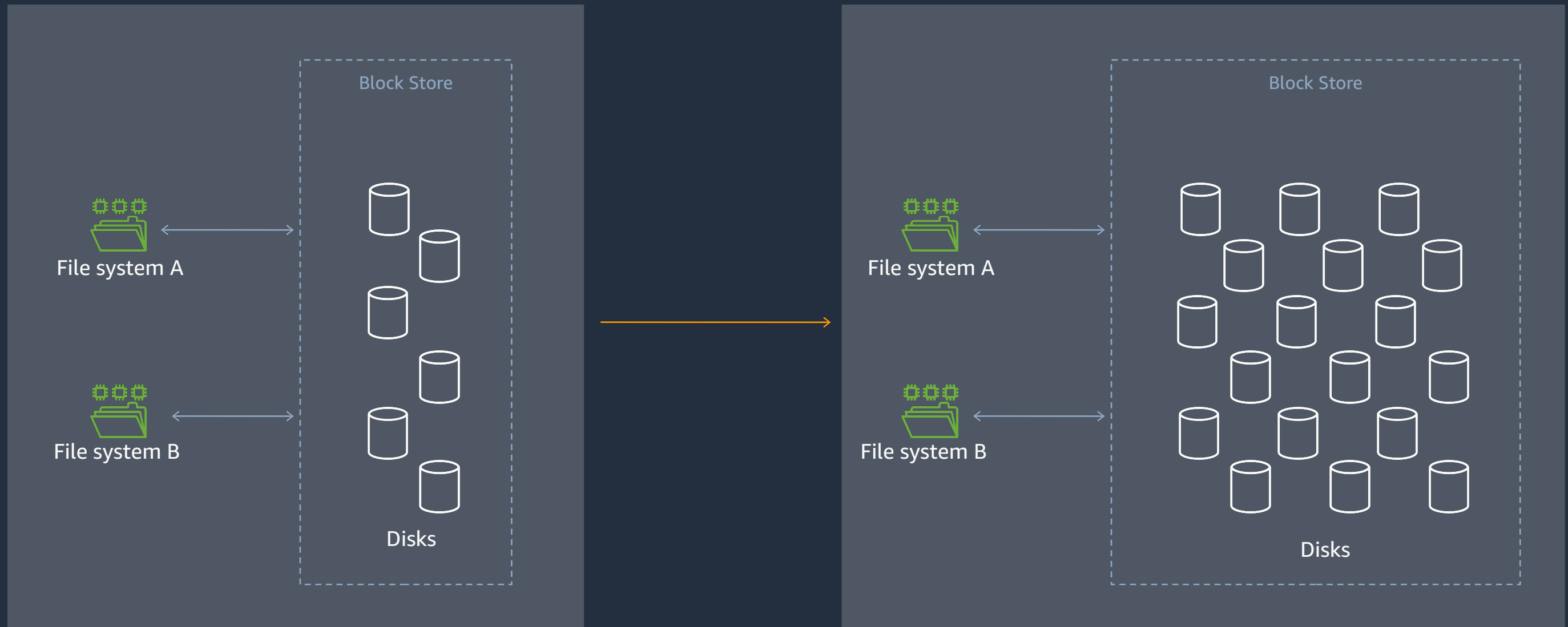
Disks

7

aws

# Block responsibilities: availability

# Block responsibilities: naming, locating, interface

- Single address space
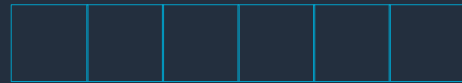- Location independence
- Read, write, allocate, delete
- Ordering

Block Store

File system A

File system B
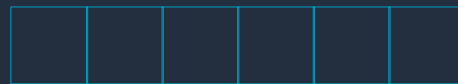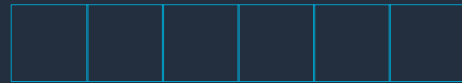
Disks

aws

# Block responsibilities: elastic scaling

# Logical data structure: Extent

Extent A

(A, 2)

Extent B

(B, 0)

aws

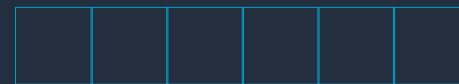# Logical data structure: Extent

Extent A

(A, 2)

Extent B

(B, 0)

Extent C

aws

# What is an extent?

## Paxos Replicated State Machine

Extent Logical View

State

Blocks

aws

# What is an extent on disk?

# Extent replicas across availability zones

# Cellular architecture



Cell A

Cell B

Cell C

File system

File system

Block Store

Block Store

Block Store

# Operation ordering: conditional writes

**Block Interface**

(data, version) ← read(block)

version ← write(block, data, version)

aws

# Operation ordering: conditional writes

Block Interface

(data, version) ← read(block)

version ← write(block, data, version)

version[] ← multiWrite(block[], data[], version[])

aws

# Operation ordering: conditional writes

**Block Interface**

(data, version) ← read(block)

version ← write(block, data, version)

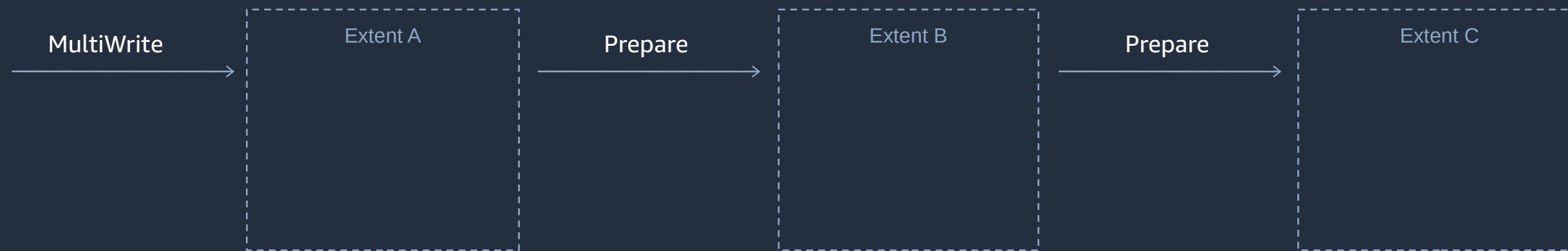version[] ← multiWrite(block[], data[], version[])

**File Operations**

Begin transaction

Read blocks

Modify blocks

Commit or start over

aws

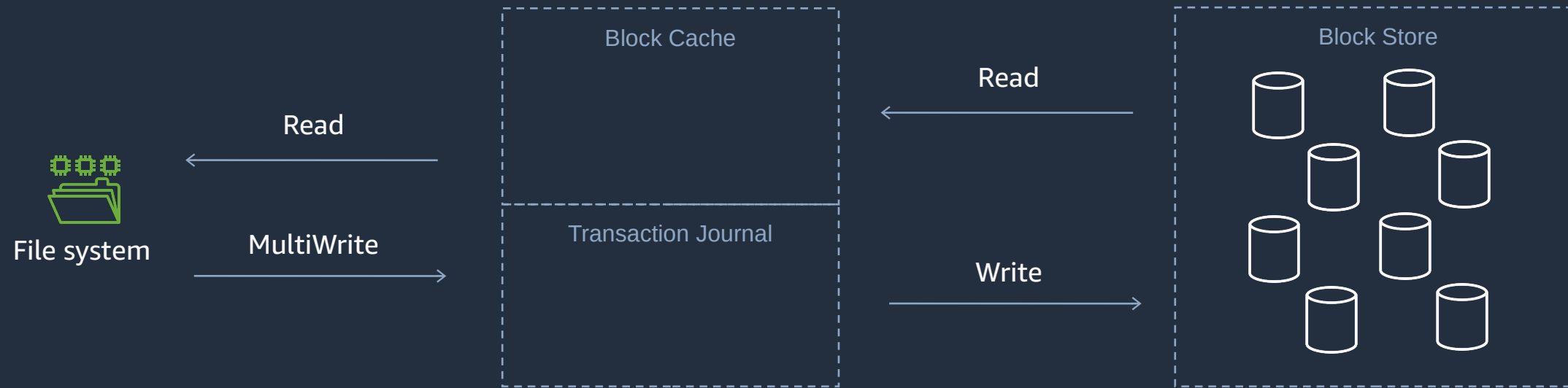# Multi-block conditional write is a 2-phase commit

MultiWrite →

Extent A

Prepare →

Extent B

Prepare →

Extent C

aws

# Multi-block conditional write is a 2-phase commit

MultiWrite → Extent A

Prepare → Extent B

Commit ← 

Prepare → Extent C

Commit ←

aws

# Reduce transaction latency



File system

Read

MultiWrite

Block Cache

Transaction Journal

Read

Write

Block Store

aws

# Cloud-native file system building blocks

### Properties

- Durable and secure
- Available
- Elastic scaling
- Simple interface
- Resilient

### Techniques

- Replicated state machine
- Multi-AZ consensus
- Cellular architecture
- Optimistic transactions
- Scale-out and lower latency

aws

# Thank You!

Jacob Strauss
jsstraus@amazon.com

aws