# Storage in the Age of AI
## Rethinking Data Infrastructure
### for the Next Wave of Applications

London, UK | February 11, 2026

**Rob Pankow**

CEO, founder

Simplyblock

# Rob Pankow

/in/robertpankow

## CEO & Co-Founder at simplyblock.io

## 📈 Professional Background

🌐 **Simplyblock**
Building cloud-native storage & database technologies for teams in private clouds, regulated environments or BYOC.

⏳ **10y+ experience in start-up and scale-ups**
Ex-Rocket Internet & ex-Delivery Hero. Scaled businesses of different sizes internationally.

## 🗄 Technical Focus

🗄 **Mission-Critical Infrastructure**
Building high-performance data infrastructure for zero-downtime workloads.

🧊 **Cloud-Native Architecture**
Kubernetes-native storage and Serverless Postgres

🛡 **Private & Sovereign Cloud**
Bridging the gap between public cloud UX and private data center requirements.

# AI is breaking traditional infra

Architectures optimized for virtual machines are failing under the pressure of AI/ML workloads.

**simplyblock**

## Legacy VM-Era Storage

**IOPS Throttling**
Performance caps silently kick in during bursts, causing GPUs or workers to wait for data.

**Overprovisioning**
Teams pay for peak I/O all the time, even though AI workloads only spike during training or evaluation.

**Noisy Neighbors**
Parallel training jobs or analytics queries steal I/O from each other, creating latency spikes that slow or restart runs.

**VS**

## AI on Kubernetes

**Extreme Parallelism**
Thousands of concurrent threads and pods accessing data simultaneously.

**Tight SLOs**
Sub-millisecond latency requirements to keep expensive GPUs fed.

**Hot & Cold Data Mix**
Need for ultra-fast NVMe caching coupled with massive object stores.

# AI needs high speed & parallelism

Moving beyond legacy constraints to meet the demands of modern intelligent applications.

## Extreme Performance
**Raw Speed & Throughput**

### Microsecond Latency

Sub-millisecond access times essential for GPU saturation and real-time inference.

- Millions of IOPS per cluster
- Kernel bypass networking

## Massive Parallelism
**Concurrent Access**

### Linear Scale-Out

Add nodes to increase capacity and performance linearly without bottlenecks.

- Thousands of concurrent pods
- Shared data access patterns

## Developer/Agent Agility
**Workflow Velocity**

### Instant Ephemeral Environments

Zero-copy cloning for isolated dev/test environments.

- Ephemeral environments for CI/CD
- Point-in-Time Recovery (PITR)

## Unified Architecture
**Simplicity & Efficiency**

### OLTP + OLAP

Unified platform for transactional and analytical workloads.

- Performance-optimized storage topology
- Cost-efficient tiering.

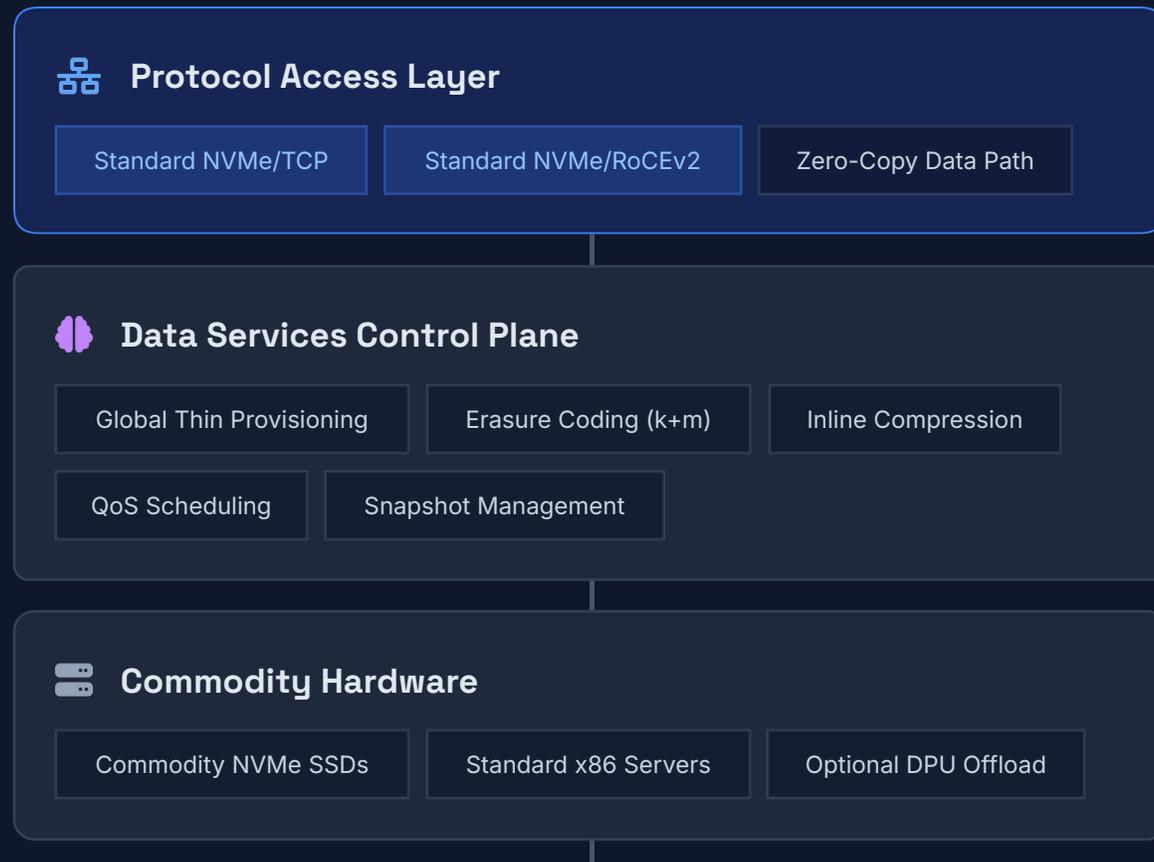simplyblock

# Traditional storage can't keep up

**Storage should feel like local disk.** Simple, fast and native to the application.



## Legacy Virtualized Storage
Kernel-heavy data path designed for HDDs

- Application
- Filesystem Layer Overhead — **Overhead**
- Block Layer (Kernel) Latency — **Latency**
- SCSI/iSCSI Driver Context Switch — **Context Switch**
- Hypervisor / vSwitch Interrupts — **Interrupts**
- Physical Storage

**Context Switching Heavy**       **Coupled to Node**

## NVMe-oF
Modern, efficient, scalable architecture

**AI Application (Kubernetes)**

**Kernel Bypass**
Direct access from userspace, skipping kernel interrupts

**Zero-Copy I/O**
Data moves directly from storage to network

**Disaggregated and/or Hyper-converged**
Storage scales independently from compute nodes

**Shared NVMe Storage Pool**

**Independent & Elastic**       **Direct & Efficient**

**Kubernetes has become the control plane where AI, data, and storage converge**

simplyblock

# Whole data stack needs to be refreshed

Key characteristics of next-generation software-defined storage platforms.

simplyblock

## Protocol Access Layer

| Standard NVMe/TCP | Standard NVMe/RoCEv2 | Zero-Copy Data Path |

## Data Services Control Plane

| Global Thin Provisioning | Erasure Coding (k+m) | Inline Compression |

| QoS Scheduling | Snapshot Management |

## Commodity Hardware

| Commodity NVMe SSDs | Standard x86 Servers | Optional DPU Offload |

**Dual Protocol Flexibility**

Modern stacks leverage **NVMe/TCP** for ubiquity and **NVMe/RoCEv2** for ultra-low latency <200µs.

**Storage Efficiency**

Thin provisioning and erasure coding maximize raw storage utilization (75%+ savings) without the overhead of traditional RAID.

**Instant Data Agility**

Zero-copy snapshots enable rapid cloning for CI/CD and databases, decoupling data size from operational speed.

**Strict Tenant Isolation**

QoS guarantees enforce IOPS/throughput limits per volume, preventing "noisy neighbors" in shared multi-tenant clusters.

# It's worth it

Comparing simplyblock NVMe-oF against traditional Ceph clusters.

**4x**

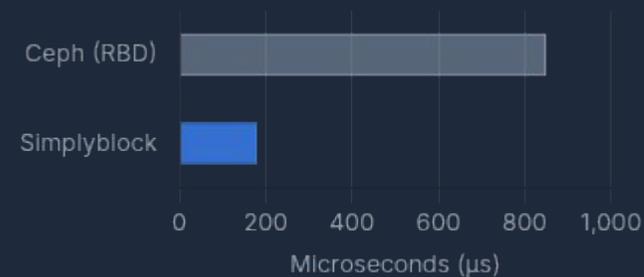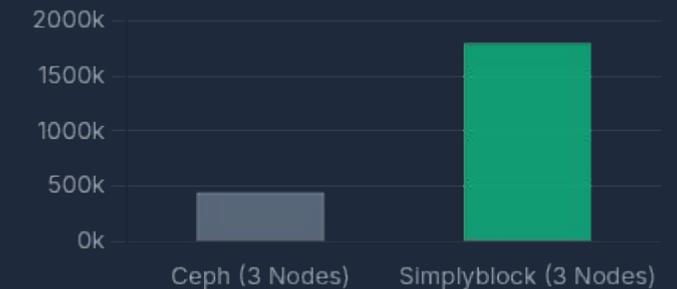Faster Throughput

↑ vs Legacy SDS

**<200μs**

Read Latency

✓ Consistent P99

**Millions**

IOPS Per Cluster

🗄 Linear Scaling

## ⏱ Latency (Lower is Better)

Ceph (RBD)

Simplyblock

0  200  400  600  800  1,000

Microseconds (μs)

**75% Less Hardware**

## 🎛 IOPS Performance

2000k
1500k
1000k
500k
0k

Ceph (3 Nodes)   Simplyblock (3 Nodes)

## 🖲 Hardware Efficiency (Resources to achieve 1M IOPS)

Legacy SDS (Ceph)                                    4+ Nodes Required

Simplyblock                                          1 Node Sufficient

ⓘ Kernel bypass and NVMe-optimization allow simplyblock to saturate hardware limits, requiring significantly fewer CPU cores and servers for the same workload.

# The story of modern storage

Modern AI pipelines require storage that behaves like local flash but scales like cloud.

## NVMe-oF (TCP & RoCE)

**Why it matters**

Delivers local NVMe-like performance over the network, essential for keeping expensive GPUs fully utilized without I/O wait.

## Kernel Bypass & Zero Copy

**Why it matters**

Eliminates CPU context switching overhead, maximizing throughput and minimizing tail latency for inference requests.

## QoS Service Classes

**Why it matters**

Guarantees IOPS and bandwidth for critical training jobs in multi-tenant clusters, preventing noisy neighbor issues.

## Live Migration

**Why it matters**

Allows legacy VM-based workloads to run alongside containers and enables node maintenance without disrupting long training jobs. Needs RWX at block storage level.

## Instant Snapshots & Clones

**Why it matters**

Accelerates CI/CD pipelines and enables data scientists to experiment on production-sized datasets instantly without duplication costs.

## Linear Scalability

**Why it matters**

**Shared-everything architecture** allows storage capacity and performance to grow independently of compute resources.

simplyblock

# Evaluating CSI drivers for AI workloads

| Capabilities | Ceph (Rook) | Longhorn | OpenEBS (Mayastor) | Simplyblock |
|---|---|---|---|---|
| NVMe-oF (TCP/RoCE) | Limited | ✗ | ✓ | ✓ |
| Kernel Bypass | ✗ | ✗ | ✓ | ✓ |
| QoS Service Classes | Limited | ✗ | ✗ | ✓ |
| Live Migration | ✓ | ✓ | Limited | ✓ |
| Instant Snapshots | ✓ | ✓ | ✓ | ✓ |
| Multi-tenancy | ✓ | Limited | ✗ | ✓ |
| Performance (IOPS) | Medium | Low | High | Very High |
| DPU/SmartNIC Support | ✗ | ✗ | ✗ | ✓ |

ⓘ For detailed CSI driver comparisons and benchmarks, visit **storageclass.info/drivers**

StorageClass.info **CSI Drivers**   Drivers   What is StorageClass?   Glossary   Sponsors

simplyblock

# CSI Drivers Directory

**storageclass.info/drivers**

🔍 Search drivers by name, provider, or description...

**Filters**                                      Reset

**Storage Types**                          ⌃

☐ Block

☐ File

☐ Object

**Capabilities**                              ⌃

☐ Dynamic

☐ Snapshot

☐ Raw

☐ Expansion

☐ Clone

☐ Topology

☐ Tracking

☐ QoS

☐ NVMe-oF

☐ iSCSI

**Access Modes**                            ⌃

☐ Read Only Many

☐ Read Write Once

Showing **136** of **150** CSI drivers

### AlibabaCloud Disk

CSI Driver for an AlibabaCloud Disk

`diskplugin.csi.alibabacloud.com`

`dynamic`  `snapshot`  `raw`  `expansion`  `topology`
`block`

Source ⧉                              Details >

### AlibabaCloud Nas

CSI Driver for AlibabaCloud Network Attached Storage (NAS)

`nasplugin.csi.alibabacloud.com`

`dynamic`  `block`

Source ⧉                              Details >

### AlibabaCloud Oss

CSI Driver for AlibabaCloud Object Storage Service (OSS)

`ossplugin.csi.alibabacloud.com`

`object`

Source ⧉                              Details >

### Alluxio

CSI Driver for Alluxio File System)

`csi.alluxio.com`

`dynamic`  `file`

Source ⧉                              Details >

### ArStor CSI

CSI Driver for Huayun Storage Service (ArStor)

`arstor.csi.huayun.io`

`dynamic`  `snapshot`  `raw`  `expansion`  `clone`
`block`

Source ⧉                              Details >

### AWS Elastic Block Storage

CSI Driver for AWS Elastic Block Storage (EBS)

`ebs.csi.aws.com`

`dynamic`  `snapshot`  `raw`  `expansion`  `block`

Source ⧉                              Details >

### AWS Elastic File System

CSI Driver for AWS Elastic File System (EFS)

`efs.csi.aws.com`

`dynamic`  `file`

Source ⧉                              Details >

### AWS FSx for Lustre

CSI Driver for AWS FSx for Lustre (EBS)

`fsx.csi.aws.com`

`dynamic`  `block`

Source ⧉                              Details >

### Azure Blob

CSI Driver for Azure Blob storage

`blob.csi.azure.com`

`dynamic`  `expansion`  `object`

Source ⧉                              Details >

# Why storage matters for developers and agents?

Traditional managed cloud databases feel like 1980's

## Instant Database Branching

**Requirement: Zero-Copy Cloning**

Developers need isolated environments for every feature branch without waiting for slow, expensive data copying operations.

- ✓ Parallel development workflows
- ✓ Eliminate storage duplication (CoW)
- ✓ Test with production-like data safely

## Visual Data Observability

**Requirement: Schema Transparency**

Complex microservices architectures require clear visibility into data models and relationships without digging through raw SQL.

- ✓ Auto-generated ER diagrams
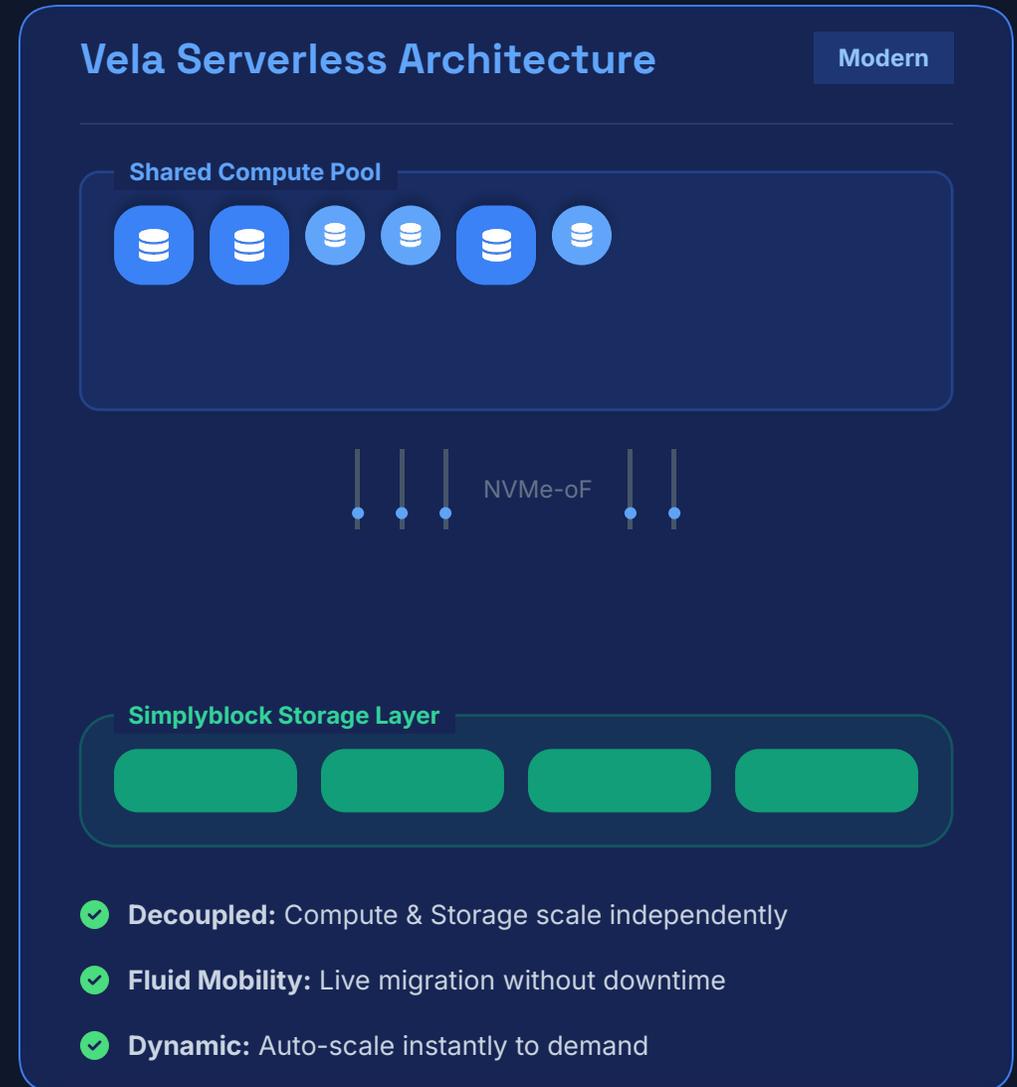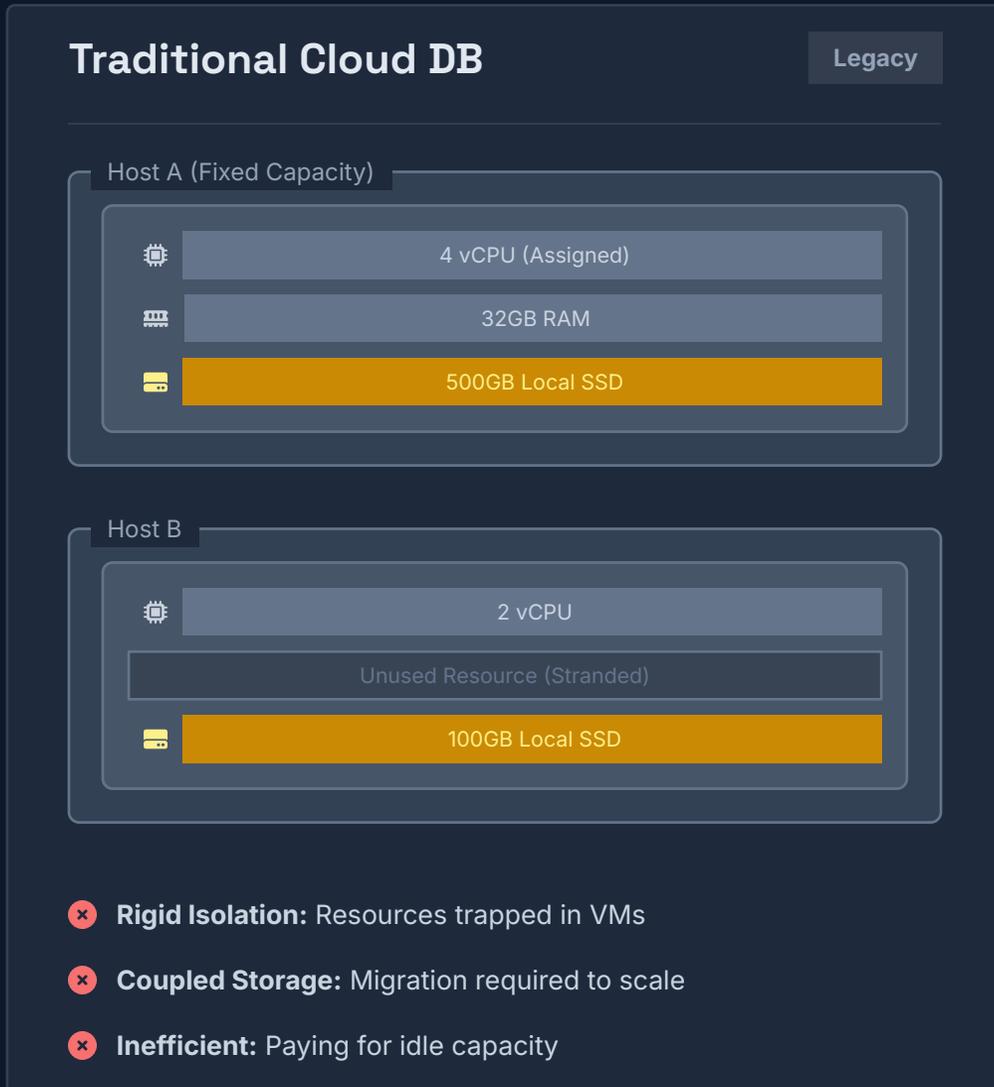- ✓ Visual relationship mapping
- ✓ Self-documenting infrastructure

## Point-in-Time Recovery

**Requirement: Granular Resilience**

Modern CI/CD pipelines require the ability to instantly rewind databases to specific states for testing or disaster recovery.

- ✓ Second-level recovery granularity
- ✓ Instant rollback from testing errors
- ✓ Reduced RTO for critical failures

simplyblock

# AI-ready storage enables data stack

**simplyblock**

## Traditional Cloud DB

### Host A (Fixed Capacity)

| | |
|---|---|
| ▦ | 4 vCPU (Assigned) |
| ▦ | 32GB RAM |
| ▤ | 500GB Local SSD |

### Host B

| | |
|---|---|
| ▦ | 2 vCPU |
| | Unused Resource (Stranded) |
| ▤ | 100GB Local SSD |

❌ **Rigid Isolation:** Resources trapped in VMs

❌ **Coupled Storage:** Migration required to scale

❌ **Inefficient:** Paying for idle capacity

→

## Vela Serverless Architecture

Modern

### Shared Compute Pool

NVMe-oF

### Simplyblock Storage Layer

✅ **Decoupled:** Compute & Storage scale independently

✅ **Fluid Mobility:** Live migration without downtime

✅ **Dynamic:** Auto-scale instantly to demand

The main problem Postgres has, as I see it, is that it isn't a native cluster and there's no way to make it into one. It's a single database engine in which replicas exploit continuous crash recovery in a way that's very loosely coupled. As a result, there is an entire cottage industry of HA tools, of which only Patroni really "got it right". Yet Patroni and Kubernetes solutions like CloudNativePG are just workarounds for that fundamental shortcoming.

One way around that is to introduce all the missing pieces: quorum, cluster metadata, connection routing, fencing, etc., to the Postgres project itself. Another option is to decouple storage from the engine itself, and use the engine as an interchangeable compute node. Provided you're on a distributed storage system and have a quorum-mediated write layer, any compute node can access any data across the entire storage fabric. No more worries about which node has what locally reproduced data. The question becomes: is the storage mounted? Done.

It looks like Vela is following in Neon's footsteps and choosing option two. Given extensions like pg_lake which leverage DuckDB for Parquet and Iceberg compatibility, that seems like the direction things are ultimately going anyway.

# Example: Vela Postgres Architecture

Best practice: Decouple storage, orchestration, and applications.

**Database** — Postgres

- Instant Branching
- PITR
- Dev API

**Virtualization** — K8s+ Vela OS

- QoS Mapping
- Snapshots
- Replication

**Storage Foundation** — Simplyblock

- NVMe-oF
- Thin Provisioning
- <200µs

## Instant Environments
Storage-level copy-on-write enables creating isolated database branches in seconds rather than minutes.
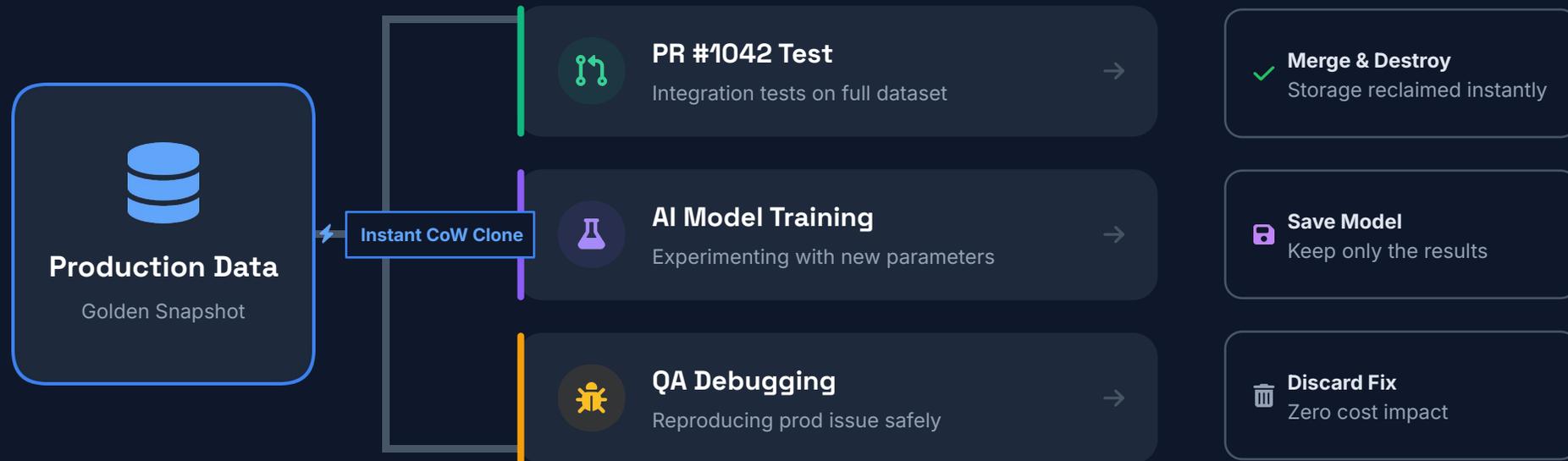
## Noisy Neighbor Protection
Orchestration maps app priorities to strict storage QoS classes, guaranteeing IOPS for critical workloads.

## Live Mobility
Decoupled architecture enables compute auto-scaling and zero-copy data migration.

simplyblock

# Ephemeral storage for ephemeral environments

Enabling CI/CD and AI experimentation with instant, low-cost data clones.

simplyblock

**Production Data**
Golden Snapshot

Instant CoW Clone

**PR #1042 Test**
Integration tests on full dataset

**AI Model Training**
Experimenting with new parameters

**QA Debugging**
Reproducing prod issue safely

✓ **Merge & Destroy**
Storage reclaimed instantly

**Save Model**
Keep only the results

**Discard Fix**
Zero cost impact

## Cost Efficiency

Copy-on-Write (CoW) ensures clones consume zero storage initially. You only pay for the changes made (delta).

## Developer Velocity

Spin up databases in seconds, not hours. No waiting for data hydration or restore scripts.

## Data Safety

Clones are isolated. Developers can break things in ephemeral environments without affecting production.

# Key Takeaways

Path to next-generation data infrastructure

## AI Demands a Storage Rethink

Legacy VM-centric storage cannot support the massive parallelism, throughput, and low-latency requirements of modern AI training and inference workloads.

**01**

## The Performance Formula

**NVMe-oF + Decoupling of Storage/Compute** is the only architecture that delivers local-flash speeds while maintaining the efficiency and scalability of cloud storage.

**02**

## Developer Experience Wins

Infrastructure must empower developers. Instant branching, zero-copy clones, and Point-in-Time Recovery are essential for high-velocity engineering teams.

**03**

## Cloud Experience On-Prem

Wherever you're building, the flexibility and ease of public cloud should be available to your teams. If your developer needs to contact DBA or Platform to spin up a database instance, you're doing it wrong.

**04**