

The Chronicles of Clusterville: A Kubernetes Tale

A five-chapter epic where containers come to life and orchestration becomes adventure

Chapter 1: The Birth of Clusterville and The Great Pod Housing Crisis

In the beginning, there was chaos. Applications ran wild on bare metal servers, fighting for resources, crashing into each other, and generally making a mess of things. Then came **Master Node**, the wise overseer who would bring order to this digital wilderness.

Master Node wasn't alone. He had three loyal advisors who formed the **Control Plane**:

API Server - The Grand Communicator, a fast-talking diplomat who spoke every language. Every request, every command, every whisper in Clusterville had to pass through him. He was the only one who could talk to **etcd**, the ancient librarian.

Spec Deep Dive: API Server runs on port 6443 (typically) and is the central management entity. It validates and configures data for pods, services, and controllers. When you run `kubectl`, you're talking to API Server. He's stateless and horizontally scalable.

etcd - The Memory Keeper, an obsessive record-keeper who wrote everything down in a distributed key-value diary. "Pod-27 lives at 10.244.1.5," she'd mutter while scribbling. "Deployment-Frontend wants 3 replicas." She never forgot anything and kept multiple backup copies because trust issues.

Spec Deep Dive: etcd is the source of truth. It stores cluster state, configuration data, and metadata. It uses the Raft consensus algorithm - needs a quorum ($N/2 + 1$) to function. If you have 3 etcd instances, you can lose 1. With 5, you can lose 2. Odd numbers are preferred. Default port: 2379 for client requests, 2380 for peer communication.

Scheduler - The Tetris Master, who decided which **Worker Node** (the actual neighborhoods) should house which **Pod** (the apartment buildings). He was obsessed with efficiency and fairness, always calculating, always optimizing.

Spec Deep Dive: The Scheduler watches for unscheduled pods through API Server. It filters nodes (predicates) - checking if nodes have enough CPU/memory, if ports are available, if labels match. Then it scores remaining nodes (priorities) based on resource balance, affinity rules, taints/tolerations. The highest-scoring node wins. This happens in milliseconds.

Controller Manager - The Anxious Parent, who constantly worried about everyone's wellbeing. He ran multiple sub-personalities: ReplicaSet Controller (counts pods obsessively), Deployment Controller (manages rollouts), Node Controller (checks if worker nodes are alive), and more. If something wasn't matching the "desired state," he'd panic and fix it immediately.

Spec Deep Dive: Controller Manager runs control loops that watch the cluster state through API Server and make changes to move current state toward desired state. Each controller is technically a separate process, but they're compiled into a single binary for efficiency. Reconciliation loop typically runs every 10 seconds.

The Worker Nodes: Where Life Happens

While Master Node and the Control Plane lived in their ivory tower, the real action happened in the **Worker Nodes** - the bustling neighborhoods of Clusterville. Each Worker Node had three essential residents:

Kubelet - The Neighborhood Mayor, a middle manager who took orders from the Control Plane and made sure they happened locally. He was on every Worker Node, constantly reporting back to API Server: "Yes, Pod-27 is running. Yes, Container-A inside Pod-27 is healthy. Yes, yes, I'm checking!"

Spec Deep Dive: Kubelet is the primary node agent. It registers the node with API Server, receives PodSpecs (Pod definitions), and ensures containers are running and healthy by working with the container runtime. It reports node and pod status back to API Server every 10 seconds (default). It also manages the mounting of volumes into pods. Kubelet doesn't manage containers not created by Kubernetes.

Kube-proxy - The Traffic Cop, who managed all the networking magic. He maintained network rules and made sure traffic going to a "Service" address got distributed properly to the actual Pods behind it. He knew every shortcut, every route, and never caused a traffic jam.

Spec Deep Dive: Kube-proxy maintains network rules on nodes using iptables, IPVS, or eBPF (depending on mode). When a Service is created, kube-proxy creates rules to forward traffic from the Service's ClusterIP to backend Pods. In iptables mode (most common), it translates Service IPs to Pod IPs. In IPVS mode, it provides better performance for large services with load balancing algorithms like round-robin, least connection, etc.

Container Runtime - The Construction Crew (Docker, containerd, CRI-O), who actually built and maintained the containers. They were the muscle, the ones doing the physical work while everyone else managed and coordinated.

Spec Deep Dive: The Container Runtime Interface (CRI) allows Kubernetes to use different runtimes. containerd is now the most common (Docker is deprecated in K8s but containerd was always doing the real work under Docker). The runtime pulls images, creates container namespaces, manages cgroups for resource limits, and reports container status back to Kubelet.

Pods: The Apartment Buildings

In Clusterville, applications didn't live in houses - they lived in **Pods**, which were like apartment buildings. But here's the weird thing: most Pods only had ONE apartment (container). Sometimes they'd have 2-3 apartments (containers) that were so closely related they needed to share the same building - same address (IP), same utilities (volumes), same mailbox (network namespace).

Pod-Frontend-X7B3 was a typical resident. She was a cozy one-bedroom (single nginx container) sitting on Worker-Node-1.

Spec Deep Dive:

```
yaml
apiVersion: v1
kind: Pod
metadata:
  name: pod-frontend-x7b3
  labels:
    app: frontend
    tier: web
spec:
  containers:
  - name: nginx
    image: nginx:1.19
    resources:
      requests:
        memory: "64Mi"
        cpu: "250m"
      limits:
        memory: "128Mi"
        cpu: "500m"
    ports:
    - containerPort: 80
```

Resources explained:

- **requests:** "I need at least this much to function" - used by Scheduler to place the Pod
- **limits:** "Don't let me use more than this" - enforced by the node
- CPU measured in millicores (m): 1000m = 1 CPU core
- Memory in bytes: Mi (mebibytes), Gi (gibibytes)

If a container exceeds memory limits: **OOMKilled** (Out of Memory Killed) - one of the most common errors you'll see. The container gets terminated and might restart based on restartPolicy.

If a container exceeds CPU limits: It gets throttled (CPU cycles restricted), but not killed. You'll see high CPU throttling metrics.

The Great Crash of Pod-Frontend

One day, Pod-Frontend-X7B3 crashed. Just stopped responding. Citizens panicked. But Controller Manager, ever vigilant, noticed immediately.

"The desired state says we should have 3 Frontend pods," he muttered, checking his notes. "But we only have 2 running. UNACCEPTABLE!"

He spun up **Pod-Frontend-K9M2** on Worker-Node-2 within seconds. Crisis averted.

But this crash revealed something important: How did anyone know the pod had crashed? Enter our first set of special characters...

Chapter 2: The Three Wise Probes and the Health Inspector

After the Great Crash, the Council of Clusterville decided they needed a better early warning system. They couldn't just wait for pods to completely die before taking action. They needed... **Probes**.

Three types of probes were appointed, each with a different philosophy:

Liveness Probe - "The Grim Reaper"

Liveness was morbid. His only job was to check if containers were alive or dead. Not "functioning well" - just "are you literally alive?" If you failed his check, you got restarted. No second chances. No negotiations.

"Are you alive?" he'd ask every 10 seconds, knocking on Pod-Frontend's door.

Spec Deep Dive:

```
yaml

livenessProbe:
  httpGet:
    path: /healthz
    port: 8080
  httpHeaders:
  - name: Custom-Header
    value: LivenessCheck
  initialDelaySeconds: 15 # Wait this long before first check
  periodSeconds: 10      # Check every 10 seconds
  timeoutSeconds: 1      # Wait 1 second for response
  failureThreshold: 3    # Fail 3 times before restarting
  successThreshold: 1    # Only need 1 success to be "alive"
```

Probe Types:

1. **httpGet**: Makes HTTP GET request, expects 200-399 status code
2. **tcpSocket**: Tries to open TCP connection
3. **exec**: Runs a command in the container, expects exit code 0

Common Liveness Errors:

- **CrashLoopBackOff:** Container keeps failing liveness checks, gets restarted, fails again, repeat. The backoff increases exponentially (0s, 10s, 20s, 40s... up to 5 minutes).
- **initialDelaySeconds too short:** Container hasn't finished starting, liveness check fails, container gets killed. Set this LONGER than your startup time.
- **Checking dependencies:** Never check external dependencies (database, APIs) in liveness. You don't want to restart your app because your database is slow. Only check internal health.

One day, Pod-Backend-Database was under heavy load. His liveness probe checked if the database could respond in 1 second. Under load, it took 3 seconds. Liveness Probe declared him DEAD and restarted him. This made things WORSE. All connections were lost, causing more load on other pods, triggering more restarts. This is called a **restart cascade** or **thundering herd problem**.

The fix: Use Readiness Probe for traffic management, save Liveness for actual deadlock/hung process detection.

Readiness Probe - "The Bouncer"

Readiness was more diplomatic. Her job was to determine if a pod was ready to receive traffic. She didn't kill anyone - she just told the Service, "Don't send traffic here yet."

This was crucial during startups, during heavy load, or during temporary issues.

Spec Deep Dive:

```
yaml

readinessProbe:
  httpGet:
    path: /ready
    port: 8080
  initialDelaySeconds: 5
  periodSeconds: 5
  failureThreshold: 3
```

When Readiness Probe failed, the pod was removed from the Service's **Endpoints** list. Kube-proxy stopped sending traffic. But the pod kept running - it got a chance to recover.

Real-world scenario: Pod-Frontend needs to load a 500MB cache file at startup. Takes 30 seconds.

- Liveness: `initialDelaySeconds: 40` (give it time to start)
- Readiness: `initialDelaySeconds: 5`, checking every 5 seconds

For the first 30 seconds, Readiness fails, so no traffic is routed. Once the cache loads and `/ready` returns 200, Readiness passes, and traffic flows in. Liveness doesn't interfere because it's not checking yet.

Startup Probe - "The Patient Mentor"

Startup was the newest probe, added because some applications were slow starters. Legacy Java apps that took 90 seconds to boot. Machine learning models that loaded huge weights at startup.

Before Startup Probe existed, you had to set huge `initialDelaySeconds` on Liveness, which meant you wouldn't detect a hung process for 90+ seconds even after the pod was running.

Startup Probe solved this: "I'll handle the slow startup. Once I succeed, Liveness and Readiness take over."

Spec Deep Dive:

```
yaml

startupProbe:
  httpGet:
    path: /startup
    port: 8080
  periodSeconds: 10
  failureThreshold: 30    # 30 failures * 10 seconds = 5 minutes max startup
                          # No initialDelaySeconds needed

livenessProbe:
  httpGet:
    path: /healthz
    port: 8080
  periodSeconds: 10
  failureThreshold: 3
```

The genius: Liveness and Readiness are DISABLED until Startup Probe succeeds once. This gives slow-starting apps up to `periodSeconds * failureThreshold` time to start (5 minutes in this example), but once started, Liveness catches issues within 30 seconds (10s * 3 failures).

The Tale of the Misconfigured Probe

Pod-Payment-Gateway had a developer who thought they were being clever:

```
yaml
```

```
livenessProbe:
  httpGet:
    path: /health
    port: 8080
  periodSeconds: 1    # Check every second! So responsive!
  timeoutSeconds: 5    # Give it 5 seconds to respond
  failureThreshold: 1  # One strike and you're out!
```

Disaster struck. A minor network blip caused one health check to fail. Pod-Payment-Gateway was immediately restarted. In the middle of processing payments. Transactions lost. Customers angry. Developers sad.

Probe Best Practices:

1. **Liveness:** Set `failureThreshold` to 3+. Give benefit of doubt.
2. **Readiness:** Can be more sensitive, but not too aggressive
3. **periodSeconds:** 10 seconds is usually fine. Don't check every second (wasteful).
4. **timeoutSeconds:** 1-2 seconds unless you have good reason
5. **Separate endpoints:** `/healthz` (liveness) vs `/ready` (readiness) with different logic
6. **Never check dependencies in Liveness:** Only check internal health

Chapter 3: The Love Triangle - Ingress, Services, and the DNS Matchmaker

Pod-Frontend and Pod-Backend had a problem. They were in love (they needed to communicate), but they kept changing addresses. Every time a pod restarted, it got a new IP. How could they find each other?

Enter **Service** - the matchmaker who provided stable identities.

Service: The Stable Address

Service-Frontend was like a phone number that always worked, even if you changed phones. She had a fixed ClusterIP (internal IP address) and a name (`frontend-service`). When anyone wanted to talk to Frontend pods, they called the Service, and she'd route them to healthy pods.

Spec Deep Dive:

```
yaml
```

```
apiVersion: v1
kind: Service
metadata:
  name: frontend-service
spec:
  selector:
    app: frontend      # Matches pods with this label
  ports:
    - protocol: TCP
      port: 80          # Service port
      targetPort: 8080  # Container port
  type: ClusterIP      # Default, internal only
```

How it works:

1. Service is created, API Server assigns it a ClusterIP (e.g., 10.96.0.15)
2. CoreDNS (the DNS service) creates a record: `frontend-service.default.svc.cluster.local` → 10.96.0.15
3. Kube-proxy on each node creates iptables rules: Traffic to 10.96.0.15:80 → distribute to Pod IPs on port 8080
4. When Pod-Backend wants to reach Frontend, it just calls `http://frontend-service` (DNS resolves it)

Service Types:

ClusterIP (default): Only accessible within cluster

```
yaml

type: ClusterIP
```

NodePort: Accessible from outside cluster on a specific port (30000-32767) on ANY node

```
yaml

type: NodePort
ports:
  - port: 80
    targetPort: 8080
    nodePort: 30080  # Optional, auto-assigned if omitted
```

Access via `http://<any-node-ip>:30080`. Traffic goes to Service, then to Pods (might be on different nodes).

LoadBalancer: Creates an external load balancer (cloud provider specific)

```
yaml
```


type: LoadBalancer

On AWS, this creates an ELB/ALB. On GCP, a Load Balancer. Gets an external IP.

ExternalName: Maps to an external DNS name (no proxying)

yaml

type: ExternalName

externalName: my-database.example.com

Ingress: The Grand Gateway

But LoadBalancer Services had a problem: Each one created a separate (expensive) cloud load balancer. For 10 services, you paid for 10 load balancers.

Ingress was the solution - a sophisticated host/path-based router. One load balancer, many services.

Spec Deep Dive:

yaml

```

apiVersion: networking.k8s.io/v1
kind: Ingress
metadata:
  name: main-ingress
  annotations:
    nginx.ingress.kubernetes.io/rewrite-target: /
    cert-manager.io/cluster-issuer: "letsencrypt-prod"
spec:
  ingressClassName: nginx  # Which controller to use
  tls:
  - hosts:
    - myapp.example.com
    secretName: myapp-tls  # Where TLS cert is stored
  rules:
  - host: myapp.example.com
    http:
      paths:
      - path: /api
        pathType: Prefix
        backend:
          service:
            name: backend-service
            port:
              number: 80
      - path: /
        pathType: Prefix
        backend:
          service:
            name: frontend-service
            port:
              number: 80

```

How Ingress Works:

1. You deploy an **Ingress Controller** (nginx-ingress, traefik, HAProxy, etc.) - this is a pod that watches for Ingress resources
2. The Ingress Controller creates a LoadBalancer Service (one for all Ingresses)
3. You create Ingress resources defining rules
4. Ingress Controller configures itself based on these rules
5. Traffic flow: Internet → LoadBalancer → Ingress Controller Pod → Service → Backend Pods

pathType Options:

- **Prefix:** Matches /api, /api/, /api/users, /api/users/123

- **Exact:** Only matches exact path
- **ImplementationSpecific:** Depends on controller

Common Ingress Issues:

Issue 1: 404 Not Found despite correct Ingress config

- Check: Is the Ingress Controller installed? `kubectl get pods -n ingress-nginx`
- Check: Does the Service exist and have Endpoints? `kubectl get endpoints backend-service`
- Check: Is the path rewrite correct? If backend expects `/` but gets `/api/users`, it fails.

Issue 2: Default backend - 404

- The Ingress needs a default backend for unmatched requests
- Usually you deploy a default-backend service that shows a nice 404 page

Issue 3: TLS/HTTPS not working

- Check: Does the Secret exist? `kubectl get secret myapp-tls`
- Check: Is cert-manager running and issuing certificates?
- Check: Are DNS records pointing to the LoadBalancer IP?

The Great Ingress Saga

One day, Traffic Controller wanted to update the Ingress rules. He wanted `myapp.example.com/v2` to route to the new Backend-V2 service.

```
yaml
- path: /v2
  pathType: Prefix
  backend:
    service:
      name: backend-v2-service
      port:
        number: 80
```

He applied the config. Immediately, half the traffic started failing!

What happened?: Path priority matters. His config had:

```
yaml
```

```
- path: /      # This matches EVERYTHING
...
- path: /v2    # This never gets hit!
...
```

Fix: Most specific paths first:

```
yaml

- path: /v2
...
- path: /api
...
- path: /
...
```

Another common gotcha - rewrite-target:

Without rewrite:

- Request: `GET /api/users`
- Ingress forwards to backend: `GET /api/users`
- Backend expects `/users` (doesn't know about `/api` prefix)
- Result: 404

With rewrite:

```
yaml

annotations:
  nginx.ingress.kubernetes.io/rewrite-target: /$2
...
path: /api(/$)(.*)
```

- Request: `GET /api/users`
- Ingress forwards: `GET /users`
- Backend happy!

Chapter 4: The Vault Heist - Secrets, ConfigMaps, and the PV/PVC Storage Wars

Every application in Clusterville had secrets. Database passwords, API keys, TLS certificates. Where should they store them?

ConfigMaps and Secrets: The Dynamic Duo

ConfigMap was the town crier - information everyone could know. Configuration settings, non-sensitive environment variables.

Secret was the spy - information that needed to be protected (though not as well as everyone thought - secrets in Kubernetes are only base64 encoded, not encrypted at rest by default).

ConfigMap Spec:

```
yaml

apiVersion: v1
kind: ConfigMap
metadata:
  name: app-config
data:
  database_url: "postgres://db.example.com:5432"
  log_level: "info"
  feature_flags.json: |
    {
      "new_feature": true,
      "beta_mode": false
    }
```

Secret Spec:

```
yaml

apiVersion: v1
kind: Secret
metadata:
  name: app-secrets
type: Opaque
data:
  db_password: cGFzc3dvcmQxMjM= # base64 encoded
  api_key: c2VjcmlV0a2V5NDU2    # base64 encoded
```

Using them in Pods:

As environment variables:

```
yaml
```

```
spec:
  containers:
  - name: app
    env:
    - name: DATABASE_URL
      valueFrom:
        configMapKeyRef:
          name: app-config
          key: database_url
    - name: DB_PASSWORD
      valueFrom:
        secretKeyRef:
          name: app-secrets
          key: db_password
```

As mounted files:

```
yaml

spec:
  containers:
  - name: app
    volumeMounts:
    - name: config-volume
      mountPath: /etc/config
    - name: secret-volume
      mountPath: /etc/secrets
      readOnly: true
  volumes:
  - name: config-volume
    configMap:
      name: app-config
  - name: secret-volume
    secret:
      secretName: app-secrets
```

Now `/etc/config/database_url` contains the database URL, and `/etc/secrets/db_password` contains the password.

Secret Types:

- **Opaque:** Generic secret (default)
- **kubernetes.io/service-account-token:** Service account token
- **kubernetes.io/dockerconfigjson:** Docker registry credentials

```
yaml
```

```
type: kubernetes.io/dockerconfigjson
data:
  .dockerconfigjson: <base64-encoded-docker-config>
```

- **kubernetes.io/tls**: TLS certificate and key

```
yaml

type: kubernetes.io/tls
data:
  tls.crt: <base64-cert>
  tls.key: <base64-key>
```

Enter HashiCorp Vault: The Real Spy

But Secret wasn't really secure. Anyone with access to etcd could read secrets. Anyone with kubectl access could decode them (`kubectl get secret app-secrets -o jsonpath='{.data.db_password}' | base64 -d`).

The town hired **Vault** - a professional secret keeper who lived outside Clusterville (external service) but worked closely with the residents.

Vault's Approach:

1. Secrets are encrypted at rest in Vault (not Kubernetes)
2. Applications authenticate to Vault to retrieve secrets
3. Secrets can be dynamically generated (like database credentials with TTL)
4. Full audit logging of who accessed what secret when

Integration Methods:

Method 1: Vault Agent Sidecar

```
yaml

spec:
  containers:
    - name: app
    ...
    - name: vault-agent
      image: vault:1.12.0
      # Vault agent authenticates, fetches secrets, writes to shared volume
```

Method 2: Vault CSI Provider Uses the Secrets Store CSI driver to mount Vault secrets as volumes:

```
yaml

spec:
  volumes:
    - name: secrets-store
      csi:
        driver: secrets-store.csi.k8s.io
        readOnly: true
        volumeAttributes:
          secretProviderClass: "vault-database"
```

Method 3: External Secrets Operator Syncs secrets from Vault to Kubernetes Secrets:

```
yaml

apiVersion: external-secrets.io/v1beta1
kind: ExternalSecret
metadata:
  name: vault-example
spec:
  refreshInterval: 1h
  secretStoreRef:
    name: vault-backend
    kind: SecretStore
  target:
    name: app-secrets
  data:
    - secretKey: db_password
      remoteRef:
        key: database/config
        property: password
```

The Heist Gone Wrong:

One day, Pod-API-Gateway needed a database password. The developer hardcoded it:

```
yaml

env:
  - name: DB_PASSWORD
    value: "SuperSecret123" # DON'T DO THIS!
```

Anyone running `kubectl describe pod` could see it. Anyone with access to the Git repo where this YAML was stored could see it. It got checked into version control history forever.

The right way:

1. Store in Vault or external secret manager
 2. Use ExternalSecrets or CSI driver to bring into K8s
 3. Reference in Pod spec as a Secret
 4. Never commit secrets to Git (use tools like git-secrets, detect-secrets)
-

PersistentVolumes and PersistentVolumeClaims: The Storage Real Estate Market

Pod-Database had a problem: Every time he restarted, his data vanished! Pods are ephemeral - their storage is temporary.

He needed **PersistentVolume** - real estate that existed independently of pods.

The Storage Hierarchy:

1. **PersistentVolume (PV)**: The actual storage (like a warehouse)
2. **PersistentVolumeClaim (PVC)**: A request for storage (like renting warehouse space)
3. **StorageClass**: The real estate agency that can provision PVs on-demand

PersistentVolume Spec:

```
yaml

apiVersion: v1
kind: PersistentVolume
metadata:
  name: pv-database
spec:
  capacity:
    storage: 10Gi
  volumeMode: Filesystem
  accessModes:
    - ReadWriteOnce      # RWO: one node can mount read-write
  persistentVolumeReclaimPolicy: Retain
  storageClassName: fast-ssd
  hostPath:              # DON'T USE IN PRODUCTION
    path: /data/pv-database
```

AccessModes:

- **ReadWriteOnce (RWO)**: One node mounts read-write (most common)
- **ReadOnlyMany (ROX)**: Many nodes mount read-only
- **ReadWriteMany (RWX)**: Many nodes mount read-write (needs NFS, or cloud provider support)

ReclaimPolicy:

- **Retain:** PV remains after PVC deletion (manual cleanup required)
- **Delete:** PV automatically deleted when PVC deleted (data lost!)
- **Recycle:** Deprecated (basic scrub `rm -rf /volume/*`)

PersistentVolumeClaim Spec:

```
yaml
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: pvc-database
spec:
  accessModes:
    - ReadWriteOnce
  resources:
    requests:
      storage: 5Gi
  storageClassName: fast-ssd
```

Kubernetes matches PVCs to PVs based on:

1. StorageClass name
2. Access modes compatibility
3. Capacity (PV must be \geq PVC request)

Using in Pod:

```
yaml
spec:
  containers:
    - name: postgres
      volumeMounts:
        - mountPath: /var/lib/postgresql/data
          name: database-storage
  volumes:
    - name: database-storage
      persistentVolumeClaim:
        claimName: pvc-database
```

StorageClass: Dynamic Provisioning

Creating PVs manually is tedious. **StorageClass** enables dynamic provisioning.

yaml

```
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: fast-ssd
provisioner: kubernetes.io/aws-ebs
parameters:
  type: gp3
  iopsPerGB: "50"
  fsType: ext4
allowVolumeExpansion: true
volumeBindingMode: WaitForFirstConsumer
```

volumeBindingMode:

- **Immediate:** PV created immediately when PVC is created
- **WaitForFirstConsumer:** PV created when a pod using the PVC is scheduled (better for topology constraints)

Now when you create a PVC with `storageClassName: fast-ssd`, Kubernetes automatically:

1. Calls the AWS EBS provisioner
2. Creates an EBS volume
3. Creates a PV pointing to that EBS volume
4. Binds the PVC to the PV

Common PV/PVC Issues:

Issue 1: PVC stuck in Pending

- Check: `kubectl describe pvc pvc-database`
- Common causes:
 - No PV available matching the request
 - StorageClass provisioner not installed
 - Insufficient capacity in the StorageClass
 - Access mode not supported

Issue 2: Pod stuck in ContainerCreating

- Check: `kubectl describe pod database-pod`
- Often: "Unable to attach or mount volumes"
- Causes:

- PVC doesn't exist
- PV is bound to a different node (RWO) and pod scheduled on wrong node
- Permissions issue (fsGroup, runAsUser mismatch)

Issue 3: Data persistence not working

- Check: Is volumeMount path correct?
- Check: Is the app actually writing to that path?
- For databases, check subPath usage:

```
yaml

volumeMounts:
- mountPath: /var/lib/postgresql/data
  name: database-storage
  subPath: postgres    # Creates /postgres subdirectory
```

The Great Data Loss of Node-3:

Pod-Database was running happily on Node-3 with a hostPath PV (local storage on the node). Then Node-3 crashed. Pod-Database was rescheduled to Node-4. But the data was on Node-3!

Lesson: **Never use hostPath for production data.** Use:

- **Cloud provider volumes:** AWS EBS, GCP Persistent Disk, Azure Disk
- **Network storage:** NFS, CephFS, GlusterFS
- **Distributed storage:** Rook/Ceph, Longhorn, OpenEBS

These allow pods to move between nodes while keeping data accessible.

Chapter 5: The Politics of Placement - Affinity, Taints, and Tolerations

Not all pods could live anywhere. Some were picky. Some nodes were exclusive. This created Clusterville's most complex social dynamics.

Node Affinity: "I Want to Live Near..."

Pod-GPU-Trainer needed GPUs for machine learning. He couldn't just live anywhere. He used **nodeSelector** (the simple way) or **nodeAffinity** (the sophisticated way).

Simple nodeSelector:

```
yaml
```

```
spec:
  nodeSelector:
    gpu: "true"
    instance-type: "g4dn.xlarge"
```

This is inflexible - either a node matches ALL labels or the pod won't schedule.

Advanced nodeAffinity:

```
yaml

spec:
  affinity:
    nodeAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        nodeSelectorTerms:
          - matchExpressions:
              - key: gpu
                operator: In
                values:
                  - "true"
            - key: gpu-type
                operator: In
                values:
                  - nvidia-tesla-v100
                  - nvidia-tesla-a100
          preferredDuringSchedulingIgnoredDuringExecution:
            - weight: 100
              preference:
                matchExpressions:
                  - key: zone
                    operator: In
                    values:
                      - us-east-1a
            - weight: 50
              preference:
                matchExpressions:
                  - key: instance-type
                    operator: In
                    values:
                      - g4dn.2xlarge
```

requiredDuringScheduling...: MUST match or pod won't schedule **preferredDuringScheduling...:** TRY to match, but not required (weighted preferences) **...IgnoredDuringExecution:** If node labels change after pod is running, pod stays (not evicted)

Operators:

- **In:** Label value must be in the list
 - **NotIn:** Label value must NOT be in the list
 - **Exists:** Label key must exist (any value)
 - **DoesNotExist:** Label key must not exist
 - **Gt:** Greater than (numeric)
 - **Lt:** Less than (numeric)
-

Pod Affinity: "I Want to Live Near My Friends"

Pod-Cache wanted to be on the same node as **Pod-API** to reduce network latency. This is **podAffinity**.

```
yaml

spec:
  affinity:
    podAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        - labelSelector:
            matchExpressions:
              - key: app
                operator: In
                values:
                  - api-server
          topologyKey: kubernetes.io/hostname
```

topologyKey: Defines the "zone" for affinity

- `kubernetes.io/hostname`: Same node
- `topology.kubernetes.io/zone`: Same availability zone
- `topology.kubernetes.io/region`: Same region

Pod Anti-Affinity: "I DON'T want to live near..."

```
yaml
```

```
spec:
  affinity:
    podAntiAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
      - labelSelector:
          matchExpressions:
            - key: app
              operator: In
              values:
                - frontend
        topologyKey: kubernetes.io/hostname
```

This ensures no two frontend pods end up on the same node (high availability).

Real-world example - Database cluster:

```
yaml

# PostgreSQL primary and replicas should be on DIFFERENT nodes
spec:
  affinity:
    podAntiAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
      - labelSelector:
          matchLabels:
            app: postgres
        topologyKey: kubernetes.io/hostname
    # But should be in the SAME availability zone
    podAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
      - weight: 100
        podAffinityTerm:
          labelSelector:
            matchLabels:
              app: postgres
          topologyKey: topology.kubernetes.io/zone
```

Taints and Tolerations: "You Shall Not Pass... Unless..."

Taints were like "Keep Out" signs on nodes. **Tolerations** were like special passes that let certain pods ignore those signs.

Node-GPU-1 had expensive GPUs. He didn't want random pods wasting his resources:

```
bash
```

```
kubectl taint nodes node-gpu-1 gpu=true:NoSchedule
```

Now regular pods couldn't schedule there. But Pod-GPU-Trainer had a toleration:

```
yaml

spec:
  tolerations:
  - key: gpu
    operator: Equal
    value: "true"
    effect: NoSchedule
```

He could schedule on Node-GPU-1!

Taint Effects:

- **NoSchedule:** Don't schedule new pods (existing pods stay)
- **PreferNoSchedule:** Try to avoid scheduling here (soft)
- **NoExecute:** Don't schedule AND evict existing pods that don't tolerate

NoExecute Example:

```
bash

kubectl taint nodes node-3 maintenance=true:NoExecute
```

All pods without a matching toleration are immediately evicted!

```
yaml

tolerations:
- key: maintenance
  operator: Equal
  value: "true"
  effect: NoExecute
  tolerationSeconds: 300  # Stay for 5 minutes, then evict
```

Common Built-in Taints:

Kubernetes automatically adds these:


```
node.kubernetes.io/not-ready:NoExecute    # Node not ready
node.kubernetes.io/unreachable:NoExecute  # Node unreachable
node.kubernetes.io/disk-pressure:NoSchedule # Low disk space
node.kubernetes.io/memory-pressure:NoSchedule # Low memory
node.kubernetes.io/pid-pressure:NoSchedule # Too many processes
node.kubernetes.io/network-unavailable:NoSchedule
```

DaemonSets automatically tolerate these (except PID pressure) so they run everywhere.

The Great Scheduling Disaster

One day, the admins created a new node pool for batch jobs:

```
bash

kubectl taint nodes batch-node-1 workload=batch:NoSchedule
kubectl taint nodes batch-node-2 workload=batch:NoSchedule
```

Then they deployed a critical database without checking:

```
yaml

apiVersion: apps/v1
kind: StatefulSet
metadata:
  name: postgres
spec:
  replicas: 3
  # No tolerations specified!
```

The postgres pods stayed in **Pending** state forever! Why? ALL nodes were either:

1. Regular nodes (full, no capacity)
2. Batch nodes (tainted, postgres couldn't tolerate)

Error in events:

```
0/10 nodes available: 5 node(s) had insufficient CPU,
5 node(s) had taint {workload: batch}, that the pod didn't tolerate.
```

The Fix:

```
yaml
```

```
spec:
  template:
    spec:
      tolerations:
      - key: workload
        operator: Equal
        value: batch
        effect: NoSchedule
      # OR use nodeSelector to force it to regular nodes
      nodeSelector:
        workload: general
```

Common Affinity Errors

Error 1: "0/5 nodes available: 5 node(s) didn't match pod affinity rules"

Cause: Required podAffinity can't be satisfied

```
yaml

podAffinity:
  requiredDuringScheduling...:
  - labelSelector:
      matchLabels:
        app: non-existent-app  # No such pod exists!
```

Fix: Change to `preferredDuringScheduling` or ensure the target pods exist first.

Error 2: "0/5 nodes available: 5 node(s) didn't match pod anti-affinity rules"

Cause: You have 3 replicas with required anti-affinity on hostname, but only 2 nodes!

```
yaml

replicas: 3  # Need 3 different nodes
podAntiAffinity:
  requiredDuringScheduling...:
    topologyKey: kubernetes.io/hostname  # But only 2 nodes!
```

Fix: Add more nodes, reduce replicas, or change to `preferredDuringScheduling`.

Error 3: Pods all scheduling on same node despite anti-affinity

Cause: Anti-affinity is `preferred`, not `required`. Under resource pressure, scheduler ignores preferences.

Fix: Use `requiredDuringScheduling` if you really need it enforced.

The Perfect Storm: Combining Everything

Here's a real-world example combining multiple concepts:

```
yaml
```

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: ml-training
spec:
  replicas: 3
  template:
    spec:
      # Need GPU nodes
      nodeSelector:
        gpu: "true"

      # Prefer specific GPU types
      affinity:
        nodeAffinity:
          preferredDuringSchedulingIgnoredDuringExecution:
            - weight: 100
              preference:
                matchExpressions:
                  - key: gpu-type
                    operator: In
                    values:
                      - nvidia-a100

            # Spread across different nodes (HA)
          podAntiAffinity:
            requiredDuringSchedulingIgnoredDuringExecution:
              - labelSelector:
                  matchLabels:
                    app: ml-training
                topologyKey: kubernetes.io/hostname

            # But prefer same zone (reduce inter-zone traffic)
          podAffinity:
            preferredDuringSchedulingIgnoredDuringExecution:
              - weight: 50
                podAffinityTerm:
                  labelSelector:
                    matchLabels:
                      app: ml-training
                  topologyKey: topology.kubernetes.io/zone

            # Tolerate GPU node taints
          tolerations:
            - key: gpu
              operator: Equal
```

```
value: "true"
effect: NoSchedule
```

```
# Resource limits
```

```
containers:
```

```
- name: trainer
```

```
resources:
```

```
requests:
```

```
  nvidia.com/gpu: 1
```

```
  memory: 16Gi
```

```
  cpu: 4
```

```
limits:
```

```
  nvidia.com/gpu: 1
```

```
  memory: 32Gi
```

```
  cpu: 8
```

```
# Health checks
```

```
livenessProbe:
```

```
  httpGet:
```

```
    path: /health
```

```
    port: 8080
```

```
  initialDelaySeconds: 60
```

```
  periodSeconds: 10
```

```
readinessProbe:
```

```
  httpGet:
```

```
    path: /ready
```

```
    port: 8080
```

```
  initialDelaySeconds: 10
```

```
  periodSeconds: 5
```

This deployment:

1. ☒ Only schedules on GPU nodes
 2. ☒ Prefers A100 GPUs if available
 3. ☒ Spreads 3 replicas across 3 different nodes
 4. ☒ But tries to keep them in the same availability zone
 5. ☒ Can schedule on tainted GPU nodes
 6. ☒ Has proper resource requests/limits
 7. ☒ Has health checks configured
-

Epilogue: The Day Clusterville Fell (And How etcd Saved It)

It was a Tuesday morning when disaster struck. A rogue script ran `kubectl delete all --all` in production. Within seconds, all Deployments, Services, Pods vanished. Chaos erupted. Applications went dark. The citizens of Clusterville panicked.

But **etcd**, the Memory Keeper, remained calm. "I remember everything," she whispered from her distributed vault. "Everything."

This is the story of backup, disaster recovery, and the resurrection of Clusterville.

Understanding What Lives Where: The Clusterville Census

Before we talk about backup, you need to understand what lives where in Clusterville:

In etcd's Memory (the source of truth):

- All Kubernetes resources: Pods, Deployments, Services, ConfigMaps, Secrets, Ingresses, PVs, PVCs, RBAC rules, etc.
- Cluster state: Which pods are running where, which nodes are healthy
- Configuration: Cluster settings, API server configuration

NOT in etcd:

- **Application data** in PersistentVolumes (lives on storage backend - EBS, NFS, etc.)
- **Container images** (lives in container registry - Docker Hub, ECR, GCR)
- **Logs** (if using external logging like ELK, Loki)
- **Metrics** (if using external monitoring like Prometheus)

In your laptop/CI system:

- **Kubeconfig file** (`~/.kube/config`) - Your "passport" to Clusterville
- **YAML manifests** (if you use GitOps) - Your "blueprints"

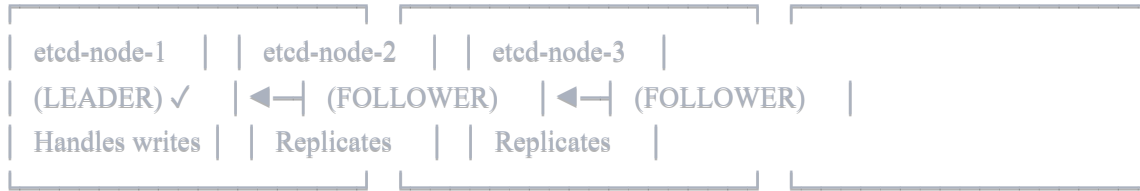
Think of it like a city:

- **etcd** = City Hall records (who owns what, all official documents)
 - **PVs** = Actual houses and their contents (your data)
 - **Images** = Building materials warehouse (reusable components)
 - **Kubeconfig** = Your ID card (proves you're authorized to enter)
 - **Manifests** = Architectural blueprints (how to rebuild if needed)
-

The Anatomy of etcd: Inside the Memory Keeper's Mind

etcd is a distributed key-value store. Let's see what "distributed" really means:

etcd Cluster (3 nodes for quorum):



Raft Consensus Protocol:

1. One node is elected LEADER
2. All write requests go to the leader
3. Leader replicates to followers
4. Once MAJORITY (2 out of 3) acknowledge, write is committed
5. If leader fails, followers elect a new leader (30-50ms)

Why odd numbers?

- 3 nodes: Can tolerate 1 failure (need 2 for quorum)
- 5 nodes: Can tolerate 2 failures (need 3 for quorum)
- 4 nodes: Can still only tolerate 1 failure (need 3 for quorum) - so 4 nodes is wasteful, use 3 or 5

What happens if etcd loses quorum?

Before: 3 nodes, quorum = 2

Node-1 ✓ (leader)

Node-2 ✓

Node-3 ✗ (dies)

Status: HEALTHY (2/3 still works)

Then: Node-2 also dies

Node-1 ✓ (alone)

Node-2 ✗

Node-3 ✗

Status: NO QUORUM! Cluster READ-ONLY!

Without quorum, Kubernetes enters **read-only mode**:

- ✗ Can't create new pods
- ✗ Can't delete resources

- ❌ Can't update anything
 - ✅ Existing pods keep running (kubelet doesn't need etcd for running pods)
 - ✅ Services keep working (kube-proxy has cached rules)
 - But you're flying blind - can't make any changes!
-

The Great Backup: Saving Clusterville's Soul

etcd stores data as a **snapshot** - a point-in-time backup of all cluster state.

Backup Strategy 1: etcdctl Snapshot

```
bash
```

```
# Install etcdctl (etcd client)
```

```
ETCD_VER=v3.5.9
```

```
wget https://github.com/etcd-io/etcd/releases/download/${ETCD_VER}/etcd-${ETCD_VER}-linux-amd64.tar.gz
```

```
tar xzf etcd-${ETCD_VER}-linux-amd64.tar.gz
```

```
sudo mv etcd-${ETCD_VER}-linux-amd64/etcdctl /usr/local/bin/
```

```
# Take a snapshot (run on etcd node or master node)
```

```
ETCDCTL_API=3 etcdctl snapshot save /backup/etcd-snapshot-$(date +%Y%m%d-%H%M%S).db \
```

```
--endpoints=https://127.0.0.1:2379 \
```

```
--cacert=/etc/kubernetes/pki/etcd/ca.crt \
```

```
--cert=/etc/kubernetes/pki/etcd/server.crt \
```

```
--key=/etc/kubernetes/pki/etcd/server.key
```

```
# Verify snapshot
```

```
ETCDCTL_API=3 etcdctl snapshot status /backup/etcd-snapshot-20251007-143000.db --write-out=table
```

Output:

```
+-----+-----+-----+-----+
| HASH | REVISION | TOTAL KEYS | TOTAL SIZE |
+-----+-----+-----+-----+
| 8e9b0f2a | 123456 | 2847 | 89 MB |
+-----+-----+-----+-----+
```

What's in this snapshot?

- Every Pod definition
- Every Service configuration
- Every ConfigMap and Secret
- Every Deployment, StatefulSet, DaemonSet

- Every Ingress rule
- Every PV and PVC definition (but NOT the actual data)
- RBAC roles and bindings
- Custom Resource Definitions (CRDs)
- Literally EVERYTHING in `kubectrl get all --all-namespaces`

Automated Backup with CronJob:

yaml

apiVersion: batch/v1

kind: CronJob

metadata:

name: etcd-backup

namespace: kube-system

spec:

schedule: "0 2 * * *" *# Daily at 2 AM*

jobTemplate:

spec:

template:

spec:

hostNetwork: true

nodeName: master-node-1 *# Run on master with etcd*

containers:

- name: etcd-backup

image: k8s.gcr.io/etcd:3.5.9

command:

- /bin/sh

- -c

- |

ETCDCTL_API=3 etcdctl snapshot save /backup/etcd-\$(date +%Y%m%d-%H%M%S).db \

--endpoints=https://127.0.0.1:2379 \

--cacert=/etc/kubernetes/pki/etcd/ca.crt \

--cert=/etc/kubernetes/pki/etcd/server.crt \

--key=/etc/kubernetes/pki/etcd/server.key

Upload to S3

aws s3 cp /backup/etcd-*.db s3://my-etcd-backups/

Keep only last 7 days locally

find /backup -name "etcd-*.db" -mtime +7 -delete

volumeMounts:

- name: etcd-certs

mountPath: /etc/kubernetes/pki/etcd

readOnly: true

- name: backup-dir

mountPath: /backup

volumes:

- name: etcd-certs

hostPath:

path: /etc/kubernetes/pki/etcd

- name: backup-dir

hostPath:

path: /var/backups/etcd

restartPolicy: OnFailure

Disaster Strikes: The Great Deletion

Back to our disaster. Someone ran:

```
bash  
  
kubectl delete all --all --all-namespaces
```

Within seconds:

- 247 Pods terminated
- 83 Services deleted
- 42 Deployments gone
- 28 ConfigMaps vanished
- 19 Ingresses disappeared

But wait! `kubectl delete all` doesn't actually delete "all". It deletes:

- Pods, Services, Deployments, ReplicaSets, StatefulSets, DaemonSets, Jobs, CronJobs

It does NOT delete:

- ☒ ConfigMaps (not in "all")
- ☒ Secrets (not in "all")
- ☒ PersistentVolumes/PersistentVolumeClaims (not in "all")
- ☒ Ingresses (not in "all")
- ☒ RBAC (not in "all")
- ☒ Namespaces (not in "all")

So the deletion was bad, but not catastrophic. However, let's pretend it was worse - someone deleted EVERYTHING including PVCs, ConfigMaps, Secrets.

The Resurrection: Restoring from etcd Backup

Step 1: Stop the Kubernetes Control Plane

```
bash
```

On master node(s)

```
sudo systemctl stop kube-apiserver
```

```
sudo systemctl stop kube-controller-manager
```

```
sudo systemctl stop kube-scheduler
```

Or if using static pods (kubeadm)

```
sudo mv /etc/kubernetes/manifests /etc/kubernetes/manifests.backup
```

This stops API server; controller-manager, scheduler

Step 2: Stop etcd

bash

```
sudo systemctl stop etcd
```

Or for kubeadm

```
sudo mv /etc/kubernetes/manifests/etcd.yaml /etc/kubernetes/
```

Step 3: Backup Current (Broken) etcd Data

bash

```
sudo mv /var/lib/etcd /var/lib/etcd.broken
```

Step 4: Restore from Snapshot

bash

```
ETCDCTL_API=3 etcdctl snapshot restore /backup/etcd-snapshot-20251007-020000.db \
```

```
--data-dir=/var/lib/etcd \
```

```
--name=master-node-1 \
```

```
--initial-cluster=master-node-1=https://10.0.1.10:2380 \
```

```
--initial-cluster-token=etcd-cluster-1 \
```

```
--initial-advertise-peer-urls=https://10.0.1.10:2380
```

For a 3-node etcd cluster, you need to restore on ALL nodes:

bash

On master-node-1:

```
ETCDCTL_API=3 etcdctl snapshot restore /backup/etcd-snapshot.db \
--data-dir=/var/lib/etcd \
--name=master-node-1 \
--initial-cluster=master-node-1=https://10.0.1.10:2380,master-node-2=https://10.0.1.11:2380,master-node-3=https://10.0.1.12:2380 \
--initial-cluster-token=etcd-cluster-1 \
--initial-advertise-peer-urls=https://10.0.1.10:2380
```

On master-node-2:

```
ETCDCTL_API=3 etcdctl snapshot restore /backup/etcd-snapshot.db \
--data-dir=/var/lib/etcd \
--name=master-node-2 \
--initial-cluster=master-node-1=https://10.0.1.10:2380,master-node-2=https://10.0.1.11:2380,master-node-3=https://10.0.1.12:2380 \
--initial-cluster-token=etcd-cluster-1 \
--initial-advertise-peer-urls=https://10.0.1.11:2380
```

On master-node-3:

```
ETCDCTL_API=3 etcdctl snapshot restore /backup/etcd-snapshot.db \
--data-dir=/var/lib/etcd \
--name=master-node-3 \
--initial-cluster=master-node-1=https://10.0.1.10:2380,master-node-2=https://10.0.1.11:2380,master-node-3=https://10.0.1.12:2380 \
--initial-cluster-token=etcd-cluster-1 \
--initial-advertise-peer-urls=https://10.0.1.12:2380
```

Step 5: Fix Permissions

bash

```
sudo chown -R etcd:etcd /var/lib/etcd
```

Or for kubeadm

```
sudo chown -R root:root /var/lib/etcd
```

Step 6: Start etcd

bash

```
sudo systemctl start etcd
```

Or for kubeadm

```
sudo mv /etc/kubernetes/etcd.yaml /etc/kubernetes/manifests/
```

Step 7: Start Control Plane

bash

```
# For kubeadm
sudo mv /etc/kubernetes/manifests.backup/* /etc/kubernetes/manifests/

# Wait for pods to come up
watch kubectl get pods -n kube-system
```

Step 8: Verify

```
bash

kubectl get all --all-namespaces
kubectl get configmaps --all-namespaces
kubectl get secrets --all-namespaces
kubectl get pvc --all-namespaces
kubectl get ingress --all-namespaces
```

🎉 **Clusterville is back!** All resources from the backup time are restored!

What Got Restored and What Didn't: The Fine Print

✅ Restored from etcd backup:

- All Kubernetes resources (Pods, Deployments, Services, etc.)
- ConfigMaps - YES! They're stored in etcd
- Secrets - YES! They're stored in etcd
- Ingress rules - YES! They're stored in etcd
- PV and PVC definitions - YES! But...

❌ NOT restored (lives elsewhere):

- **Actual data in PersistentVolumes:** If someone deleted the PVC and the reclaim policy was "Delete", the underlying storage (EBS volume, NFS share) might be deleted. etcd restoration brings back the PVC definition, but the data is gone if the storage was deleted.
- **Container images:** Not in etcd. If your registry is down, pods won't start even though definitions are restored.
- **Application state:** If your app had in-memory state, it's gone.
- **Logs:** Unless stored in PVs or external system.

⚠️ Time Travel Paradox:

Your backup was from 2 AM. It's now 3 PM. You restore. What happens?

- Resources created between 2 AM - 3 PM: **GONE** (not in backup)

- Resources deleted between 2 AM - 3 PM: **BACK** (they exist in backup)
- ConfigMaps/Secrets updated between 2 AM - 3 PM: **OLD VALUES** (backup has old values)

This is why backup frequency matters!

The ConfigMap Conundrum: What If You Only Need One Thing?

Sometimes you don't need to restore the entire cluster. Maybe you just accidentally deleted a critical ConfigMap:

```
bash

kubectl delete configmap app-config -n production # Oops!
```

Option 1: Restore from Git (if you practice GitOps)

```
bash

git log --all -- production/app-config.yaml
# Find the file
git checkout <commit-hash> -- production/app-config.yaml
kubectl apply -f production/app-config.yaml
```

This is why **GitOps** (storing all YAML in Git) is powerful - Git is your backup!

Option 2: Extract from etcd backup without full restore

You can't easily extract one resource from an etcd snapshot, but you can restore it to a temporary location and query it:

```
bash

# Restore snapshot to temp location
ETCDCTL_API=3 etcdctl snapshot restore /backup/etcd-snapshot.db \
  --data-dir=/tmp/etcd-restore

# Start a temporary etcd instance
etcd --data-dir=/tmp/etcd-restore --listen-client-urls=http://localhost:12379 --advertise-client-urls=http://localhost:12379 &

# Query it
ETCDCTL_API=3 etcdctl get /registry/configmaps/production/app-config \
  --endpoints=http://localhost:12379 --print-value-only > /tmp/configmap.yaml

# Convert from etcd's internal format (it's protobuf, tricky!)
# Easier: Use a tool like etcd-browser or restore to a test cluster and export
```

Honestly, this is complex. Better approaches:

1. **Use Velero** (backup tool that understands Kubernetes resources)
 2. **Store critical configs in Git**
 3. **Keep backups at application level** (not just cluster level)
-

The Kubeconfig Saga: Your Passport to Clusterville

Kubeconfig is NOT stored in etcd. It's YOUR key to access the cluster. If you lose it, you can't access the cluster even if it's running perfectly!

Anatomy of kubeconfig (`~/.kube/config`):

yaml


```
apiVersion: v1
kind: Config
current-context: prod-cluster  # Which cluster you're currently using

clusters:  # List of clusters you can access
- name: prod-cluster
  cluster:
    server: https://k8s-api.example.com:6443  # API server address
    certificate-authority-data: LS0tLS1CRUdJT... # CA cert (base64)
    # OR
    # certificate-authority: /path/to/ca.crt

users:  # List of user credentials
- name: admin-user
  user:
    client-certificate-data: LS0tLS1CRUdJT... # Your cert (base64)
    client-key-data: LS0tLS1CRUdJT... # Your private key (base64)
    # OR for service account token
    # token: eyJhbGciOiJSUzI1NiIsInR5cCI6I...
```

```
contexts:  # Combines cluster + user + namespace
- name: prod-cluster
  context:
    cluster: prod-cluster
    user: admin-user
    namespace: production  # Default namespace for this context

- name: dev-cluster
  context:
    cluster: dev-cluster
    user: dev-user
    namespace: development
```

Switching contexts:

```
bash

kubectl config get-contexts
kubectl config use-context prod-cluster
kubectl config current-context
```

The Disaster: Lost Kubeconfig

Your laptop died. Kubeconfig gone. Cluster is fine, but you can't access it!

Recovery if you have access to master node:

```
bash
```

```
# SSH to master node
```

```
ssh ubuntu@master-node-1
```

```
# Get admin kubeconfig (kubeadm clusters)
```

```
sudo cat /etc/kubernetes/admin.conf > ~/kubeconfig
```

```
scp ubuntu@master-node-1:~/kubeconfig ~/.kube/config
```

```
# Or generate new one
```

```
kubectl config set-cluster prod-cluster \
```

```
--server=https://k8s-api.example.com:6443 \
```

```
--certificate-authority=/etc/kubernetes/pki/ca.crt \
```

```
--embed-certs=true
```

```
kubectl config set-credentials admin-user \
```

```
--client-certificate=/etc/kubernetes/pki/admin.crt \
```

```
--client-key=/etc/kubernetes/pki/admin.key \
```

```
--embed-certs=true
```

```
kubectl config set-context prod-cluster \
```

```
--cluster=prod-cluster \
```

```
--user=admin-user
```

```
kubectl config use-context prod-cluster
```

Recovery using Service Account Token (if you can't access master):

```
bash
```

```
# If you know a service account name
```

```
TOKEN=$(kubectl --kubeconfig=/etc/kubernetes/admin.conf \
```

```
-n kube-system get secret \
```

```
$(kubectl --kubeconfig=/etc/kubernetes/admin.conf \
```

```
-n kube-system get sa admin-user -o jsonpath='{.secrets[0].name}') \
```

```
-o jsonpath='{.data.token}' | base64 -d)
```

```
kubectl config set-credentials admin-user --token=$TOKEN
```

Best Practice:

1. **Backup kubeconfig** to secure location (password manager, encrypted drive)
2. **Use multiple authentication methods:** Client certs AND OIDC/SSO
3. **Document how to regenerate** credentials
4. **Use RBAC carefully:** Don't give everyone cluster-admin!

The Ingress Restoration Story

Ingresses **ARE** stored in **etcd**, so they'll be restored. But there's a catch...

Ingress depends on:

1. **Ingress Controller** (nginx-ingress, traefik, etc.) - a pod that needs to be running
2. **LoadBalancer Service** - created by the Ingress Controller
3. **DNS records** - pointing to the LoadBalancer IP
4. **TLS Secrets** - certificates for HTTPS

After restoration:

```
bash

kubectl get ingress -A
# Shows all Ingresses (restored from etcd)

kubectl get pods -n ingress-nginx
# Ingress Controller pod exists (restored)

kubectl get svc -n ingress-nginx
# LoadBalancer service exists BUT...
```

The LoadBalancer IP might change!

Before disaster:

```
ingress-nginx-controller LoadBalancer 10.100.200.50 52.1.2.3
```

After restore:

```
ingress-nginx-controller LoadBalancer 10.100.200.50 52.1.2.99 # NEW IP!
```

The cloud provider assigned a **NEW** external IP! Your