# Distributed Key-Value Stores:
# Performance and Scalability for Flash Media

## Richard Elling

# Key-Value (KV) Databases Overview

- **Simple is important**
  - Fast, easy to use

- **Consistency is a differentiator**
  - Strong consistency
    - Important for control planes
    - Examples: etcd, consul, zookeeper
  - Weak consistency (aka eventually consistent, NoSQL)
    - Popular for large-scale applications
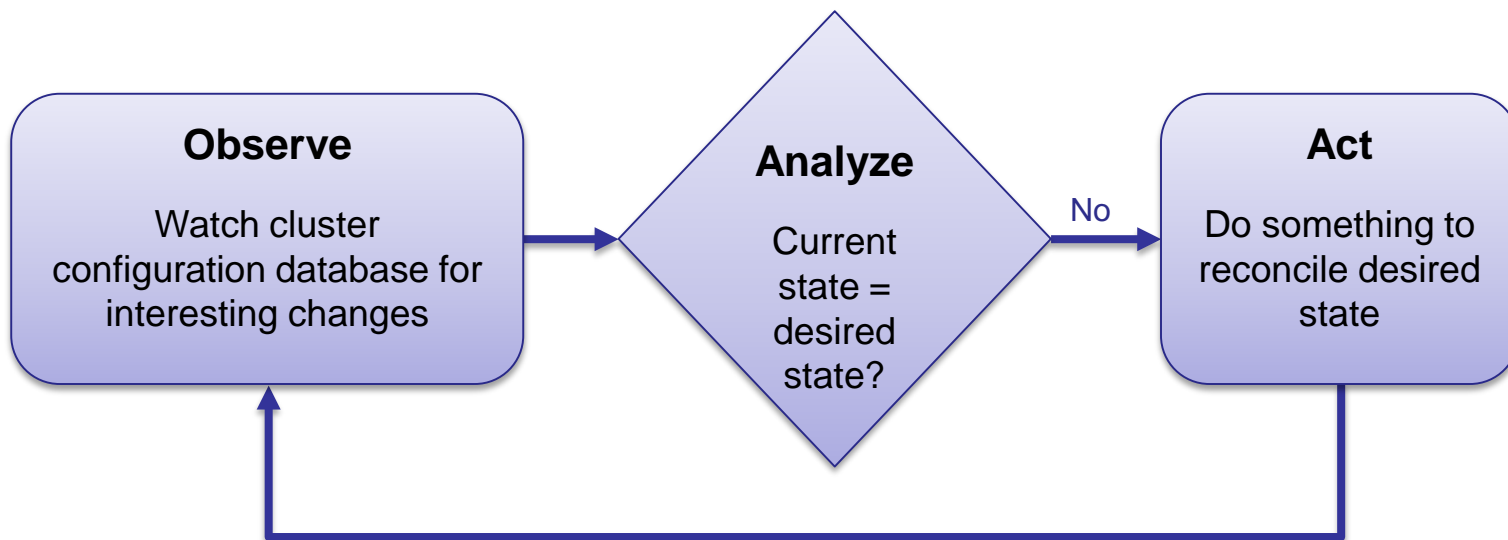    - Examples: dynamo, cosmos db

# Kubernetes and KV

- Distributed system scheduler for containers
- How it works
  - Declare desired state in key-value database
  - Controllers (aka operators) watch the state and act to change current state to desired state
- End result
  - Billions of applications served reliably
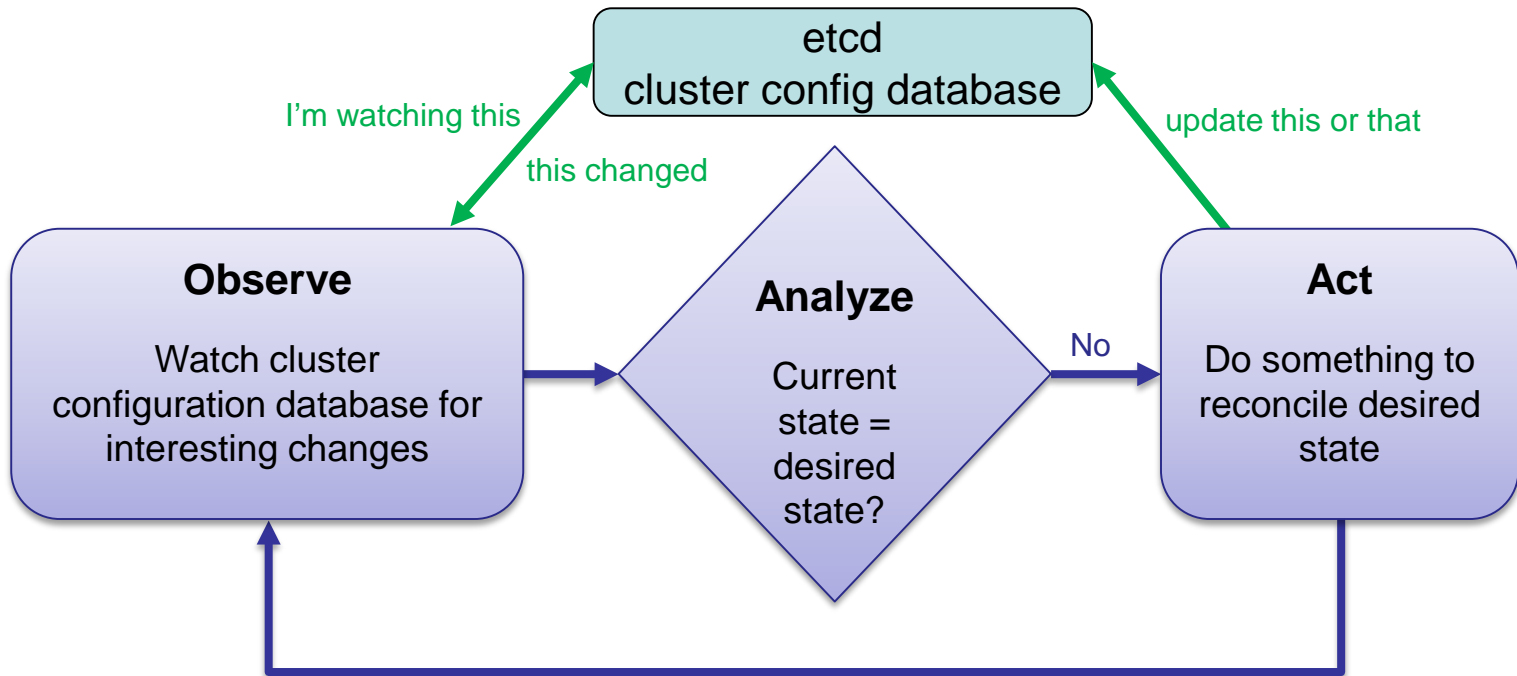  - Clusters up to 5,000 nodes each

VIKING™
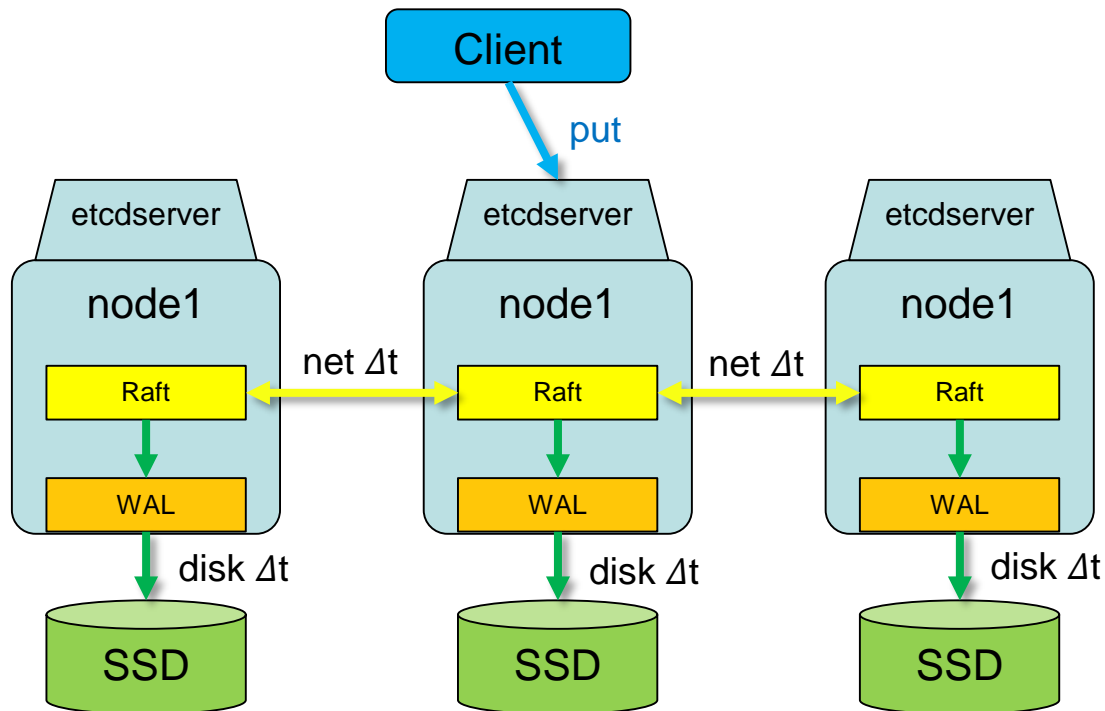Enterprise Solutions

# Kubernetes Controllers/Operators

```
┌─────────────────────────┐        ╱╲                    ┌─────────────────────────┐
│       Observe           │       ╱  ╲                   │         Act             │
│                         │      ╱Analyze╲    No         │                         │
│  Watch cluster          │─────▶ Current  ─────────────▶│  Do something to        │
│  configuration database │      ╲ state = ╱             │  reconcile desired      │
│  for interesting changes│       ╲desired╱              │  state                  │
│                         │        ╲state?╱              │                         │
└─────────────────────────┘         ╲  ╱                 └─────────────────────────┘
        ▲                            ╲╱                            │
        │                                                         │
        └─────────────────────────────────────────────────────────┘
```

**Observe** — Watch cluster configuration database for interesting changes

**Analyze** — Current state = desired state?

**Act** — Do something to reconcile desired state

No

# Kubernetes Controllers/Operators

# Etcd Cluster Configuration Database

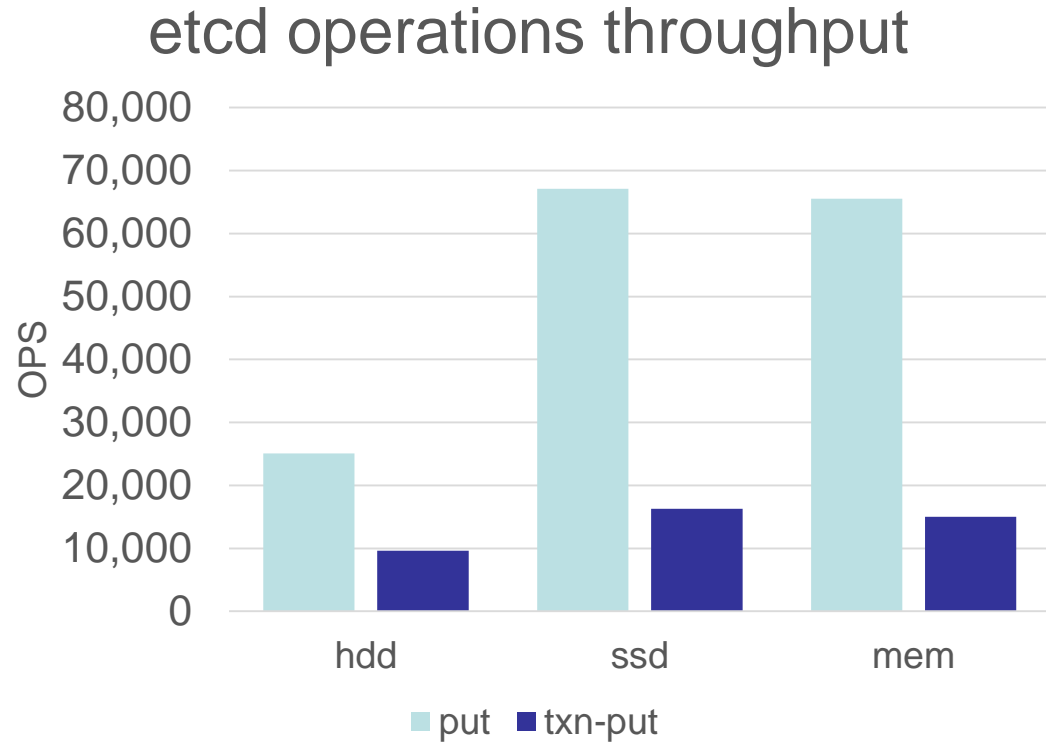| Requirements | Etcd in Kubernetes |
| --- | --- |
| Consistency | Single-writer<br>Updates acknowledged by quorum of masters |
| Availability | Multiple-masters |
| Partition tolerance | Raft protocol ensures all masters are consistent |
| Performance | Writes committed to persistent storage log<br>Reads satisfied by any master |

# Etcd Persistent Storage

# Kubernetes Practical Implementations

- Many nodes + fast SSD storage =
  - Many events to watch
  - When things break, many changes to process

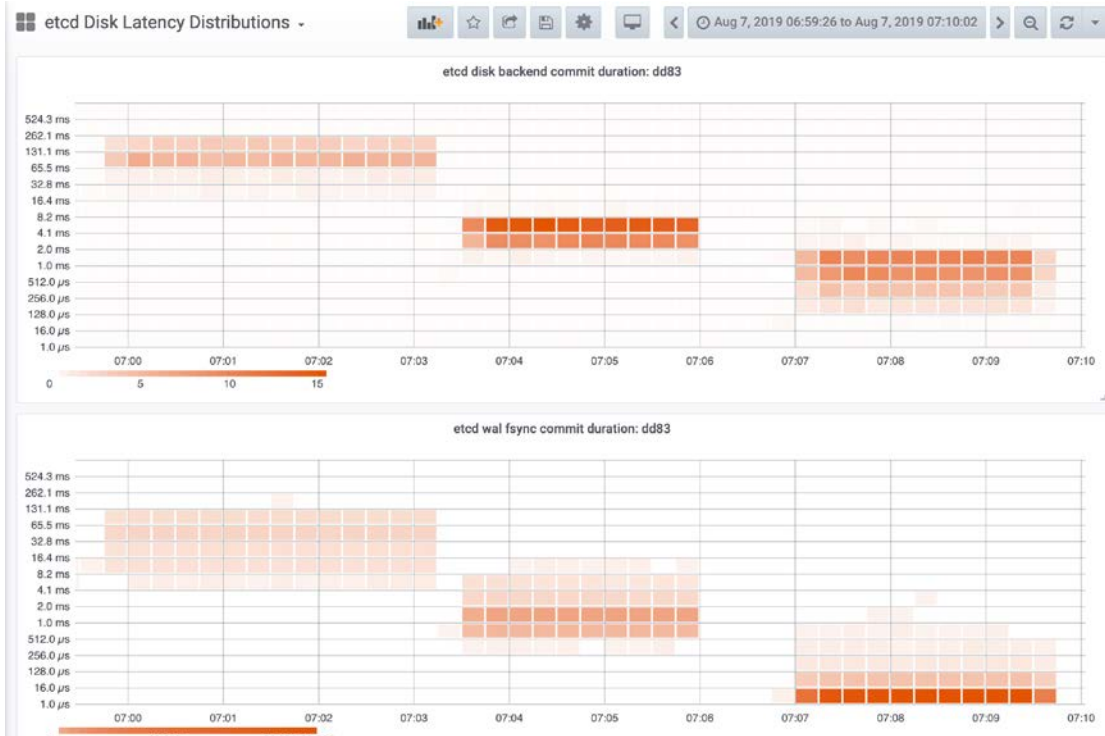- Can we improve etcd write workload scalability by adding faster storage?

# etcd Benchmark

## etcd operations throughput

# Transaction Latency Analysis

# Thank You

Richard.Elling@VikingEnterprise.com