

# Operating a Self-Healing Bare-Metal Kubernetes Platform at Global Scale

Theme: Open Source Ecosystem in Synergy

---

Presenters - Aparna Prabhu and Nikhil Pathak, DigitalOcean

# What We Operate

5000+

Pods

Running across 366  
Kubernetes nodes

15

Production Regions

~4M

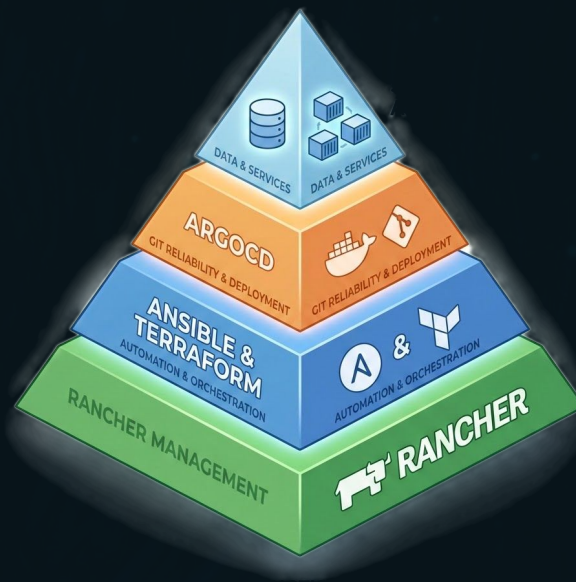
Average Requests/Hour

Hardware ranges from 40 cores to 4,500 cores, 78 GiB RAM to 16 TiB RAM, and 10 GiB to 7 TiB per PV

SLA and SLO - 99.9% # of customers - 640k

# An All Open Source Stack

Our operational strategy relies heavily on a diverse array of open-source initiatives to maintain a robust, production-grade infrastructure.



## Rancher

K8s clusters run directly on bare metal servers and are orchestrated through the open-source flavour of Rancher.

## StatefulSets as a Service

Workloads of multiple tenants and open-source statefulsets including Database clusters, Consul Clusters, and Redis Clusters.

## Ansible & Terraform

Configuration management through Ansible and Terraform for infrastructure-as-code practices.

## ArgoCD

Deployment management through ArgoCD for GitOps-driven continuous delivery.

# Why Bare Metal?

It's all about performance, control, and cost efficiency at scale!



No Hypervisor Overhead



Direct Hardware Access



Lower Latency



Custom Storage Backends



No Vendor Lock-in



Full Control over Control Plane & Auditing Powers

# Bare Metal Kubernetes Clusters

## Cross-cutting concerns

HA quorum · automation · observability · security

Container runtime

**Operating system**  
Host configuration

RBAC

Kernel & sysctl

CNI plugin

**Networking**  
Connectivity

CIDR planning

API server HA

CSI driver

**Storage**  
Persistent volumes

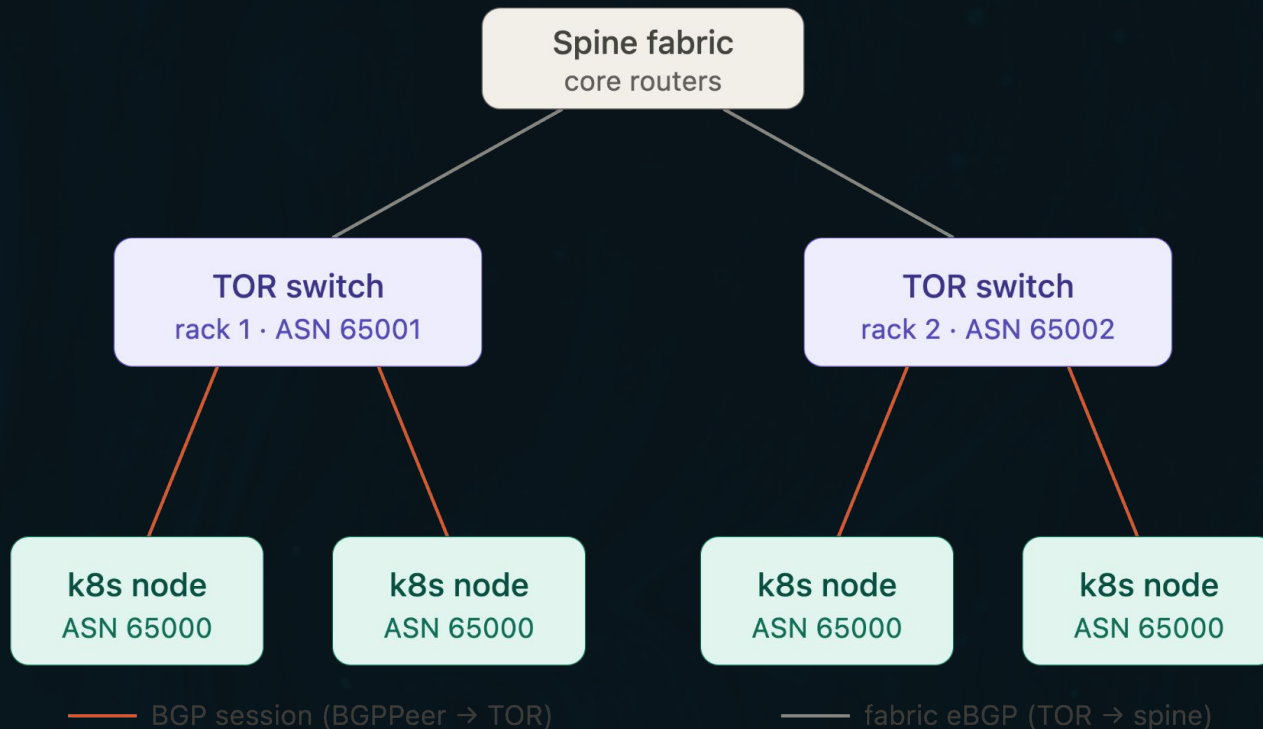
High performance of etcd

Backups

## Bare metal nodes

Physical servers · NICs · disks · BIOS/firmware

# Networking Considerations for Bare Metal Clusters



# Key Architectural Decisions

01

---

Cluster-of-Clusters

02

---

Intent-Driven Placement

03

---

Blast-Radius Containment

04

---

Automated Recovery

05

---

Continuous Reconciliation

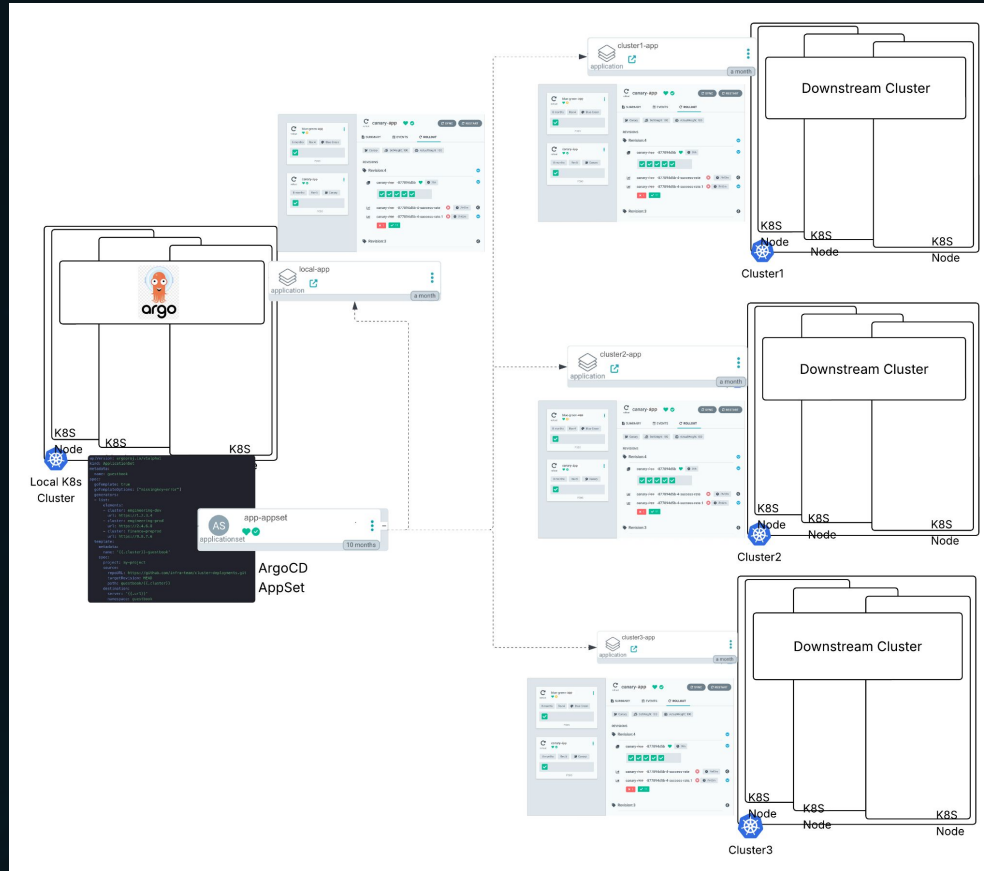
# Cluster-of-Clusters Topology

Automating Everything with Open source

- **Cluster Provisioning Automation**
- **Persistent Volume Creation**
- **Application Deployment**

☑ Spinning up a new region is a matter of just **6 hours!**

# Application Deployment



# Intent-Driven Placement



Choosing the Right  
Hardware



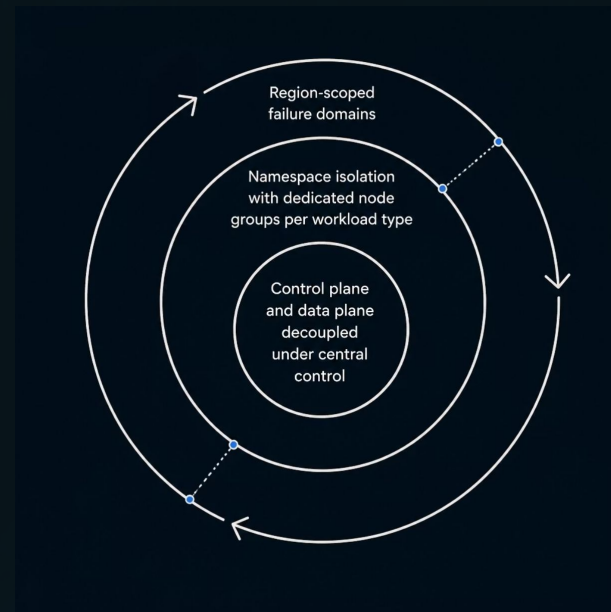
Physical Server  
Reusability



Least-Hop Network  
Path

# Blast-Radius Containment

Limiting fallout through dedicated failure domains



Central Control, Regional Service Provider

Control Plane & Data Plane Decoupled

Namespace Isolation  
Dedicated node groups per workload type with region-scoped failure domains.

# Automated Recovery

1

Single AWX Job

2

Intelligent Maintenance Logic

3

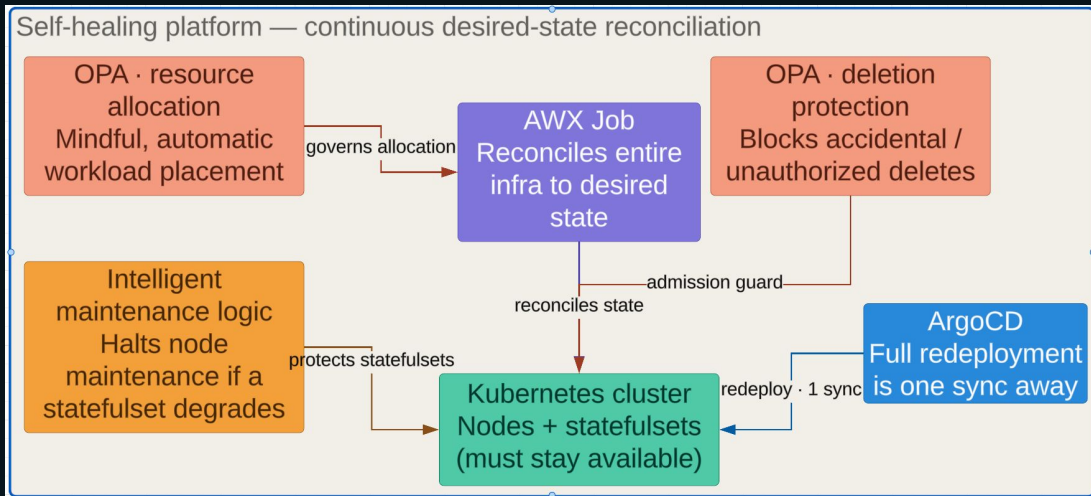
Policy-Driven Allocation

4

Enhanced Deletion Protection

5

Rapid Redeployment



# Continuous Reconciliation



Infrastructure & Machine  
Reconciliation



Application Lifecycle Management



Secret & Security Reconciliation



Self-Correction with  
Stork8s-Checker



Automated Data Protection



Automated Recovery

# Fitting the pieces of a giant puzzle

01

---

Unified Lifecycle  
Automation

02

---

Decoupled Architecture  
for Resilience

03

---

Hardware-First  
Performance

04

---

Policy-Driven  
Operations



# Small Team Handling Big Stuff



Questions?