

Kubernetes Operator: Design Document

Overview

This system implements a Kubernetes operator that extends the cluster API with custom resources and automated reconciliation logic. The key architectural challenge is building a robust controller that watches cluster state, maintains desired configuration, and handles the complexities of distributed system failures while integrating seamlessly with Kubernetes' declarative model.

This guide is meant to help you understand the big picture before diving into each milestone. Refer back to it whenever you need context on how components connect.

Context and Problem Statement

Milestone(s): All milestones - provides foundational understanding for the entire operator pattern

Managing applications in Kubernetes can quickly become overwhelming as systems grow in complexity. While Kubernetes provides excellent primitives for deploying and running containers, it lacks built-in knowledge about how specific applications should behave, how they should respond to failures, or how they should be configured for optimal performance. This gap between Kubernetes' generic container orchestration and application-specific operational knowledge creates significant challenges for platform teams and application developers.

The fundamental problem is that **operational knowledge** - the accumulated wisdom about how to run, monitor, upgrade, backup, and troubleshoot applications - typically exists only in documentation, runbooks, and the minds of experienced operators. When incidents occur at 3 AM, or when applications need to be scaled or upgraded, human operators must manually translate this knowledge into a series of kubectl commands, configuration changes, and monitoring checks. This approach doesn't scale as organizations grow their Kubernetes footprint or as applications become more sophisticated.

Consider a typical scenario: deploying a stateful application like a database cluster requires not just creating Deployments and Services, but also configuring persistent volumes, setting up replication, managing leader election, coordinating rolling updates without data loss, implementing backup schedules, monitoring cluster health, and handling various failure scenarios. Each of these operations requires deep understanding of both Kubernetes APIs and the application's specific requirements. When done manually, this process is error-prone, time-consuming, and difficult to reproduce consistently across environments.

Mental Model: The Tireless System Administrator

To understand what Kubernetes operators provide, imagine the most experienced system administrator on your team - someone who has managed complex applications for years, knows every configuration nuance, never makes typos, never forgets steps in procedures, and never needs sleep. This administrator sits in front of a terminal 24/7, continuously monitoring the state of all applications under their care.

When this tireless administrator notices that a database replica has failed, they immediately spring into action: checking the health of remaining replicas, determining if the failure affects the primary, deciding whether to promote a new replica or wait for the failed one to recover, updating load balancer configurations if needed, sending alerts to the appropriate teams, and documenting the incident. They perform these actions not by following a rigid script, but by understanding the current state of the system, comparing it to the desired state, and taking the minimal necessary actions to close the gap.

This administrator never gets overwhelmed by multiple simultaneous issues because they process each situation methodically, always working from first principles: "What should the system look like?" and "What's the smallest change I can make to get closer to that goal?" They don't panic during outages, they don't skip steps during routine maintenance, and they learn from each incident to improve their procedures.

A **Kubernetes operator** is essentially this tireless system administrator encoded as software. It continuously watches the actual state of resources in your cluster, compares that state to your declared intentions, and takes corrective actions when gaps are detected. It never gets tired, never forgets to check on things, and can manage dozens or hundreds of application instances simultaneously with perfect consistency.

The operator doesn't just monitor and react - it also embodies years of operational wisdom about the applications it manages. It knows that when scaling up a database cluster, replicas should be added one at a time with health checks between each addition. It understands that certain configuration changes require rolling restarts while others can be applied dynamically. It recognizes the difference between temporary network hiccups and genuine failures that require intervention.

Most importantly, this tireless administrator is **declarative** rather than **imperative**. Instead of being told "scale the database to 5 replicas, then update the configuration, then restart the monitoring agent," the operator is simply told "I want a database cluster with 5 replicas running version 2.4.1 with monitoring enabled." The operator figures out the sequence of operations needed to make that desire a reality, handles any errors or unexpected conditions that arise, and continues working until the actual state matches the declared intent.

Existing Approaches Comparison

Before diving into why operators represent a superior approach, it's important to understand the landscape of existing solutions and their limitations. Each approach addresses part of the operational complexity problem, but none provides the complete solution that operators enable.

Manual kubectl Operations

The most basic approach to managing Kubernetes applications involves manually crafting YAML manifests and applying them with kubectl. This approach gives operators complete control and transparency - every action is explicit and visible.

Aspect	Strengths	Weaknesses
Control	Complete visibility into every operation; no hidden abstractions	Requires deep Kubernetes expertise; prone to human error
Debugging	Easy to trace exactly what resources exist and how they were created	No history of changes; difficult to understand why resources are in current state
Scalability	Works well for simple applications with few resources	Doesn't scale beyond small numbers of applications or environments
Consistency	Can achieve perfect consistency if procedures are followed exactly	Relies entirely on human discipline; no enforcement of procedures
Error Handling	Operator can immediately see and respond to errors	All error detection and recovery is manual; easy to miss problems

The fundamental limitation of manual operations is that they require constant human attention and perfect execution. Consider managing a multi-environment deployment: ensuring that the same sequence of operations is performed correctly in development, staging, and production environments requires meticulous documentation and discipline. When problems arise - a pod fails to start, a service becomes unreachable, or a configuration change has unexpected side effects - human operators must diagnose the issue and determine appropriate remediation steps.

As the number of applications and environments grows, manual operations become increasingly unsustainable. Teams spend more time on routine operational tasks and less time on developing new features or improving system reliability.

Helm Charts and Package Managers

Helm represents a significant improvement over manual operations by introducing templating, dependency management, and release lifecycle management. Helm charts package application resources into reusable units that can be customized for different environments through values files.

Aspect	Strengths	Weaknesses
Reusability	Charts can be shared and reused across teams and organizations	Limited to deployment-time configuration; no runtime adaptation
Templating	Values files enable environment-specific customization	Complex applications require complex templating logic
Versioning	Release management with rollback capabilities	Rollbacks are coarse-grained; can't handle partial failures gracefully
Ecosystem	Large ecosystem of pre-built charts for common applications	Charts often make assumptions that don't fit specific use cases
Day-2 Operations	Good for initial deployment and major updates	No ongoing lifecycle management; no automatic problem resolution

Helm excels at solving the **deployment** problem - packaging applications in a way that makes them easy to install, configure, and upgrade. However, Helm charts are fundamentally **installation tools**, not **operational tools**. Once a Helm release is deployed, it doesn't continue monitoring the application, responding to failures, or adapting to changing conditions.

For example, a Helm chart can deploy a database cluster with appropriate initial configuration, but it can't automatically scale the cluster based on load, replace failed replicas, perform routine maintenance tasks, or coordinate complex upgrade procedures that require specific sequences of operations.

Custom Scripts and Automation

Many organizations bridge the gap between manual operations and full automation by developing custom scripts, often integrated with CI/CD pipelines or scheduled as cron jobs. These scripts encode operational procedures and can perform routine maintenance tasks automatically.

Aspect	Strengths	Weaknesses
Customization	Scripts can encode organization-specific procedures perfectly	High maintenance overhead; scripts break as systems evolve
Integration	Can integrate with existing monitoring and alerting systems	Often fragile and difficult to debug when problems occur
Automation	Reduces manual toil for routine operations	Usually imperative; difficult to ensure idempotency
Knowledge Capture	Procedures are documented in code rather than runbooks	Knowledge is scattered across many scripts with no unified model
Error Recovery	Can include error handling and retry logic	Error handling is ad-hoc; difficult to handle all edge cases

The primary challenge with custom scripts is that they tend to be **imperative** rather than **declarative**. A script might contain logic like "if the database has fewer than 3 replicas, create a new one," but scripts struggle with more complex scenarios like "ensure the database cluster is healthy and properly configured according to this specification."

Scripts also suffer from the **coordination problem**: multiple scripts running simultaneously can interfere with each other, leading to race conditions and inconsistent state. Without a unified control plane, it's difficult to ensure that all automation is working toward the same goals.

Why the Operator Pattern

The operator pattern addresses the fundamental limitations of existing approaches by treating **operational knowledge as code** and integrating deeply with Kubernetes' declarative architecture. Rather than viewing automation as a collection of scripts or deployment tools, operators embody the principle that application management should be as declarative and self-healing as the underlying Kubernetes platform.

Design Insight: The operator pattern succeeds because it aligns with Kubernetes' core philosophy: desired state should be declared, and the system should continuously work to make reality match that declaration. Operators extend this philosophy from basic Kubernetes resources to application-specific concerns.

Encoding Operational Knowledge

The most significant advantage of operators is their ability to **encode operational knowledge directly into the control plane**. Instead of documenting procedures in runbooks that humans must execute, operators embed this knowledge into software that can execute procedures consistently and reliably.

Consider the operational knowledge required to manage a database cluster:

Operational Concern	Traditional Approach	Operator Approach
Scaling	Runbook: "Add replicas one at a time, wait for sync, update load balancer"	Code: Reconciliation loop that creates replicas, monitors readiness, updates services
Failure Recovery	Alert → Human investigates → Manual recovery steps	Automated detection → Self-healing → Human notified of resolution
Upgrades	Maintenance window → Manual coordination → Hope nothing breaks	Rolling upgrade with automated rollback on failure detection
Configuration Changes	Change management → Careful manual application → Monitor for issues	Declarative configuration → Automatic application → Drift detection
Backup Management	Scheduled scripts → Manual verification → Hope backups work	Integrated backup scheduling → Automatic validation → Recovery testing

This knowledge encoding provides several critical benefits:

Consistency: The same operational procedures are executed identically every time, regardless of which team member initiated the action or what time of day it occurs. There's no variation in procedure execution based on operator experience or stress levels.

Auditability: All operational actions are taken through the Kubernetes API and can be logged, monitored, and audited. The decision-making process is transparent and reproducible.

Testability: Operational procedures can be tested in development environments before being applied to production. Complex scenarios like failure recovery can be validated through chaos engineering approaches.

Evolution: Operational knowledge improves over time as the operator code is refined. Lessons learned from incidents can be encoded into the operator, making the entire system more resilient.

Self-Healing Infrastructure

Operators enable **self-healing infrastructure** by continuously monitoring actual state and taking corrective actions when deviations from desired state are detected. This goes far beyond simple resource creation and deletion to include application-specific health checks, performance optimization, and failure recovery procedures.

Decision: Continuous Reconciliation vs Event-Driven Reactions

- **Context:** Systems can respond to changes either by continuously checking state or by reacting to specific events
- **Options Considered:** Event-driven (webhook-based), polling-based, hybrid approach with informers
- **Decision:** Kubernetes informer pattern with continuous reconciliation loops
- **Rationale:** Event-driven systems can miss events or become inconsistent during network partitions; continuous reconciliation provides eventual consistency guarantees and can detect/correct drift from any source
- **Consequences:** Higher resource usage but much stronger consistency guarantees and better failure recovery

The self-healing capabilities of operators manifest in several ways:

Automatic Problem Detection: Operators continuously monitor not just whether resources exist, but whether they're functioning correctly according to application-specific criteria. A database operator might check not only that pods are running, but that replication is working, that performance metrics are within acceptable ranges, and that backup processes are completing successfully.

Intelligent Recovery: When problems are detected, operators can take application-aware corrective actions. Rather than simply restarting failed components, an operator might determine whether a failure affects service availability, attempt graceful recovery procedures first, and coordinate complex recovery processes that involve multiple components.

Drift Prevention: Operators detect and correct configuration drift, ensuring that the actual system configuration always matches the declared specification. This prevents the gradual degradation that often occurs in manually managed systems as small changes accumulate over time.

Adaptive Behavior: Advanced operators can adapt their behavior based on observed conditions. For example, an operator might automatically adjust resource allocations based on load patterns, or modify backup frequencies based on data change rates.

Integration with Kubernetes Ecosystem

Operators integrate seamlessly with the broader Kubernetes ecosystem, leveraging existing tooling, monitoring, and operational practices rather than requiring entirely new approaches.

Integration Point	Benefit	Example
kubectl and API	Operators are managed using standard Kubernetes tools	<code>kubectl get databases</code> shows all database instances with custom columns
RBAC	Standard Kubernetes security model applies	Database operator can be restricted to specific namespaces or resource types
Monitoring	Prometheus metrics, events, and logging work naturally	Operator exposes application-specific metrics alongside infrastructure metrics
GitOps	Custom resources can be managed through Git workflows	Database configurations versioned and deployed through standard GitOps practices
Networking	Application networking integrates with Kubernetes services and ingress	Operator manages service discovery and load balancing using native mechanisms

This integration means that adopting operators doesn't require teams to learn entirely new toolsets or abandon existing operational practices. Instead, operators enhance and extend existing Kubernetes workflows with application-specific intelligence.

The operator pattern represents a fundamental shift from **reactive operational practices** to **proactive declarative management**. Rather than waiting for problems to occur and then responding, operators continuously work to maintain desired state and prevent problems from arising in the first place. This shift enables organizations to manage complex applications at scale while reducing operational toil and improving system reliability.

Key Insight: Operators don't just automate existing manual procedures - they enable entirely new approaches to system management that aren't practical with manual operations. The continuous reconciliation model allows for levels of consistency and reliability that are impossible to achieve through human-driven processes.

Implementation Guidance

Understanding the operator pattern conceptually is the first step; implementing operators successfully requires careful technology choices and project structure decisions. This section provides concrete guidance for teams building their first Kubernetes operator.

Technology Recommendations

Component	Simple Option	Advanced Option
Operator Framework	Operator SDK with Helm or Ansible	Kubebuilder with controller-runtime
Language Runtime	Go 1.21+ with controller-runtime	Go with custom client-go usage
Custom Resource Generation	Kubebuilder markers and code generation	Hand-written OpenAPI v3 schemas
Testing Framework	Envtest with Ginkgo/Gomega	Custom test harness with kind clusters
Admission Webhooks	controller-runtime webhook builder	Custom HTTP servers with admission review handling
Certificate Management	cert-manager with self-signed issuer	External certificate authority integration

For teams new to operator development, **Kubebuilder with controller-runtime** provides the optimal balance of power and ease of use. It generates substantial boilerplate code, provides excellent testing utilities, and follows Kubernetes community best practices.

Recommended Project Structure

Organizing operator code properly from the beginning prevents architectural debt and makes the codebase maintainable as features are added:

```
database-operator/
├── cmd/
│   └── main.go           ← Entry point and flag parsing
├── api/
│   └── v1/
│       ├── database_types.go    ← Custom resource type definitions
│       ├── groupversion_info.go  ← API group and version metadata
│       └── zz_generated.deepcopy.go ← Generated code (do not edit)
├── controllers/
│   ├── database_controller.go    ← Main reconciliation logic
│   └── database_controller_test.go ← Controller unit tests
├── internal/
│   ├── resources/              ← Resource creation and management
│   │   ├── deployment.go        ← Database deployment logic
│   │   ├── service.go          ← Service creation and updates
│   │   └── configmap.go        ← Configuration management
│   ├── status/                 ← Status calculation and updates
│   │   └── conditions.go      ← Status condition management
│   └── webhooks/               ← Admission webhook implementations
        ├── validator.go         ← Validation webhook logic
        └── defaulter.go         ← Defaulting webhook logic
└── config/
    ├── crd/                  ← Custom resource definitions
    ├── rbac/                 ← Role and binding manifests
    ├── webhook/               ← Webhook configuration
    └── samples/               ← Example custom resources
└── hack/
    └── boilerplate.go.txt     ← License header for generated files
```

This structure separates concerns clearly: API definitions remain stable in the `api/` directory, business logic lives in `controllers/`, and implementation details are encapsulated in `internal/` packages.

Infrastructure Starter Code

The following code provides a complete foundation for operator development. Copy this code and modify it according to your specific requirements rather than building everything from scratch.

Main Entry Point (cmd/main.go):

```
package main

import (
    "context"
    "log"
    "os"

    "github.com/argoproj/operator/v2/pkg/controller"
    "github.com/argoproj/operator/v2/pkg/resource"
    "github.com/argoproj/operator/v2/pkg/version"
    "github.com/argoproj/operator/v2/pkg/webhook"
)

var (
    controllerName = "operator"
    controllerPath = "operator"
    controllerVersion = version.V2
    controllerResource = resource.Resource{Group: "operator", Version: "v2", Kind: "Operator"}
    controllerWebhook = webhook.Webhook{Path: "/operator", Version: "v2", Controller: controllerName}
)

func main() {
    log.Println("Starting operator")
    controller.Run(controllerName, controllerPath, controllerVersion, controllerResource, controllerWebhook)
}
```

```
package main
```

GO

```
import (
    "flag"
    "fmt"
    "os"

    "k8s.io/apimachinery/pkg/runtime"
    utilruntime "k8s.io/apimachinery/pkg/util/runtime"
    clientgoscheme "k8s.io/client-go/kubernetes/scheme"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/healthz"
    "sigs.k8s.io/controller-runtime/pkg/log/zap"

    databasev1 "github.com/yourorg/database-operator/api/v1"
    "github.com/yourorg/database-operator/controllers"
)

var (
    scheme     = runtime.NewScheme()
    setupLog = ctrl.Log.WithName("setup")
)
```

func init() {

```
    utilruntime.Must(clientgoscheme.AddToScheme(scheme))
    utilruntime.Must(databasev1.AddToScheme(scheme))
}
```

func main() {

```
    var metricsAddr string
    var enableLeaderElection bool
    var probeAddr string

    flag.StringVar(&metricsAddr, "metrics-bind-address", ":8080", "The address the metric endpoint binds to.")
    flag.StringVar(&probeAddr, "health-probe-bind-address", ":8081", "The address the probe endpoint binds to.")
    flag.BoolVar(&enableLeaderElection, "leader-elect", false, "Enable leader election for controller manager.")
```

```
opts := zap.Options{Development: true}

opts.BindFlags(flag.CommandLine)

flag.Parse()

ctrl.SetLogger(zap.New(zap.UseFlagOptions(&opts)))

mgr, err := ctrl.NewManager(ctrl.GetConfigOrDie(), ctrl.Options{
    Scheme:           scheme,
    MetricsBindAddress: metricsAddr,
    Port:             9443,
    HealthProbeBindAddress: probeAddr,
    LeaderElection:   enableLeaderElection,
    LeaderElectionID: "database-operator-lock",
})

if err != nil {
    setupLog.Error(err, "unable to start manager")
    os.Exit(1)
}

if err = (&controllers.DatabaseReconciler{
    Client: mgr.GetClient(),
    Scheme: mgr.GetScheme(),
}).SetupWithManager(mgr); err != nil {
    setupLog.Error(err, "unable to create controller", "controller", "Database")
    os.Exit(1)
}

if err := mgr.AddHealthzCheck("healthz", healthz.Ping); err != nil {
    setupLog.Error(err, "unable to set up health check")
    os.Exit(1)
}

if err := mgr.AddReadyzCheck("readyz", healthz.Ping); err != nil {
```

```
    setupLog.Error(err, "unable to set up ready check")

    os.Exit(1)

}

setupLog.Info("starting manager")

if err := mgr.Start(ctrl.SetupSignalHandler()); err != nil {

    setupLog.Error(err, "problem running manager")

    os.Exit(1)

}
}
```

Basic Custom Resource Type (api/v1/database_types.go):

```
package v1

import (
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
)

// DatabaseSpec defines the desired state of Database

type DatabaseSpec struct {

    // TODO: Define your custom resource specification fields here

    // Examples:

    // Replicas      *int32  `json:"replicas,omitempty"`
    // Version       string   `json:"version"`
    // StorageSize   string   `json:"storageSize,omitempty"`

}

// DatabaseStatus defines the observed state of Database

type DatabaseStatus struct {

    // TODO: Define your status fields here

    // Examples:

    // Conditions     []metav1.Condition `json:"conditions,omitempty"`
    // ReadyReplicas   int32            `json:"readyReplicas,omitempty"`
    // ObservedGeneration int64          `json:"observedGeneration,omitempty"`

}

//+kubebuilder:object:root=true

//+kubebuilder:subresource:status

// Database is the Schema for the databases API

type Database struct {

    metav1.TypeMeta   `json:",inline"`
    metav1.ObjectMeta `json:"metadata,omitempty"`

    Spec  DatabaseSpec `json:"spec,omitempty"`
    Status DatabaseStatus `json:"status,omitempty"`

}

//+kubebuilder:object:root=true
```

```
// DatabaseList contains a list of Database

type DatabaseList struct {

    metav1.TypeMeta `json:",inline"`

    metav1.ListMeta `json:"metadata,omitempty"`

    Items          []Database `json:"items"`

}

func init() {
    SchemeBuilder.Register(&Database{}, &DatabaseList{})
}
```

Core Logic Skeleton

The following skeleton provides the structure for implementing your reconciliation logic. Fill in the TODO sections based on your specific application requirements:

```
package controllers

import (
    "context"
    "time"

    "k8s.io/apimachinery/pkg/runtime"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/log"

    databasev1 "github.com/yourorg/database-operator/api/v1"
)

// DatabaseReconciler reconciles a Database object

type DatabaseReconciler struct {
    client.Client
    Scheme *runtime.Scheme
}

//+kubebuilder:rbac:groups=database.example.com,resources=databases,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups=database.example.com,resources=databases/status,verbs=get;update;patch
//+kubebuilder:rbac:groups=database.example.com,resources=databases/finalizers,verbs=update

// Reconcile implements the main reconciliation logic for Database resources

func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
    log := log.FromContext(ctx)

    // TODO 1: Fetch the Database instance from the API server
    // Use r.Get() to retrieve the resource by namespaced name
    // Handle not-found errors gracefully (resource may have been deleted)

    // TODO 2: Check if the resource is being deleted (DeletionTimestamp is set)
    // If deleting and has finalizers, perform cleanup operations
    // Remove finalizers when cleanup is complete to allow deletion
}
```

```

    // TODO 3: Add finalizer if not already present

    // Finalizers ensure cleanup logic runs before resource deletion

    // Use controllerutil.AddFinalizer() helper function

    // TODO 4: Reconcile owned resources (Deployments, Services, ConfigMaps, etc.)

    // Compare desired state (from spec) with actual state (from cluster)

    // Create, update, or delete resources as needed to match desired state

    // TODO 5: Update the Database status based on owned resource states

    // Calculate readiness, error conditions, and other status information

    // Use r.Status().Update() to write status back to API server

    // TODO 6: Determine requeue behavior

    // Return ctrl.Result{} for no requeue

    // Return ctrl.Result{RequeueAfter: 30*time.Second} for periodic reconciliation

    // Return ctrl.Result{Requeue: true} for immediate requeue on recoverable errors

    log.Info("reconciliation completed successfully")

    return ctrl.Result{}, nil
}

// SetupWithManager sets up the controller with the Manager

func (r *DatabaseReconciler) SetupWithManager(mgr ctrl.Manager) error {
    return ctrl.NewControllerManagedBy(mgr).
        For(&databasev1.Database{}).

        // TODO: Add Owns() calls for resources created by this controller

        // Example: Owns(&appsv1.Deployment{}).

        Complete(r)
}

```

Language-Specific Hints

Go-Specific Best Practices:

- Use `sigs.k8s.io/controller-runtime/pkg/client` instead of raw client-go for better ergonomics
- Always check `errors.NotFound()` when fetching resources that might not exist
- Use `controllerutil.SetControllerReference()` to establish owner relationships

- Leverage `kubebuilder` markers for RBAC and webhook generation rather than hand-writing manifests
- Use structured logging with `logr.Logger` instead of standard library logging

Error Handling Patterns:

- Distinguish between permanent errors (bad configuration) and transient errors (network timeouts)
- Use `ctrl.Result{Requeue: true}` for transient errors that should retry immediately
- Use `ctrl.Result{RequeueAfter: duration}` for rate-limited retries or periodic reconciliation
- Log errors with sufficient context for debugging, but don't log expected conditions like "resource not found"

Performance Considerations:

- Cache frequently accessed data in the controller struct rather than fetching repeatedly
- Use `client.Reader` for read-only operations to leverage informer caches
- Batch related operations when possible to reduce API server load
- Consider using predicates to filter events that don't require reconciliation

Milestone Checkpoints

After implementing each major component, verify functionality with these checkpoints:

Custom Resource Definition Checkpoint:

```
# Apply your CRD to a cluster
bash
kubectl apply -f config/crd/bases/
# Verify the CRD is registered
kubectl get crd databases.database.example.com
# Test custom resource creation
kubectl apply -f config/samples/database_v1_database.yaml
# Verify custom columns are displayed
kubectl get databases
```

BASH

Controller Checkpoint:

```
# Run the controller locally
bash
make install run
# In another terminal, create a test resource
kubectl apply -f config/samples/database_v1_database.yaml
# Verify reconciliation occurs (check controller logs)
# Should see "reconciliation completed successfully" messages
# Verify status is updated
kubectl get database sample-database -o yaml
```

BASH

Webhook Checkpoint:

```
# Install cert-manager for webhook certificates
kubectl apply -f https://github.com/cert-manager/cert-manager/releases/download/v1.13.0/cert-manager.yaml
# Deploy webhooks
make deploy-webhooks

# Test validation webhook with invalid resource
# Should receive admission denied error

# Test defaulting webhook by omitting optional fields
# Resource should be created with default values populated
```

BASH

Expected signs that implementation is working correctly:

- Controller logs show reconciliation events for resource changes
- Custom resources display appropriate status conditions and custom columns
- Owned resources (Deployments, Services) are created and updated appropriately
- Webhook validation prevents invalid resources from being created
- Resource deletion triggers proper cleanup and finalizer removal

Goals and Non-Goals

Milestone(s): All milestones - defines the scope and boundaries that guide implementation decisions across CRD design, controller logic, webhooks, testing, and deployment

Building a Kubernetes operator requires clear boundaries around what problems it will solve and what complexities it will deliberately avoid. Think of this goals definition as **drawing the blueprint before breaking ground** - it prevents scope creep, guides architectural decisions, and sets realistic expectations for both implementers and users. Without clear goals, operators often evolve into sprawling systems that try to solve every operational problem, leading to complexity that undermines the very automation they were meant to provide.

The operator pattern excels at codifying repetitive operational procedures, but not every operational task belongs in an operator. Some problems are better solved with existing Kubernetes primitives, external tools, or manual intervention. By explicitly defining what this operator will and will not handle, we create focused software that does its intended job well rather than attempting to be a universal solution.

Functional Goals

The **functional goals** define the specific automation capabilities our operator will provide. These represent the core value proposition - the repetitive operational tasks that the operator will handle automatically, reducing manual intervention and human error.

Primary Resource Management Capabilities

Our operator will provide automated lifecycle management for database instances through custom resources. The primary functional goals center around the `Database` custom resource and its associated operational workflows.

Capability	Description	Automation Level	User Interaction
Database Provisioning	Create new database instances with specified configuration	Fully Automated	User creates Database resource
Replica Management	Scale database replicas up/down based on spec changes	Fully Automated	User updates replicas field
Version Upgrades	Perform rolling upgrades to new database versions	Semi-Automated	User updates version field, operator handles rollout
Storage Management	Provision and resize persistent volumes for database storage	Fully Automated	User specifies storageSize, operator manages PVCs
Configuration Updates	Apply configuration changes without service interruption	Semi-Automated	User updates config, operator validates and applies
Health Monitoring	Continuously monitor database health and update status	Fully Automated	Status conditions updated automatically

Decision: Database-Focused Domain Model

- **Context:** Operators work best with a specific domain rather than generic resource management. We need concrete use cases to drive meaningful automation.
- **Options Considered:** Generic application operator, database-specific operator, multi-application operator
- **Decision:** Database-specific operator with concrete database management workflows
- **Rationale:** Database management has well-understood operational patterns (provisioning, scaling, backup, upgrade) that translate naturally to declarative automation. This provides concrete requirements for testing reconciliation logic without getting lost in abstract scenarios.
- **Consequences:** Focused scope enables deeper automation but limits reusability across application types

The operator will manage **owned resources** - Kubernetes resources created and controlled by the operator to implement the desired database configuration. These owned resources form the actual infrastructure that delivers the database service.

Owned Resource Type	Purpose	Lifecycle Management	Status Reflection
StatefulSet	Manages database pods with stable identity and storage	Created, updated, and scaled based on Database spec	Ready replicas count reflected in Database status
Service	Provides stable network endpoint for database access	Created with appropriate ports and selectors	Service endpoints reflected in Database status
ConfigMap	Holds database configuration files and initialization scripts	Updated when Database spec.config changes	Configuration hash tracked in status
PersistentVolumeClaim	Provides durable storage for database data	Created and resized based on spec.storageSize	Storage capacity reflected in status
Secret	Manages database credentials and TLS certificates	Created with generated passwords, rotated on demand	Credential generation timestamp in status

Declarative Configuration Management

The operator will support **declarative configuration** where users specify the desired end state through the `DatabaseSpec`, and the operator continuously works to achieve and maintain that state through reconciliation.

Configuration Aspect	Declarative Behavior	Validation Rules	Default Handling
Replica Count	Scale StatefulSet to match spec.replicas	Must be positive integer, maximum of 10	Defaults to 1 if omitted
Database Version	Perform rolling upgrade to spec.version	Must be supported version from allowed list	Defaults to latest stable version
Storage Size	Resize PVC to match spec.storageSize	Cannot decrease, must use valid resource quantities	Defaults to "10Gi"
Resource Requirements	Update container resources in StatefulSet	Must fit within namespace resource quotas	Defaults to moderate CPU/memory requests
Configuration Parameters	Update ConfigMap with spec.config values	Validate against database-specific schema	Applies sensible production defaults

The Desired State Principle

The operator embodies Kubernetes' core principle that users declare what they want, not how to achieve it. When a user creates a `Database` resource requesting 3 replicas with 50Gi storage, they shouldn't need to know about StatefulSets, PVCs, or rolling update strategies. The operator translates high-level intent into low-level Kubernetes resource management.

Status Reporting and Observability

The operator will provide comprehensive status reporting through the `DatabaseStatus` subresource, giving users visibility into the current state of their database instances without requiring knowledge of the underlying Kubernetes resources.

Status Information	Update Frequency	Status Conditions	Troubleshooting Value
Ready Replica Count	Every reconciliation cycle	Ready, Progressing, Degraded	Indicates scaling progress and health
Current Version	On version change completion	VersionUpgrading, VersionReady	Shows upgrade progress and rollback status
Storage Utilization	Periodic polling (configurable)	StorageWarning, StorageCritical	Enables proactive capacity planning
Connection Information	On service endpoint changes	ServiceReady, ServiceUnavailable	Provides connection details for applications
Last Reconciliation	Every reconciliation attempt	ReconcileSuccess, ReconcileError	Shows operator activity and error states

Non-Functional Goals

The **non-functional goals** define quality attributes that constrain how the operator implements its functional capabilities. These goals ensure the operator behaves reliably in production environments and integrates smoothly with Kubernetes operational practices.

Reliability and Fault Tolerance

The operator must handle the distributed system challenges inherent in Kubernetes environments, where network partitions, node failures, and API server unavailability are routine operational realities.

Reliability Aspect	Target Behavior	Failure Handling	Recovery Mechanism
Controller Availability	99.9% uptime during cluster operation	Leader election prevents split-brain scenarios	Automatic failover to standby controller replica
Reconciliation Consistency	All reconciliation operations are idempotent	Partial failures don't leave resources in inconsistent states	Full state reconstruction from cluster resources
API Server Connectivity	Graceful handling of connection interruptions	Exponential backoff on API failures	Automatic reconnection with informer resync
Resource Conflict Resolution	Handle concurrent updates from multiple sources	Optimistic locking prevents conflicting changes	Retry with latest resource version on conflicts
Data Durability	Never lose user-specified configuration	All spec changes persisted before processing	Status updates use separate API calls to prevent overwrites

Decision: Eventual Consistency Model

- **Context:** Kubernetes operates on eventual consistency - cluster state changes propagate asynchronously through informers and caches
- **Options Considered:** Strong consistency with distributed locking, eventual consistency with convergence, manual consistency checking
- **Decision:** Embrace eventual consistency with robust reconciliation convergence
- **Rationale:** Fighting Kubernetes' eventual consistency model leads to complex synchronization code that's prone to deadlocks. Designing for convergence through idempotent operations aligns with Kubernetes patterns and provides better fault tolerance.
- **Consequences:** Simpler controller logic and better failure recovery, but requires careful status reporting to avoid user confusion during convergence periods

Performance and Scalability

The operator should handle reasonable numbers of database instances without consuming excessive cluster resources or causing performance degradation.

Performance Metric	Target Value	Measurement Method	Scaling Behavior
Reconciliation Latency	< 30 seconds for spec changes	Time from Database update to owned resource convergence	Linear scaling with number of owned resources per Database
Resource Consumption	< 100MB memory, < 0.1 CPU cores per 100 Database resources	Metrics exported via Prometheus endpoints	Memory usage grows with informer cache size
API Server Load	< 10 QPS per controller replica	Rate limiting on client requests	Batched status updates reduce API call frequency
Webhook Response Time	< 5 seconds for validation/mutation	Admission review processing time	Independent of cluster size for validation rules
Controller Startup Time	< 60 seconds to ready state	Time from pod start to processing first reconciliation	Informer cache sync time dominates startup

Operational Integration

The operator must integrate seamlessly with standard Kubernetes operational practices and tooling used by platform teams.

Integration Aspect	Standard Compliance	Tooling Compatibility	Operational Benefit
RBAC Authorization	Principle of least privilege for all permissions	Compatible with namespace-scoped and cluster-scoped deployments	Security audit compliance and permission isolation
Resource Labeling	Consistent labels for owned resources with operator identification	Works with label selectors in monitoring and debugging tools	Clear ownership tracking and troubleshooting
Event Recording	Kubernetes events for all significant operator actions	Visible in kubectl describe and event aggregation systems	Audit trail and debugging information
Metrics Export	Prometheus metrics for controller performance and resource counts	Integrates with existing monitoring infrastructure	Operational visibility and alerting capability
Logging Structure	Structured JSON logs with consistent fields	Compatible with log aggregation systems like ELK or Loki	Centralized logging and log-based alerting

Operational Transparency Principle

Platform teams need visibility into operator behavior to debug issues, plan capacity, and maintain security compliance. The operator should export its internal state and decision-making process through standard Kubernetes observability mechanisms rather than requiring specialized knowledge or tools.

Explicit Non-Goals

The **explicit non-goals** define complexities and capabilities that this operator implementation will deliberately exclude. These boundaries prevent scope creep and ensure the operator remains focused on its core value proposition while avoiding areas where other tools provide better solutions.

Data Management and Backup Operations

While the operator manages database infrastructure, it will not implement data-level backup, restore, or migration capabilities.

Excluded Capability	Rationale	Recommended Alternative	Interface Boundary
Automated Backup Scheduling	Database backup strategies vary significantly by database type and organizational requirements	Use dedicated backup operators like Velero or database-specific backup tools	Operator provisions storage; external tools handle data backup
Point-in-Time Recovery	Recovery procedures require database-specific knowledge and operational context	Database-native tools with organization-specific recovery procedures	Operator provides persistent storage; DBA tools handle recovery
Cross-Cluster Data Migration	Migration involves network, security, and data governance concerns beyond infrastructure	Specialized migration tools with proper data governance workflow	Operator manages instances; migration tools handle data movement
Data Encryption at Rest	Encryption requirements vary by compliance framework and key management infrastructure	Kubernetes storage encryption or database-native encryption features	Operator specifies storage requirements; platform provides encryption

Decision: Infrastructure vs Data Layer Separation

- **Context:** Operators can manage infrastructure (pods, services, storage) or data-level operations (backup, migration, schema changes)
- **Options Considered:** Full data lifecycle management, infrastructure-only management, hybrid approach with pluggable data operations
- **Decision:** Infrastructure-only management with clear boundaries at the data layer
- **Rationale:** Data operations require database-specific expertise, compliance considerations, and organizational procedures that vary significantly across environments. Infrastructure management has consistent patterns across database types.
- **Consequences:** Simpler operator implementation and clearer separation of concerns, but users need additional tools for data operations

Advanced Deployment Strategies

The operator will not implement complex deployment patterns that require sophisticated orchestration beyond basic rolling updates.

Excluded Pattern	Complexity Reason	Standard Alternative	Operator Boundary
Blue-Green Deployments	Requires maintaining parallel full environments and traffic switching logic	Use service mesh or ingress controllers for traffic management	Operator manages single environment; traffic tools handle switching
Canary Deployments	Requires traffic percentage routing and automated rollback based on metrics	Use progressive delivery tools like Flagger or Argo Rollouts	Operator performs rolling updates; specialized tools handle canary logic
Multi-Region Replication	Cross-region networking, latency, and consistency concerns	Use database-native replication or service mesh multi-cluster features	Operator manages single-cluster instances; platform handles multi-region
Disaster Recovery Automation	Requires integration with backup systems, DNS switching, and recovery procedures	Use comprehensive DR orchestration tools with organizational runbooks	Operator provisions instances; DR tools coordinate full recovery

Security and Access Control

While the operator follows security best practices, it will not implement application-level security features that belong in specialized security tools.

Security Aspect	Exclusion Rationale	Recommended Approach	Responsibility Boundary
User Authentication	Database user management involves organizational identity systems	Use external authentication providers or database-native user management	Operator creates database infrastructure; identity systems handle users
Network Policies	Network security requirements vary by organizational policy and compliance framework	Use Kubernetes NetworkPolicy resources or service mesh security features	Operator labels resources for policy targeting; network tools enforce policies
Secret Rotation	Credential rotation involves integration with organizational secret management systems	Use external secret operators like External Secrets Operator or Vault integration	Operator consumes secrets; secret management tools handle rotation
Compliance Auditing	Audit requirements vary by regulatory framework and organizational procedures	Use specialized compliance tools like OPA Gatekeeper or audit logging systems	Operator generates audit events; compliance tools enforce policies

Pitfall: Scope Creep Through Security Requirements

Security requirements often drive scope expansion as stakeholders request "just one more security feature." Establishing clear boundaries early prevents the operator from becoming a monolithic security solution. Instead, design clean interfaces for integration with specialized security tools.

Performance Tuning and Optimization

The operator will not implement database-specific performance tuning or optimization features that require deep database expertise.

Performance Aspect	Exclusion Reason	Alternative Solution	Integration Point
Query Performance Tuning	Requires database-specific knowledge and workload analysis	Use database monitoring tools and DBA expertise	Operator provisions resources; DBAs tune performance
Automatic Index Management	Index strategies depend on application query patterns and data characteristics	Use database-native optimization features or specialized tuning tools	Operator manages database infrastructure; applications manage schema
Workload-Based Resource Scaling	Requires metrics analysis and workload prediction beyond basic replica scaling	Use Horizontal Pod Autoscaler or Vertical Pod Autoscaler with custom metrics	Operator handles manual scaling; autoscaling tools handle automatic scaling
Connection Pool Management	Connection pooling strategies vary by application architecture and traffic patterns	Use application-level connection pooling or proxy tools like PgBouncer	Operator provides database endpoints; applications handle connection management

Implementation Scope Summary

These goals and non-goals create a **focused operator implementation** that provides real value while maintaining reasonable complexity. The operator will excel at infrastructure lifecycle management - the operational tasks that follow predictable patterns and can be reliably automated through Kubernetes primitives.

In Scope	Out of Scope	Interface/Integration Point
Database infrastructure provisioning and scaling	Database schema management and migrations	Operator provides stable endpoints for schema tools
Resource lifecycle management (pods, services, storage)	Application deployment and configuration	Applications connect to operator-managed database services
Basic health monitoring and status reporting	Advanced performance monitoring and alerting	Operator exports metrics for monitoring systems
Rolling updates for version changes	Complex deployment strategies and traffic management	Operator performs updates; traffic tools handle advanced routing
Infrastructure-level configuration management	Database-specific tuning and optimization	Operator provides base configuration; DBAs handle tuning

This scope definition guides every implementation decision from CRD schema design through controller logic to testing strategies. When faced with feature requests or implementation choices, refer back to these goals to maintain focus and deliver a robust, maintainable operator that solves its intended problems well.

Implementation Guidance

The goals and non-goals defined above directly influence the technical implementation approach and technology choices throughout the operator development process. This guidance translates the high-level goals into concrete implementation decisions and project structure recommendations.

Technology Recommendations Table

Component	Simple Option	Advanced Option	Recommendation for Learning
Operator Framework	Kubebuilder with basic scaffolding	Operator SDK with advanced features and scorecard	Kubebuilder - cleaner generated code, better documentation
Custom Resource Schema	Basic OpenAPI v3 schema validation	CEL expressions with cross-field validation	Start with OpenAPI, add CEL for complex validation
Controller Runtime	Controller-runtime with basic manager	Custom controller with shared informers	Controller-runtime - handles boilerplate and follows patterns
Webhook Framework	Controller-runtime webhook builder	Custom webhook server with admission review handling	Controller-runtime webhook builder - integrated TLS and registration
Testing Framework	Fake client for unit tests only	Envtest for integration testing with real API server	Both - fake client for fast unit tests, envtest for end-to-end validation
Deployment Method	Static YAML manifests	Helm charts with configurable values	Start with static YAML, migrate to Helm for production deployment

Recommended Project Structure

The project structure should reflect the goals' emphasis on focused domain logic while maintaining clear separation between infrastructure concerns and business logic:

```

database-operator/
├── cmd/
│   └── manager/
│       └── main.go           ← Entry point with manager setup
├── apis/
│   └── database/
│       └── v1/
│           ├── database_types.go    ← Database CRD definition with spec/status
│           ├── database_webhook.go  ← Validation and mutation logic
│           └── zz_generated.deepcopy.go ← Generated by controller-gen
└── controllers/
    └── database/
        ├── database_controller.go    ← Core reconciliation logic
        ├── database_controller_test.go ← Unit tests with fake client
        └── suite_test.go             ← Envtest integration test setup
internal/
├── resources/                         ← Owned resource management (non-goals boundary)
│   ├── statefulset.go                 ← StatefulSet creation and updates
│   ├── service.go                    ← Service endpoint management
│   ├── configmap.go                 ← Configuration management
│   └── secret.go                    ← Credential handling
└── database/                          ← Database-specific logic (extensibility point)
    ├── postgres.go                  ← PostgreSQL-specific configuration
    └── mysql.go                     ← MySQL-specific configuration (future)
config/
├── crd/                                ← CRD manifests and kustomization
├── rbac/                               ← RBAC role definitions following least privilege
├── webhook/                            ← Webhook configuration and TLS setup
└── manager/                            ← Controller deployment manifests
hack/
├── boilerplate.go.txt                ← License header for generated files
└── tools.go                           ← Tool dependencies for go mod

```

This structure supports the functional goals by separating domain-specific logic (`internal/database/`) from generic Kubernetes resource management (`internal/resources/`), making it easy to extend support for additional database types without violating the infrastructure vs data management boundary.

Infrastructure Starter Code

Manager Setup (cmd/manager/main.go) - Complete working foundation:

```
package main
```

GO

```
import (
    "flag"
    "os"
    "time"

    "k8s.io/apimachinery/pkg/runtime"
    clientgoscheme "k8s.io/client-go/kubernetes/scheme"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/healthz"
    "sigs.k8s.io/controller-runtime/pkg/log/zap"

    databasev1 "github.com/example/database-operator/apis/database/v1"
    databasecontrollers "github.com/example/database-operator/controllers/database"
)

var (
    scheme = runtime.NewScheme()
    setupLog = ctrl.Log.WithName("setup")
)
```

```
func init() {
    _ = clientgoscheme.AddToScheme(scheme)
    _ = databasev1.AddToScheme(scheme)
}
```

```
func main() {
    var metricsAddr string
    var probeAddr string
    var enableLeaderElection bool
    var syncPeriod time.Duration

    flag.StringVar(&metricsAddr, "metrics-bind-address", ":8080", "The address the metric endpoint binds to.")
    flag.StringVar(&probeAddr, "health-probe-bind-address", ":8081", "The address the probe endpoint binds to.")
    flag.BoolVar(&enableLeaderElection, "leader-elect", false, "Enable leader election for controller manager.")

    flag.DurationVar(&syncPeriod, "sync-period", 10*time.Minute, "Minimum frequency at which watched resources are reconciled.")
}
```

```

opts := zap.Options{Development: true}

opts.BindFlags(flag.CommandLine)

flag.Parse()

ctrl.SetLogger(zap.New(zap.UseFlagOptions(&opts)))

mgr, err := ctrl.NewManager(ctrl.GetConfigOrDie(), ctrl.Options{
    Scheme:           scheme,
    MetricsBindAddress: metricsAddr,
    Port:             9443, // Webhook server port
    HealthProbeBindAddress: probeAddr,
    LeaderElection:   enableLeaderElection,
    LeaderElectionID: "database-operator-leader-election",
    SyncPeriod:       &syncPeriod,
})

```

```

if err != nil {
    setupLog.Error(err, "unable to start manager")
    os.Exit(1)
}

```

```

// Register controller - learner implements the reconciliation logic
if err = (&databasecontrollers.DatabaseReconciler{
    Client: mgr.GetClient(),
    Scheme: mgr.GetScheme(),
    Log:     ctrl.Log.WithName("controllers").WithName("Database"),
}).SetupWithManager(mgr); err != nil {
    setupLog.Error(err, "unable to create controller", "controller", "Database")
    os.Exit(1)
}

```

```

// Register webhooks - learner implements validation and mutation
if err = (&databasev1.Database{}).SetupWebhookWithManager(mgr); err != nil {
    setupLog.Error(err, "unable to create webhook", "webhook", "Database")
    os.Exit(1)
}

```

```
// Health and readiness endpoints

if err := mgr.AddHealthzCheck("healthz", healthz.Ping); err != nil {
    setupLog.Error(err, "unable to set up health check")
    os.Exit(1)
}

if err := mgr.AddReadyzCheck("readyz", healthz.Ping); err != nil {
    setupLog.Error(err, "unable to set up ready check")
    os.Exit(1)
}

setupLog.Info("starting manager")

if err := mgr.Start(ctrl.SetupSignalHandler()); err != nil {
    setupLog.Error(err, "problem running manager")
    os.Exit(1)
}

}
```

Resource Helper Utilities (internal/resources/common.go) - Infrastructure code supporting non-functional goals:

```
package resources

import (
    "fmt"

    corev1 "k8s.io/api/core/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/runtime"
    ctrl "sigs.k8s.io/controller-runtime"

    databasev1 "github.com/example/database-operator/apis/database/v1"
)

const (
    // Labels for owned resource identification (supports operational integration goals)
    ManagedByLabelKey = "app.kubernetes.io/managed-by"
    ManagedByLabelValue = "database-operator"
    NameLabelKey        = "app.kubernetes.io/name"
    ComponentLabelKey   = "app.kubernetes.io/component"
    InstanceLabelKey    = "app.kubernetes.io/instance"
)

// CommonLabels returns the standard labels for all owned resources
// Supports operational integration goal for resource identification

func CommonLabels(database *databasev1.Database, component string) map[string]string {
    return map[string]string{
        ManagedByLabelKey: ManagedByLabelValue,
        NameLabelKey:       "database",
        ComponentLabelKey: component,
        InstanceLabelKey:   database.Name,
    }
}

// SetOwnerReference establishes owner reference for garbage collection
// Supports reliability goal for automatic cleanup

func SetOwnerReference(owner *databasev1.Database, controlled metav1.Object, scheme *runtime.Scheme) error {
```

GO

```
    return ctrl.SetControllerReference(owner, controlled, scheme)
}

// ResourceName generates consistent naming for owned resources

func ResourceName(database *databasev1.Database, suffix string) string {
    if suffix == "" {
        return database.Name
    }
    return fmt.Sprintf("%s-%s", database.Name, suffix)
}

// CommonAnnotations returns standard annotations for owned resources

func CommonAnnotations(database *databasev1.Database) map[string]string {
    return map[string]string{
        "database.example.com/managed-by-operator": "true",
        "database.example.com/owner-generation":     fmt.Sprintf("%d", database.Generation),
    }
}
```

Core Logic Skeleton Code

Database Reconciler (controllers/database/database_controller.go) - Core logic that learners implement:

```
package database

import (
    "context"
    "time"

    appsv1 "k8s.io/api/apps/v1"
    corev1 "k8s.io/api/core/v1"
    "k8s.io/apimachinery/pkg/api/errors"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/apimachinery/pkg/types"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/log"
    databasev1 "github.com/example/database-operator/apis/database/v1"
    "github.com/example/database-operator/internal/resources"
)

const (
    REQUEUE_AFTER_DURATION = 30 * time.Second
    FINALIZER_NAME          = "database.example.com/finalizer"
)

// DatabaseReconciler reconciles a Database object

type DatabaseReconciler struct {
    client.Client
    Scheme *runtime.Scheme
    Log    logr.Logger
}

//+kubebuilder:rbac:groups=database.example.com,resources=databases,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups=database.example.com,resources=databases/status,verbs=get;update;patch
//+kubebuilder:rbac:groups=database.example.com,resources=databases/finalizers,verbs=update
//+kubebuilder:rbac:groups=apps,resources=statefulsets,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups="",resources=services;configmaps;secrets,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups="",resources=persistentvolumeclaims,verbs=get;list;watch;create;update;patch
```

```
//+kubebuilder:rbac:groups="",resources=events,verbs=create;patch

// Reconcile implements the main reconciliation logic comparing desired vs actual state

func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
    log := r.Log.WithValues("database", req.NamespacedName)

    // TODO 1: Fetch the Database resource from the API server
    // Hint: Use r.Client.Get() with the req.NamespacedName

    // Handle the case where the resource might have been deleted (errors.NotFound)

    // TODO 2: Handle deletion case if DeletionTimestamp is set
    // Check if database.DeletionTimestamp != nil

    // If deletion is in progress, call r.handleDeletion() and return

    // TODO 3: Add finalizer if not present
    // Check if FINALIZER_NAME is in database.Finalizers

    // If missing, add it using controllerutil.AddFinalizer and update the resource

    // TODO 4: Reconcile owned resources to match desired state
    // Call r.reconcileStatefulSet(), r.reconcileService(), r.reconcileConfigMap()
    // Each function should return error if reconciliation fails

    // TODO 5: Update status based on current state of owned resources
    // Fetch current StatefulSet status and update database.Status.ReadyReplicas
    // Set appropriate conditions in database.Status.Conditions
    // Use r.Status().Update() to persist status changes

    // TODO 6: Determine requeue behavior
    // Return ctrl.Result{RequeueAfter: REQUEUE_AFTER_DURATION} for periodic reconciliation
    // Return ctrl.Result{} with no requeue if everything is in desired state
    // Return error if reconciliation failed and should be retried with backoff

    return ctrl.Result{}, nil
}
```

```
// handleDeletion processes Database resource deletion with cleanup

func (r *DatabaseReconciler) handleDeletion(ctx context.Context, database *databasev1.Database) error {
    log := r.Log.WithValues("database", database.Name, "namespace", database.Namespace)

    // TODO 1: Perform cleanup operations before allowing deletion
    // This might include draining connections, backing up data, etc.
    // For this implementation, we rely on owner references for cascade deletion

    // TODO 2: Remove finalizer to allow deletion to proceed
    // Use controllerutil.RemoveFinalizer() and update the resource
    // Return any errors from the update operation

    return nil
}

// reconcileStatefulSet ensures the StatefulSet matches the Database spec

func (r *DatabaseReconciler) reconcileStatefulSet(ctx context.Context, database *databasev1.Database) error {
    // TODO 1: Define desired StatefulSet based on Database spec
    // Use resources.NewStatefulSet() helper to create desired state
    // Set replica count from database.Spec.Replicas
    // Set container image from database.Spec.Version
    // Configure storage from database.Spec.StorageSize

    // TODO 2: Check if StatefulSet already exists
    // Use r.Client.Get() to fetch existing StatefulSet
    // Handle both "exists" and "not found" cases

    // TODO 3: Create StatefulSet if it doesn't exist
    // Set owner reference using resources.SetOwnerReference()
    // Use r.Client.Create() to create the StatefulSet

    // TODO 4: Update StatefulSet if it exists but doesn't match desired state
    // Compare existing spec with desired spec
    // Use r.Client.Update() to apply changes
```

```

    // Handle update conflicts by refetching and retrying

    return nil
}

// SetupWithManager registers the controller with the manager and configures watches

func (r *DatabaseReconciler) SetupWithManager(mgr ctrl.Manager) error {
    return ctrl.NewControllerManagedBy(mgr).

        For(&databasev1.Database{}).                                // Primary resource to watch
        Owns(&appsv1.StatefulSet{}).                               // Watch owned StatefulSets
        Owns(&corev1.Service{}).                                 // Watch owned Services
        Owns(&corev1.ConfigMap{}).                               // Watch owned ConfigMaps
        Owns(&corev1.Secret{}).                                 // Watch owned Secrets
        Complete(r)

}

```

Language-Specific Hints

Go-Specific Implementation Guidance:

- Use `sigs.k8s.io/controller-runtime/pkg/controller/controllerutil` for finalizer management (`AddFinalizer`, `RemoveFinalizer`)
- Import `k8s.io/apimachinery/pkg/api/errors` for handling API server errors like `NotFound`, `Conflict`
- Use `context.Context` throughout for cancellation and timeouts - never ignore the context parameter
- Leverage `k8s.io/apimachinery/pkg/types.NamespacedName` for resource identification in client calls
- Use `sigs.k8s.io/controller-runtime/pkg/log` for structured logging with key-value pairs
- Import specific API group packages like `appsv1 "k8s.io/api/apps/v1"` to avoid conflicts

Client Usage Patterns:

- Use `r.Client.Get(ctx, namespaceName, &resource)` for fetching resources
- Use `r.Client.Create(ctx, &resource)` for creating new resources
- Use `r.Client.Update(ctx, &resource)` for updating existing resources
- Use `r.Status().Update(ctx, &resource)` for status-only updates
- Always check `errors.NotFound(err)` when fetching resources that might not exist

Milestone Checkpoint

After implementing the goals and non-goals framework:

Validation Command:

```
# Verify project structure matches recommendations
find . -name "*.go" | head -10

# Check that generated code compiles
go mod tidy && go build ./cmd/manager

# Verify RBAC markers are present
grep -r "kubebuilder:rbac" controllers/
```

BASH

Expected Behavior:

- Project compiles without errors
- Manager binary starts (may exit due to missing CRDs, but should parse flags)
- RBAC markers generate appropriate permissions
- File structure follows recommended organization

Signs of Problems:

- Import cycle errors → Check for circular dependencies between packages
- Missing scheme registration → Verify `AddToScheme` calls in main.go
- RBAC permission errors → Check that controller has necessary resource access

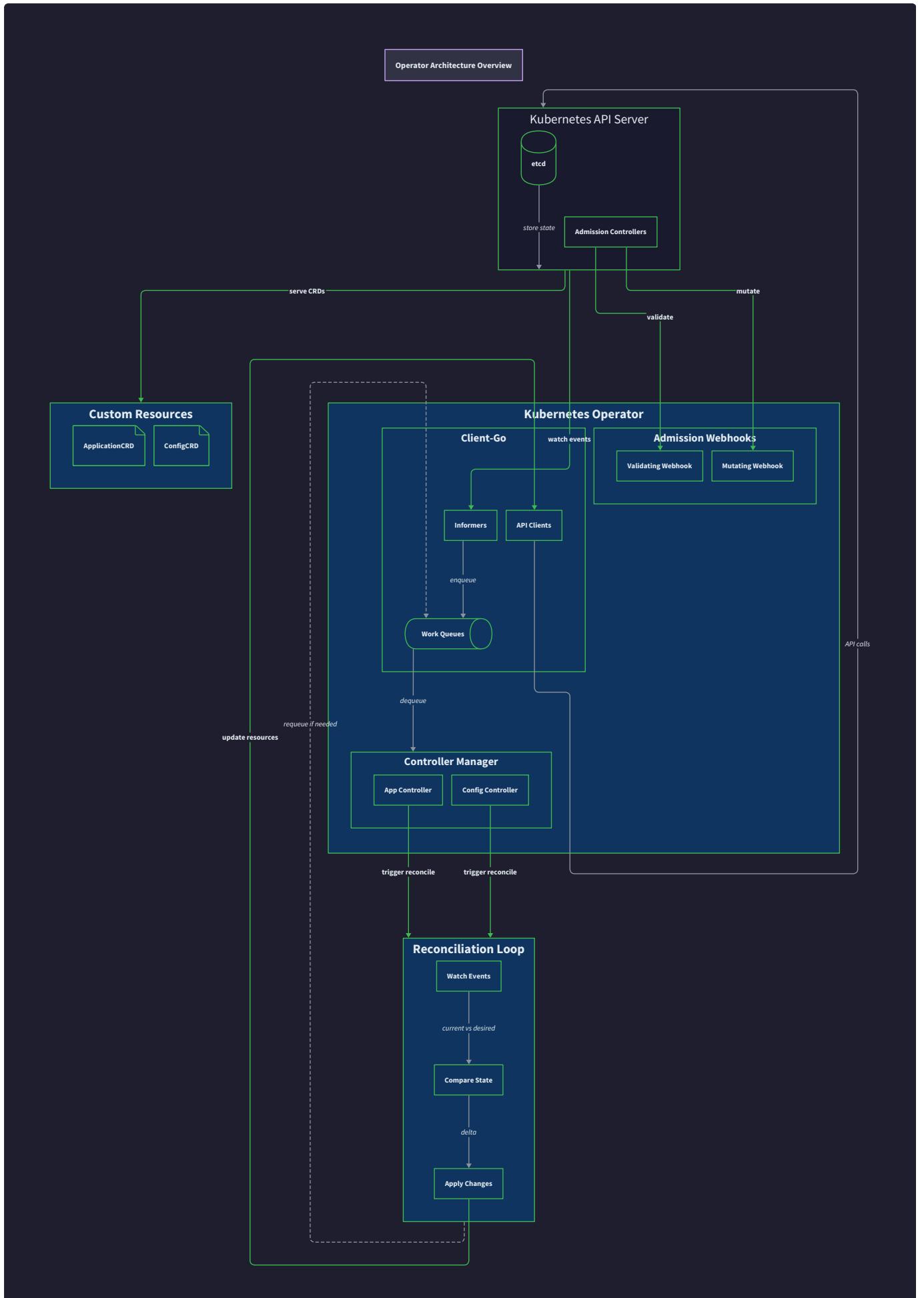
This implementation guidance translates the high-level goals into actionable development steps while maintaining the boundaries defined in the non-goals section.

High-Level Architecture

Milestone(s): Milestone 1 (Custom Resource Definition), Milestone 2 (Controller Setup) - this section establishes the foundational architecture that guides CRD design and controller implementation

Building a Kubernetes operator is like constructing a sophisticated automation system that integrates seamlessly with an existing orchestration platform. Think of it as adding a specialized department to a large corporation - the new department must understand the company's communication protocols, follow established procedures, and coordinate effectively with existing teams while bringing its own domain expertise to automate previously manual processes.

The operator architecture consists of four primary components that work together to extend Kubernetes with domain-specific automation capabilities. These components - custom resources, controllers, admission webhooks, and supporting infrastructure - form a cohesive system that integrates with the Kubernetes API machinery while maintaining separation of concerns and operational reliability.



Component Overview

The operator follows the standard Kubernetes extension pattern, leveraging the platform's built-in extensibility mechanisms to add new resource types and automated management logic. Understanding how these components interact helps establish the mental model for implementing each piece correctly.

Custom Resources serve as the API contract between users and the operator. These are not just data structures - they represent the declarative interface through which users express their desired state for managed applications. The `Database` custom resource, for example, encapsulates all the configuration parameters needed to describe a database instance, from basic settings like replicas and storage size to advanced configuration like backup schedules and monitoring endpoints.

The custom resource design follows Kubernetes conventions with distinct `spec` and `status` subresources. The `spec` contains the user's desired configuration and remains immutable once set (except for explicit updates). The `status` reflects the current observed state of the managed resources and gets updated continuously by the controller as it reconciles the actual state with the desired state.

Custom Resource Component	Purpose	Ownership	Update Pattern
<code>DatabaseSpec</code>	User's desired state declaration	User via kubectl/API	Explicit updates only
<code>DatabaseStatus</code>	Controller's observed state report	Controller exclusively	Continuous reconciliation updates
Metadata annotations	Operational hints and user labels	Mixed (user + controller)	Additive updates
Finalizers	Cleanup coordination mechanism	Controller exclusively	Added on creation, removed on cleanup completion

Controllers implement the reconciliation logic that continuously works to make the actual state match the desired state declared in custom resources. Think of a controller as a tireless system administrator who never sleeps, constantly checking that everything is configured correctly and taking corrective action when things drift from the desired configuration.

The `DatabaseReconciler` watches for changes to `Database` custom resources and responds by examining the current state of all related Kubernetes resources (StatefulSets, Services, ConfigMaps, Secrets) and taking whatever actions are necessary to align them with the specification. This might involve creating new resources, updating existing ones, or cleaning up resources that are no longer needed.

Controllers operate through an event-driven architecture built on Kubernetes' watch mechanism. Rather than polling for changes, controllers receive notifications when resources are created, updated, or deleted. This allows them to respond quickly to changes while minimizing unnecessary API server load.

Controller Component	Function	Input	Output
<code>Informer</code>	Maintains local cache of watched resources	API server events	Cached resource state
<code>WorkQueue</code>	Buffers reconciliation requests with rate limiting	Resource change events	Serialized reconcile requests
<code>Reconciler</code>	Compares desired vs actual state and takes action	Reconcile request	Resource mutations + status updates
<code>EventRecorder</code>	Provides observability into controller actions	Controller decisions	Kubernetes events

Admission Webhooks provide validation and mutation capabilities that execute during the resource admission process, before resources are stored in etcd. These webhooks act as gatekeepers, ensuring that only valid resources enter the system and automatically applying default values or transformations to improve the user experience.

Validation webhooks implement business rules that go beyond what can be expressed in OpenAPI schema validation. For example, while the schema can enforce that a storage size is a valid quantity, a validation webhook can enforce that storage size increases are allowed but decreases are prohibited (since they could cause data loss).

Mutating webhooks enhance the user experience by automatically filling in reasonable defaults and applying organizational policies. When a user creates a `Database` resource with minimal configuration, the mutating webhook can inject appropriate resource limits, security contexts, and monitoring configurations based on organizational standards.

Webhook Type	Execution Phase	Purpose	Common Use Cases
Mutating	Before validation and storage	Apply defaults and transformations	Inject resource limits, add labels, set security policies
Validating	After mutation, before storage	Enforce business rules and constraints	Validate field combinations, check quotas, verify dependencies

Supporting Infrastructure provides the operational foundation that enables the core components to function reliably in production environments. This includes observability tools, security configurations, and deployment mechanisms that ensure the operator can be operated safely at scale.

The infrastructure layer handles cross-cutting concerns like leader election (ensuring only one controller replica actively reconciles resources), TLS certificate management for webhook endpoints, RBAC configuration for secure API access, and monitoring integrations that provide visibility into operator health and performance.

Kubernetes API Integration

The operator extends Kubernetes by leveraging the platform's built-in extensibility mechanisms rather than working around them. This integration approach ensures that the operator feels like a natural part of the Kubernetes ecosystem, with custom resources appearing in `kubectl get` output, audit logs, and administrative tools just like built-in resources.

API Server Extension through Custom Resource Definitions allows the operator to register new resource types that become first-class citizens in the Kubernetes API. When a Custom Resource Definition (CRD) is applied to a cluster, the API server automatically provides REST endpoints for the new resource type, complete with validation, versioning, and all the standard Kubernetes API features.

The CRD registration process creates API endpoints following Kubernetes conventions:

`/api/v1/namespaces/{namespace}/databases` for namespaced resources or `/api/v1/databases` for cluster-scoped resources.

These endpoints support all standard HTTP verbs (GET, POST, PUT, PATCH, DELETE) and integrate with existing tooling like `kubectl`, client libraries, and administrative dashboards.

Decision: Namespaced vs Cluster-scoped Resources

- **Context:** Database resources could be deployed at namespace level (tenant isolation) or cluster level (shared infrastructure)
- **Options Considered:** Namespaced only, cluster-scoped only, hybrid with both types
- **Decision:** Namespaced resources with cluster-scoped operator
- **Rationale:** Namespaced resources provide tenant isolation and align with Kubernetes RBAC patterns, while cluster-scoped operator deployment simplifies operational management
- **Consequences:** Enables multi-tenancy and role-based access control, but requires cluster admin privileges for operator installation

Controller Integration with the Control Plane follows the established controller pattern used by Kubernetes itself. The operator's controller uses the same client-go library that built-in controllers use, ensuring consistent behavior and reliability characteristics.

The controller registers with the API server to watch specific resource types and receives a stream of events (create, update, delete) through efficient long-polling connections. When the connection is interrupted, the client-go library automatically reconnects and reestablishes the watch, ensuring the controller doesn't miss events during network partitions or API server restarts.

Controllers maintain an eventually consistent local cache of relevant cluster state through informers. This cache serves two critical purposes: it reduces load on the API server by serving read requests locally, and it provides a consistent snapshot of related resources during reconciliation, preventing race conditions that could occur if the controller made individual API calls for each resource.

API Integration Component	Purpose	Kubernetes Integration Point	Reliability Features
Custom Resource Definition	Extends API with new resource types	API server schema registration	Version conversion, validation, storage
Controller Watch Streams	Receives resource change notifications	API server watch mechanism	Automatic reconnection, bookmark events
Informer Cache	Local eventually-consistent state replica	List/watch with resync	Cache invalidation, index optimization
Admission Webhooks	Intercepts resource mutations	API server admission chain	Timeout handling, failure policies

Admission Webhook Integration occurs within the API server's request processing pipeline, specifically during the admission phase after authentication and authorization but before persistence to etcd. This timing allows webhooks to validate and transform requests while ensuring that security policies have already been enforced.

The API server sends admission review requests to webhook endpoints over HTTPS, providing the full context of the operation including the resource being created or updated, the user making the request, and relevant metadata. Webhooks respond with admission review responses that either allow the operation (possibly with modifications) or deny it with an explanatory message.

Webhook integration requires careful attention to reliability and security. The API server enforces timeout limits on webhook calls to prevent cluster availability issues, and webhook failure policies determine whether failures result in request denial (fail-closed) or request acceptance (fail-open). TLS certificate management ensures secure communication and prevents man-in-the-middle attacks.

Critical Design Insight: The operator doesn't replace Kubernetes mechanisms - it extends them. Every interaction goes through standard Kubernetes APIs, ensuring compatibility with existing tools, security policies, and operational procedures.

Event Recording and Observability integrate with Kubernetes' standard observability mechanisms to provide visibility into operator actions. When the controller takes actions like creating StatefulSets or handling errors, it records events that appear in `kubectl describe` output and cluster monitoring systems.

The controller uses the Kubernetes event recording mechanism to document its decision-making process, providing an audit trail of actions taken during reconciliation. These events help both users and operators understand why certain actions were taken and provide debugging information when things don't work as expected.

Recommended Project Structure

Organizing the operator codebase follows established Go and Kubernetes conventions that promote maintainability, testing, and collaboration. The structure separates concerns clearly while keeping related functionality together, making it easier for team members to understand and modify the system.

Top-Level Organization follows the standard Go project layout with `cmd/` for executable entry points, `internal/` for private implementation code, and `api/` for public API definitions. This structure makes it immediately clear what parts of the codebase are intended for external consumption versus internal implementation details.

```

database-operator/
├── cmd/
│   └── manager/
│       └── main.go           ← Operator entry point with flag parsing and setup
├── api/
│   └── v1/
│       ├── database_types.go    ← Database custom resource definition
│       ├── groupversion_info.go  ← API group and version metadata
│       └── zz_generated.deepcopy.go ← Generated code (do not edit)
└── internal/
    └── controller/
        ├── database_controller.go ← DatabaseReconciler implementation
        └── database_controller_test.go ← Controller unit tests
config/
├── crd/
|   └── bases/                  ← Generated CRD manifests
├── rbac/                      ← RBAC role definitions
├── webhook/                   ← Webhook configuration manifests
└── manager/                   ← Controller manager deployment
└── webhooks/
    ├── database_defaulting.go  ← Mutating webhook implementation
    ├── database_validation.go  ← Validating webhook implementation
    └── webhook_suite_test.go   ← Webhook integration tests

```

API Package Organization separates the public API definitions from the implementation logic, making it easier to generate client libraries and ensuring that API changes are deliberate and well-considered. The `api/v1/` package contains only type definitions and generated code, with no business logic.

The API package structure supports versioning from the beginning, even if only one version exists initially. This forward-thinking approach makes it much easier to add new API versions later without requiring major refactoring. The versioned packages also align with Kubernetes' own API organization patterns.

Directory	Contents	Ownership	Generation
<code>api/v1/</code>	Type definitions and API metadata	Manual implementation	Partial (deepcopy methods)
<code>config/crd/bases/</code>	CRD YAML manifests	Generated from Go types	Fully generated
<code>config/rbac/</code>	Role and RoleBinding manifests	Manual configuration	Manual
<code>config/webhook/</code>	Webhook configuration YAML	Generated with manual TLS setup	Mixed

Internal Package Structure organizes implementation code by function rather than by layer, keeping related functionality together. The `internal/controller/` package contains all controller-related logic, while webhook implementations live in a separate `webhooks/` package at the top level (since they may be shared across multiple operators).

Controller organization separates the reconciliation logic from the Kubernetes integration machinery. The controller file focuses on the business logic of comparing desired state with actual state and deciding what actions to take. Integration concerns like informer setup, work queue configuration, and manager registration are handled in setup functions that keep the core reconciliation logic clean and testable.

Decision: Single Package vs Multi-Package Controller Organization

- **Context:** Controller code can be organized in one package or split across multiple packages by concern
- **Options Considered:** Single controller package, separate packages for reconciler/setup/types, layered packages
- **Decision:** Single controller package with clear function separation
- **Rationale:** Reduces cognitive overhead for small-to-medium operators, keeps related code together, aligns with controller-runtime patterns
- **Consequences:** May need refactoring if controller grows very large, but provides simpler navigation and testing

Configuration and Deployment Structure separates generated manifests from manually maintained configuration, making it clear what should be edited directly versus what gets regenerated from code. The `config/` directory contains all Kubernetes manifests needed to deploy the operator, organized by function.

The configuration structure supports both development and production deployments through Kustomize overlays. The base configuration in `config/` provides a working deployment that can be customized for specific environments without modifying the original files. This approach enables GitOps workflows where environment-specific configurations are maintained in separate repositories or directories.

Testing Organization co-locates unit tests with the code they test while providing separate integration test suites that exercise the full system. Controller unit tests focus on the reconciliation logic using fake clients, while integration tests in the `webhooks/` package exercise the full admission webhook flow including TLS certificate handling.

The testing structure supports both fast unit testing during development and comprehensive integration testing in CI/CD pipelines. Unit tests can run without a Kubernetes cluster and complete in seconds, while integration tests use envtest to run against a real API server but still complete quickly enough for regular execution.

Test Type	Location	Dependencies	Execution Time
Unit Tests	<code>*_test.go</code> next to implementation	Fake clients only	< 1 second per test
Controller Integration	<code>internal/controller/suite_test.go</code>	envtest with real API server	< 10 seconds per test
Webhook Integration	<code>webhooks/webhook_suite_test.go</code>	envtest with webhook server	< 15 seconds per test
End-to-End Tests	<code>test/e2e/</code>	Full cluster deployment	1-5 minutes per test

Common Pitfalls

⚠ Pitfall: Mixing API definitions with implementation logic - New operator developers often put business logic directly in the API package alongside type definitions. This creates circular dependencies and makes it difficult to generate clean client libraries. Keep the `api/` package focused solely on type definitions and generated code. Business logic belongs in `internal/controller/` or other implementation packages.

⚠ Pitfall: Ignoring the controller-runtime project structure - The controller-runtime framework expects a specific project organization, and deviating from it breaks code generation and scaffolding tools. Use `kubebuilder init` to create the initial project structure and follow the conventions it establishes. This ensures compatibility with the broader Kubernetes tooling ecosystem.

⚠ Pitfall: Not planning for API versioning from the start - Even if you only need one API version initially, organizing code as if multiple versions will exist makes future evolution much easier. Create the `api/v1/` package structure immediately rather than using an unversioned `api/` package that will require refactoring later.

⚠ Pitfall: Coupling controllers too tightly to specific custom resources - Controllers should be designed to work with interface types when possible, making them more testable and reusable. Avoid directly importing custom resource types in controller business logic; instead, define interfaces that capture the required behavior and use type assertions only at the boundaries.

Implementation Guidance

The architecture implementation follows cloud-native patterns and leverages the controller-runtime framework to minimize boilerplate code while maintaining flexibility for customization.

Technology Recommendations

Component	Simple Option	Advanced Option
Project Scaffolding	Kubebuilder v3 with default settings	Operator SDK with custom templates
API Code Generation	controller-gen with basic markers	controller-gen with comprehensive OpenAPI markers
Testing Framework	Ginkgo/Gomega with envtest	Ginkgo/Gomega with real cluster integration
Configuration Management	Kustomize with simple overlays	Helm with complex value templating
Observability	Standard Go logging with controller-runtime	Structured logging with Prometheus metrics

File Structure Setup

Initialize the project structure using kubebuilder to ensure alignment with Kubernetes conventions:

```
# Initialize project with domain and repository                                BASH
kubebuilder init --domain example.com --repo github.com/yourorg/database-operator

# Create API and controller scaffolding
kubebuilder create api --group database --version v1 --kind Database

# Create webhook scaffolding
kubebuilder create webhook --group database --version v1 --kind Database --defaulting --programmatic-validation
```

Core Type Definitions

```
// api/v1/database_types.go                                     GO

package v1

import (
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
)

// DatabaseSpec defines the desired state of Database

type DatabaseSpec struct {

    // Replicas is the number of database instances to run
    // +kubebuilder:validation:Minimum=1
    // +kubebuilder:validation:Maximum=10
    // +kubebuilder:default=1

    Replicas int32 `json:"replicas,omitempty"`

    // Version specifies the database version to deploy
    // +kubebuilder:validation:Pattern="^\\d{1,3}\\.\\d{1,3}\\.\\d{1,3}$"

    Version string `json:"version"`

    // StorageSize defines the persistent storage size for each replica
    // +kubebuilder:validation:Pattern="^\\d{1,3}([KMGT]{1})$"

    StorageSize string `json:"storageSize"`

    // TODO: Add configuration fields like backup schedule, monitoring settings
    // TODO: Add networking configuration for service exposure
    // TODO: Add security settings like TLS and authentication
}

// DatabaseStatus defines the observed state of Database

type DatabaseStatus struct {

    // Conditions represent the latest available observations of the database state
    Conditions []metav1.Condition `json:"conditions,omitempty"`

    // ReadyReplicas is the number of database instances that are ready to serve requests
}
```

```

ReadyReplicas int32 `json:"readyReplicas,omitempty"`

// ObservedGeneration reflects the generation of the most recently observed Database spec
ObservedGeneration int64 `json:"observedGeneration,omitempty"`

// TODO: Add status fields for backup status, connection endpoints
// TODO: Add metrics like storage usage, connection count
}

// +kubebuilder:object:root=true
// +kubebuilder:subresource:status
// +kubebuilder:printcolumn:name="Replicas",type="integer",JSONPath=".spec.replicas"
// +kubebuilder:printcolumn:name="Ready",type="integer",JSONPath=".status.readyReplicas"
// +kubebuilder:printcolumn:name="Version",type="string",JSONPath=".spec.version"
// +kubebuilder:printcolumn:name="Age",type="date",JSONPath=".metadata.creationTimestamp"

// Database is the Schema for the databases API

type Database struct {

    metav1.TypeMeta `json:",inline"`

    metav1.ObjectMeta `json:"metadata,omitempty"`

    Spec   DatabaseSpec   `json:"spec,omitempty"`
    Status DatabaseStatus `json:"status,omitempty"`
}

// +kubebuilder:object:root=true

// DatabaseList contains a list of Database

type DatabaseList struct {

    metav1.TypeMeta `json:",inline"`

    metav1.ListMeta `json:"metadata,omitempty"`

    Items          []Database `json:"items"`
}

func init() {
    SchemeBuilder.Register(&Database{}, &DatabaseList{})
}

```

}

Controller Foundation

```
// internal/controller/database_controller.go                                     GO

package controller

import (
    "context"
    "time"

    "k8s.io/apimachinery/pkg/runtime"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/log"

    databasev1 "github.com/yourorg/database-operator/api/v1"
)

const (
    REQUEUE_AFTER_DURATION = 30 * time.Second
    FINALIZER_NAME = "database.example.com/finalizer"
)

// DatabaseReconciler reconciles a Database object

type DatabaseReconciler struct {
    client.Client
    Scheme *runtime.Scheme
}

// +kubebuilder:rbac:groups=database.example.com,resources=databases,verbs=get;list;watch;create;update;patch;delete
// +kubebuilder:rbac:groups=database.example.com,resources=databases/status,verbs=get;update;patch
// +kubebuilder:rbac:groups=database.example.com,resources=databases/finalizers,verbs=update
// +kubebuilder:rbac:groups=apps,resources=statefulsets,verbs=get;list;watch;create;update;patch;delete
// +kubebuilder:rbac:groups="",resources=services,verbs=get;list;watch;create;update;patch;delete
// +kubebuilder:rbac:groups="",resources=configmaps,verbs=get;list;watch;create;update;patch;delete

// Reconcile compares the state specified by the Database object against the actual cluster state,
// and performs operations to make the cluster state reflect the state specified by the user.
```

```
func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
    log := log.FromContext(ctx)

    // TODO 1: Fetch the Database instance from the cluster

    // TODO 2: Handle deletion by checking if DeletionTimestamp is set

    // TODO 3: Add finalizer if not present to ensure cleanup

    // TODO 4: Reconcile StatefulSet to match spec.replicas and spec.version

    // TODO 5: Reconcile Service for database connectivity

    // TODO 6: Reconcile ConfigMap with database configuration

    // TODO 7: Update Database status with current ready replicas and conditions

    // TODO 8: Return appropriate requeue strategy based on current state

    log.Info("Reconciling Database", "database", req.NamespacedName)

    return ctrl.Result{RequeueAfter: REQUEUE_AFTER_DURATION}, nil
}

// SetupWithManager sets up the controller with the Manager.

func (r *DatabaseReconciler) SetupWithManager(mgr ctrl.Manager) error {
    return ctrl.NewControllerManagedBy(mgr).
        For(&databasev1.Database{}).

        // TODO: Add watches for owned resources (StatefulSet, Service, ConfigMap)

        // TODO: Configure predicates to filter unnecessary reconciliation triggers

        // TODO: Set up event handler for better observability

        Complete(r)
}
```

Webhook Skeleton

```
// webhooks/database_validation.go                                     GO

package webhooks

import (
    "context"
    "fmt"

    "k8s.io/apimachinery/pkg/runtime"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/webhook"
    "sigs.k8s.io/controller-runtime/pkg/webhook/admission"

    databasev1 "github.com/yourorg/database-operator/api/v1"
)

// +kubebuilder:webhook:path=/validate-database-example-com-v1-
// database,mutating=false,failurePolicy=fail,groups=database.example.com,resources=databases,verbs=create;update,vers
// ions=v1,name=vdatabase.kb.io,sideEffects=None,admissionReviewVersions=v1

// DatabaseValidator validates Database resources using custom business rules

type DatabaseValidator struct{}


// ValidateCreate implements webhook.Validator
func (v *DatabaseValidator) ValidateCreate(ctx context.Context, obj runtime.Object) (admission.Warnings, error) {
    database := obj.(*databasev1.Database)

    // TODO 1: Validate that storage size is not below minimum threshold
    // TODO 2: Validate that version string matches supported database versions
    // TODO 3: Validate that replica count doesn't exceed cluster capacity
    // TODO 4: Check for naming conflicts with existing databases

    return nil, nil
}

// ValidateUpdate implements webhook.Validator
func (v *DatabaseValidator) ValidateUpdate(ctx context.Context, oldObj, newObj runtime.Object) (admission.Warnings, error) {
```

```

newDatabase := newObj.(*databasev1.Database)

oldDatabase := oldObj.(*databasev1.Database)

// TODO 1: Prevent storage size decreases that could cause data loss

// TODO 2: Validate version upgrade paths (no downgrades, supported migrations)

// TODO 3: Ensure replica count changes are within acceptable bounds

return nil, nil
}

// ValidateDelete implements webhook.Validator

func (v *DatabaseValidator) ValidateDelete(ctx context.Context, obj runtime.Object) (admission.Warnings, error) {
    // TODO 1: Check if database contains data that should be backed up

    // TODO 2: Validate that dependent applications are prepared for deletion

    return nil, nil
}

func (v *DatabaseValidator) SetupWebhookWithManager(mgr ctrl.Manager) error {
    return ctrl.NewWebhookManagedBy(mgr).
        For(&databasev1.Database{}).
        WithValidator(v).
        Complete()
}

```

Language-Specific Development Tips

Go-Specific Implementation Guidance:

- Use `sigs.k8s.io/controller-runtime/pkg/client` for all Kubernetes API interactions - it provides automatic retries and rate limiting
- Leverage `k8s.io/apimachinery/pkg/api/errors` for proper error handling that distinguishes between different failure types
- Use `controllerutil.SetControllerReference` to establish owner relationships that enable garbage collection
- Implement the `Defaulter` and `Validator` interfaces on your custom resource types for webhook functionality
- Use `k8s.io/utils/pointer` for safely creating pointers to primitive values in Kubernetes objects

RBAC Configuration Guidelines: The generated RBAC markers in the controller comments will create appropriate permissions, but verify that they include:

- Read/write access to your custom resources and their status subresources
- Read/write access to all Kubernetes resources your controller manages (StatefulSets, Services, etc.)

- Read access to resources you need to watch but don't own (Nodes, StorageClasses for validation)

Debugging and Development Workflow:

1. Run `make generate` after modifying API types to regenerate deepcopy methods
2. Run `make manifests` to regenerate CRDs and RBAC from code annotations
3. Use `make install` to apply CRDs to your development cluster
4. Test controller logic with `make run` for local development with real cluster access
5. Use `make docker-build` and `make deploy` for full integration testing in the cluster

Data Model and Custom Resources

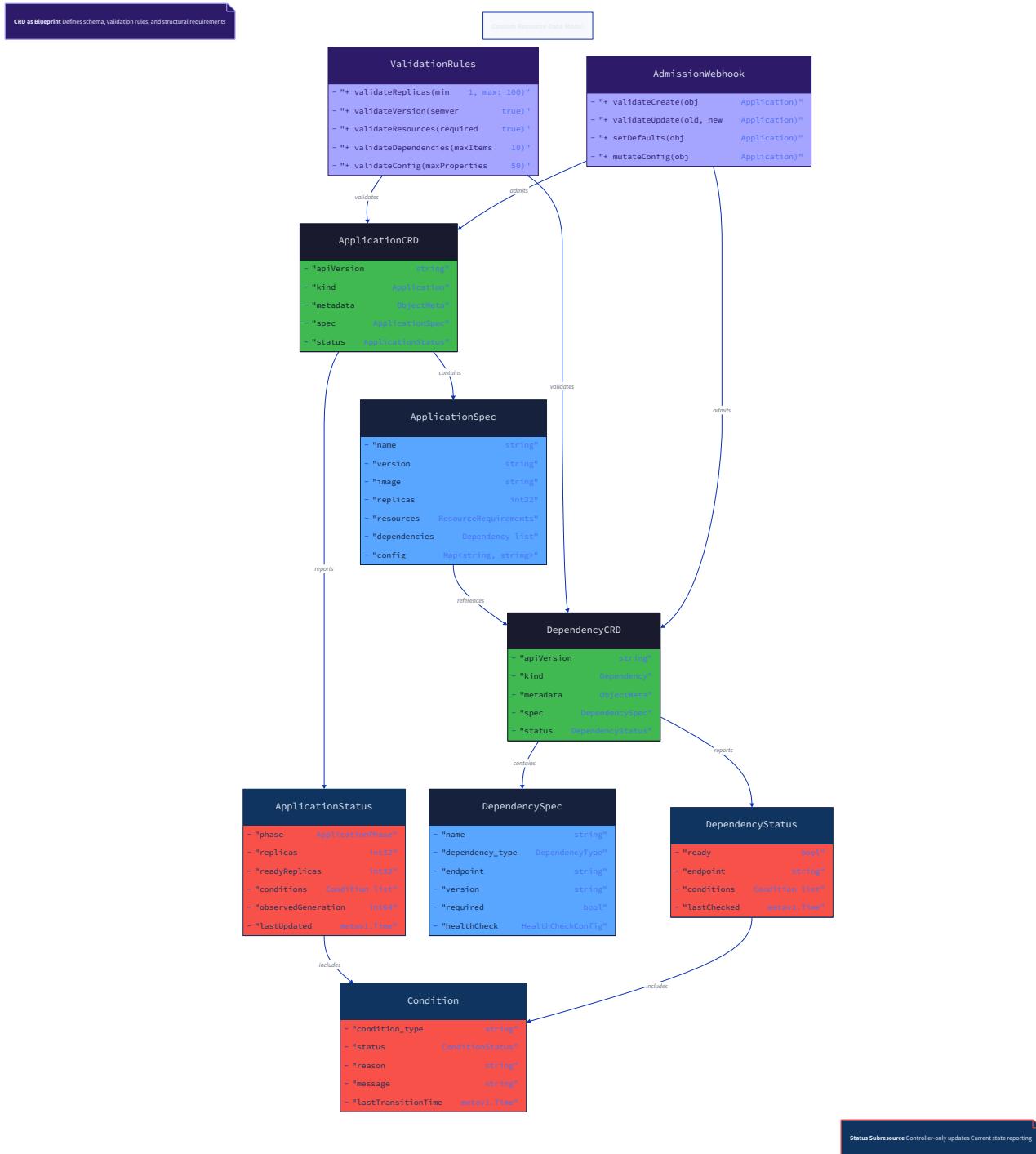
Milestone(s): Milestone 1 (Custom Resource Definition), Milestone 3 (Reconciliation Loop) - defines the API schema and status structures that the controller reconciles, and Milestone 4 (Webhooks) - establishes validation rules enforced by admission control

Mental Model: The Application Blueprint

Think of a **custom resource definition** as an architectural blueprint for buildings. Just as a blueprint defines the structure, dimensions, materials, and specifications that every building of that type must follow, a CRD defines the schema, validation rules, and structural requirements that every instance of your custom resource must conform to. The blueprint doesn't create buildings—it defines what valid buildings look like. Similarly, the CRD doesn't create application instances—it defines what valid application configurations look like.

When an architect designs a blueprint, they specify required elements (foundation, load-bearing walls), optional elements (balconies, decorative features), and validation rules (minimum ceiling height, maximum occupancy). The CRD schema works the same way: it specifies required fields that every resource must have, optional fields that provide additional configuration, and validation rules that prevent invalid configurations from being accepted.

The **status subresource** acts like a building inspector's report. While the blueprint (spec) describes what the building should look like, the inspector's report (status) describes the current state of construction: which floors are complete, what problems have been encountered, and whether the building is ready for occupancy. Just as the architect doesn't write the inspector's report, users don't directly modify the status—only the controller (acting as the inspector) updates the status based on what it observes.



The data model forms the contract between users who describe their desired application state and controllers that make that state reality. This contract must be precise enough to prevent ambiguous configurations, flexible enough to accommodate diverse use cases, and evolvable enough to support future requirements without breaking existing deployments.

Custom Resource Definition Schema

The OpenAPI v3 schema structure defines the shape and validation rules for custom resources. Unlike simple configuration files, CRD schemas provide strong typing, comprehensive validation, and automatic documentation generation. The schema serves multiple purposes: it validates user input at the API level, generates client code for type-safe programming, provides auto-completion in kubectl and IDEs, and documents the API for users and operators.

Database Custom Resource Structure:

Field	Type	Required	Description
<code>apiVersion</code>	string	Yes	Always "example.com/v1" for this resource type
<code>kind</code>	string	Yes	Always "Database" for this resource type
<code>metadata</code>	ObjectMeta	Yes	Standard Kubernetes metadata including name, namespace, labels, annotations
<code>spec</code>	DatabaseSpec	Yes	User-specified desired state configuration
<code>status</code>	DatabaseStatus	No	Controller-managed current state information

The separation between `spec` and `status` reflects Kubernetes' declarative model. Users declare their intent in the spec, while controllers report reality in the status. This separation enables optimistic concurrency control and clear ownership boundaries.

DatabaseSpec Field Definitions:

Field	Type	Required	Validation	Description
<code>replicas</code>	int32	Yes	minimum: 1, maximum: 10	Number of database instances to maintain
<code>version</code>	string	Yes	pattern: "^[0-9]+\.[0-9]+\.[0-9]+\$"	Semantic version of database software
<code>storageSize</code>	string	Yes	pattern: "^[0-9]+[KMGT]i?\${\$}"	Storage capacity using Kubernetes quantity format
<code>backupSchedule</code>	string	No	pattern: "^(@((annually yearly monthly weekly daily hourly reboot)) (@every (\d+(ns us μs ms s m h))+)((\d+ l)/*(\d+))? (\d+ l) ((\d+)? (\d+ l)/*(\d+))? (\d+ l)/*(\d+))?)\$"	Optional cron expression for automated backups
<code>resources</code>	ResourceRequirements	No	-	CPU and memory resource requests and limits
<code>config</code>	map[string]string	No	maxProperties: 20	Database-specific configuration parameters

The validation patterns serve as compile-time checks that prevent invalid configurations from entering the system. The `version` field uses semantic versioning to ensure upgrade compatibility. The `storageSize` field accepts Kubernetes quantity formats like "10Gi" or "500Mi" for consistency with other Kubernetes resources.

DatabaseStatus Field Definitions:

Field	Type	Required	Description
conditions	[]Condition	No	Array of condition objects describing the current state
readyReplicas	int32	No	Number of database instances currently serving traffic
observedGeneration	int64	No	Generation of spec that was last processed by controller
phase	string	No	High-level status: Pending, Running, Failed, Terminating
lastBackupTime	*metav1.Time	No	Timestamp of most recent successful backup
endpoints	[]string	No	Connection endpoints for client applications
message	string	No	Human-readable message describing current status

The `observedGeneration` field enables controllers to track whether they've processed the latest spec changes. When `metadata.generation` exceeds `status.observedGeneration`, the controller knows new changes are pending. The `conditions` array provides detailed status information that tools like `kubectl` can interpret and display.

Condition Structure for Status Reporting:

Field	Type	Required	Description
type	string	Yes	Condition type: Ready, Available, Progressing, Degraded
status	string	Yes	Condition status: True, False, Unknown
lastTransitionTime	metav1.Time	Yes	Time when condition last changed status
reason	string	Yes	Machine-readable reason code for the condition
message	string	Yes	Human-readable message explaining the condition
observedGeneration	int64	No	Generation when this condition was last updated

Conditions provide structured status information that both humans and machines can interpret. The `Ready` condition indicates whether the resource is fully operational. The `Available` condition indicates whether the resource is accepting traffic. The `Progressing` condition indicates whether changes are being applied. The `Degraded` condition indicates whether the resource is running but with reduced functionality.

Decision: Status Subresource Separation

- **Context:** Controllers need to update status independently from user spec changes to avoid conflicts
- **Options Considered:**
 1. Store status in spec alongside user fields
 2. Use separate status subresource with independent update path
 3. Store status in annotations or separate custom resources
- **Decision:** Use separate status subresource with independent update path
- **Rationale:** Prevents update conflicts between user changes and controller status updates, enables proper RBAC separation, follows Kubernetes conventions
- **Consequences:** Requires two API calls for create+status update, but eliminates race conditions and follows established patterns

Validation and Defaulting Strategy

Kubernetes provides multiple layers of validation and defaulting that work together to ensure resource integrity. These layers operate at different points in the request lifecycle, each serving specific purposes and handling different classes of validation problems.

Schema-Level Validation occurs first, using OpenAPI v3 schema definitions embedded in the CRD. This validation runs inside the API server and provides basic type checking, format validation, and structural constraints. Schema validation is fast and deterministic but limited to single-field constraints and simple cross-field relationships.

Admission Webhook Validation provides custom business logic validation that can examine the entire resource, query external systems, and enforce complex cross-field constraints. Validating admission webhooks receive the complete resource and can reject requests with detailed error messages. This validation is more flexible but requires additional infrastructure and introduces latency and availability dependencies.

Common Expression Language (CEL) Validation bridges the gap between schema validation and webhook validation. CEL expressions run inside the API server but can access multiple fields and perform complex logic. CEL provides better performance than webhooks while supporting more sophisticated validation than basic schema constraints.

Validation Strategy Table:

Validation Type	Performance	Complexity	Use Cases	Failure Impact
OpenAPI Schema	Fastest	Low	Type checking, format validation, required fields	Request rejected immediately
CEL Expressions	Fast	Medium	Cross-field validation, business rules	Request rejected immediately
Admission Webhooks	Slower	High	External validation, complex business logic	Request rejected or webhook timeout

Decision: Layered Validation Approach

- **Context:** Need comprehensive validation without sacrificing performance or availability
- **Options Considered:**
 1. Schema validation only for maximum performance
 2. Webhook validation only for maximum flexibility
 3. Layered approach combining schema, CEL, and selective webhook validation
- **Decision:** Layered approach combining all three validation mechanisms
- **Rationale:** Schema catches basic errors fast, CEL handles cross-field logic efficiently, webhooks reserved for complex cases requiring external data
- **Consequences:** More complex validation setup, but better user experience and system reliability

CEL Expression Examples for Database Resource:

Field Combination	CEL Expression	Purpose
Replicas and Storage	`self.replicas <= 3	
Version and Config	self.version.startsWith('5.') ? !has(self.config.legacy_mode) : true	Prevent legacy configuration on modern versions
Backup Schedule	has(self.backupSchedule) ? self.replicas >= 2 : true	Require multiple replicas when backups are enabled
Resource Limits	has(self.resources.limits.memory) ? self.resources.limits.memory >= self.resources.requests.memory : true	Ensure limits are not less than requests

CEL expressions have access to the entire resource through `self` and can use built-in functions for string manipulation, arithmetic, and collection operations. The expressions must return boolean values, where `true` means validation passes and `false` means the

request should be rejected.

Defaulting Strategy Implementation:

Field	Default Value	Condition	Rationale
<code>replicas</code>	1	Always applied if not specified	Single instance is minimum viable configuration
<code>version</code>	"5.7.0"	Only for new resources	Provides stable default for users who don't specify
<code>storageSize</code>	"10Gi"	Always applied if not specified	Reasonable default that works for development and testing
<code>resources.requests.cpu</code>	"100m"	When resources section exists but cpu not specified	Prevents resource starvation
<code>resources.requests.memory</code>	"256Mi"	When resources section exists but memory not specified	Reasonable baseline for database workloads

Defaulting occurs through mutating admission webhooks that modify incoming requests before validation. The webhook examines each field and injects appropriate defaults based on the current resource state and user-provided configuration. Defaults are applied only when fields are completely missing, not when they are explicitly set to empty values.

Validation Webhook Request Processing:

- Request Reception:** Webhook receives AdmissionReview containing the resource and operation type (CREATE, UPDATE)
- Context Extraction:** Extract the Database resource from the admission request payload
- Business Rule Validation:** Apply complex validation rules that span multiple fields or require external data
- Cross-Resource Validation:** Query API server for related resources if validation depends on cluster state
- Response Generation:** Create AdmissionResponse with allowed/denied decision and detailed messages
- Error Handling:** Return structured error responses that kubectl can display meaningfully to users

The webhook must complete validation within the configured timeout (typically 10-30 seconds) or the API server will fail the request. This timeout constraint requires careful design of validation logic to avoid external dependencies that could cause unpredictable delays.

API Versioning and Evolution

Kubernetes custom resources evolve over time as requirements change, new features are added, and operational experience reveals design improvements. The API versioning system provides backward compatibility while enabling forward progress. Unlike traditional API versioning that maintains separate endpoints, Kubernetes CRDs support multiple versions of the same resource type with automatic conversion between versions.

Version Evolution Mental Model: Think of API versioning like building renovation. When you renovate a house, you don't build a completely separate house—you modify the existing structure while keeping it livable. You might add new rooms (fields), remove outdated fixtures (deprecated fields), or reorganize the layout (restructure data). During renovation, you need temporary bridges or conversion tools so people can move between the old layout and new layout. Similarly, API versioning adds new fields and restructures existing ones while providing conversion mechanisms so existing clients continue working.

Versioning Strategy Overview:

Version	Status	Purpose	Conversion Required
v1alpha1	Deprecated	Initial experimental API	Yes - to v1beta1 and v1
v1beta1	Active	Stable API with minor additions	Yes - to v1
v1	Storage Version	Production-ready API	No - other versions convert to this

The **storage version** is the canonical representation used in etcd. All other versions are converted to and from the storage version when resources are created, updated, or retrieved. This hub-and-spoke conversion model simplifies the conversion matrix compared to requiring direct conversion between every version pair.

Decision: Hub and Spoke Conversion Model

- **Context:** Need to support multiple API versions without exponential conversion complexity
- **Options Considered:**
 1. Direct conversion between every version pair (N^2 complexity)
 2. Hub and spoke with storage version as hub ($2N$ complexity)
 3. Version deprecation forcing users to migrate immediately
- **Decision:** Hub and spoke conversion model with v1 as storage version
- **Rationale:** Minimizes conversion code maintenance, provides single source of truth, enables gradual migration
- **Consequences:** All versions must be convertible to/from storage version, but total complexity remains manageable

Version Migration Example: v1alpha1 to v1beta1:

v1alpha1 Field	v1beta1 Field	Conversion Logic
size (string)	replicas (int32)	Parse string number, convert to int32
dbVersion (string)	version (string)	Direct field rename, no transformation
storage (string)	storageSize (string)	Direct field rename, no transformation
backup (bool)	backupSchedule (string)	If true, default to "@daily"; if false, omit field
N/A	resources (ResourceRequirements)	New field, populate with empty object

Conversion Webhook Architecture:

The conversion webhook acts as a translator between API versions. When a user requests a resource in version A but it's stored in version B, the API server sends the stored version to the conversion webhook, which transforms it to the requested version. This process is transparent to clients—they see the version they requested without knowing about storage version differences.

Conversion Webhook Processing Steps:

1. **Request Analysis:** Webhook receives ConversionReview with desired version and current resource data
2. **Version Detection:** Determine source version from resource apiVersion field and target version from request
3. **Conversion Path Selection:** Choose appropriate conversion function based on source and target versions
4. **Field Mapping:** Transform fields according to version-specific mapping rules
5. **Validation:** Ensure converted resource passes target version's schema validation
6. **Response Generation:** Return ConversionResponse with converted resource or error details

Backward Compatibility Guarantees:

Compatibility Type	Guarantee Level	Description
Field Addition	Full	New optional fields can be added without breaking existing clients
Field Deprecation	Gradual	Deprecated fields continue working with warning messages for at least one version
Field Removal	Breaking	Removed fields cause validation errors and require client updates
Field Type Change	Breaking	Changing field types requires new version and conversion logic
Default Value Change	Potentially Breaking	May affect resources that rely on previous defaults

Storage Version Migration Process:

When changing the storage version from v1beta1 to v1, existing resources in etcd must be migrated to the new format. This migration happens lazily—resources are converted when they're next read or updated—or can be triggered explicitly using the storage version migrator tool.

Migration State Tracking:

Resource State	Description	Action Required
Native Storage	Resource stored in current storage version	No action needed
Conversion Pending	Resource stored in old version, needs conversion on access	Automatic conversion on read/write
Conversion Failed	Resource cannot be converted due to data incompatibility	Manual intervention required
Schema Violation	Converted resource fails new version validation	Manual correction needed

⚠ Pitfall: Breaking Conversion Assumptions

A common mistake is assuming that conversion is always reversible. When converting from a newer version to an older version, information loss may occur if the newer version contains fields that don't exist in the older version. This information loss means that converting v1 → v1alpha1 → v1 might not produce the same result as the original v1 resource.

Solution: Design conversion functions to preserve as much information as possible, use annotations to store unconvertible data, and test round-trip conversions to verify data integrity.

⚠ Pitfall: Storage Version Migration Timing

Changing the storage version in a CRD immediately affects how new resources are stored, but existing resources remain in their original format until converted. This can create inconsistency where identical resources have different storage representations.

Solution: Plan storage version migration carefully, use gradual rollout with monitoring, and consider running explicit migration jobs to convert existing resources proactively rather than waiting for lazy conversion.

Version Deprecation Process:

- Deprecation Announcement:** Mark version as deprecated in CRD status and release notes
- Warning Generation:** API server returns warning headers when deprecated versions are used
- Grace Period:** Maintain deprecated version for at least two minor releases
- Usage Monitoring:** Track deprecated version usage through metrics and audit logs
- Removal Planning:** Coordinate with users to migrate before version removal
- Final Removal:** Remove deprecated version from CRD and update conversion webhooks

This structured deprecation process ensures users have adequate time to migrate while providing clear signals about upcoming changes. The warning system allows automated detection of deprecated version usage in continuous integration pipelines and monitoring systems.

Implementation Guidance

A. Technology Recommendations:

Component	Simple Option	Advanced Option
Schema Validation	OpenAPI v3 in CRD manifest	CEL expressions + admission webhooks
Defaulting	Static defaults in CRD schema	Dynamic defaults via mutating webhook
Versioning	Single stable version	Multiple versions with conversion webhooks
Testing	kubectl dry-run validation	CRD validation test suite with controller-runtime

B. Recommended File Structure:

```
operator-project/
├── config/
│   ├── crd/
│   │   ├── bases/
│   │   │   └── example.com_databases.yaml      ← Generated CRD manifest
│   │   └── kustomization.yaml
│   └── samples/
│       └── example_v1_database.yaml          ← Example custom resource
└── api/
    └── v1/
        ├── database_types.go                ← Database struct definitions
        ├── database_webhook.go              ← Admission webhook logic
        ├── groupversion_info.go            ← API group registration
        └── zz_generated.deepcopy.go        ← Generated deep copy methods
└── internal/
    └── controller/
        └── database_controller.go        ← Controller reconciliation logic
```

C. Infrastructure Starter Code:

Complete CRD manifest generation configuration:

```
// api/v1/database_types.go
GO

package v1

import (
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/api/resource"
    corev1 "k8s.io/api/core/v1"
)

// DatabaseSpec defines the desired state of Database

type DatabaseSpec struct {

    // +kubebuilder:validation:Minimum=1
    // +kubebuilder:validation:Maximum=10
    // +kubebuilder:default=1

    Replicas int32 `json:"replicas"`

    // +kubebuilder:validation:Pattern="^[0-9]+\\.[0-9]+\\.[0-9]+$"
    // +kubebuilder:default="5.7.0"

    Version string `json:"version"`

    // +kubebuilder:validation:Pattern="^[0-9]+[KMGT]i?$"
    // +kubebuilder:default="10Gi"

    StorageSize string `json:"storageSize"`

    // +kubebuilder:validation:Pattern="^(@{annually|yearly|monthly|weekly|daily|hourly|reboot})|(@every (\\\d+(ns|us|ms|s|m|h))+)|((\\d+|\\*)|((\\d+)? (\\d+|\\*))(\\d+)? (\\d+|\\*))(\\d+)? ((\\d+|\\*))(\\d+)? ((\\d+|\\*))(\\d+)? ((\\d+|\\*))(\\d+)?)$"
    // +optional

    BackupSchedule *string `json:"backupSchedule,omitempty"`

    // +optional

    Resources *corev1.ResourceRequirements `json:"resources,omitempty"`

    // +kubebuilder:validation:MaxProperties=20
    // +optional

    Config map[string]string `json:"config,omitempty"`
}

// DatabaseStatus defines the observed state of Database
```

```

type DatabaseStatus struct {

    // +optional

    Conditions []metav1.Condition `json:"conditions,omitempty"`

    // +optional

    ReadyReplicas int32 `json:"readyReplicas,omitempty"`

    // +optional

    ObservedGeneration int64 `json:"observedGeneration,omitempty"`

    // +kubebuilder:validation:Enum=Pending;Running;Failed;Terminating

    // +optional

    Phase string `json:"phase,omitempty"`

    // +optional

    LastBackupTime *metav1.Time `json:"lastBackupTime,omitempty"`

    // +optional

    Endpoints []string `json:"endpoints,omitempty"`

    // +optional

    Message string `json:"message,omitempty"`

}

//+kubebuilder:object:root=true

//+kubebuilder:subresource:status

//+kubebuilder:resource:shortName=db

//+kubebuilder:printcolumn:name="Replicas",type=integer,JSONPath=`.spec.replicas`

//+kubebuilder:printcolumn:name="Ready",type=integer,JSONPath=`.status.readyReplicas`

//+kubebuilder:printcolumn:name="Version",type=string,JSONPath=`.spec.version`

//+kubebuilder:printcolumn:name="Phase",type=string,JSONPath=`.status.phase`

//+kubebuilder:printcolumn:name="Age",type=date,JSONPath=`.metadata.creationTimestamp`


// Database is the Schema for the databases API

type Database struct {

    metav1.TypeMeta `json:",inline"`

    metav1.ObjectMeta `json:"metadata,omitempty"`
}

```

```
Spec  DatabaseSpec `json:"spec,omitempty"`

Status DatabaseStatus `json:"status,omitempty"`

}

//+kubebuilder:object:root=true

// DatabaseList contains a list of Database

type DatabaseList struct {

    metav1.TypeMeta `json:",inline"`

    metav1.ListMeta `json:"metadata,omitempty"`

    Items          []Database `json:"items"`

}

func init() {

    SchemeBuilder.Register(&Database{}, &DatabaseList{})

}
```

D. Core Logic Skeleton Code:

Admission webhook implementation template:

GO

```

// api/v1/database_webhook.go

package v1

import (
    "context"
    "fmt"
    ap ierrors "k8s.io/apimachinery/pkg/api/errors"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/apimachinery/pkg/runtime/schema"
    "k8s.io/apimachinery/pkg/util/validation/field"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/webhook"
)

// SetupWebhookWithManager registers the webhook with the manager

func (r *Database) SetupWebhookWithManager(mgr ctrl.Manager) error {
    return ctrl.NewWebhookManagedBy(mgr).
        For(r).
        Complete()
}

//+kubebuilder:webhook:path=/mutate-example-com-v1-
//+kubebuilder:webhook:path=/validate-example-com-v1-
database,mutating=true,failurePolicy=fail,sideEffects=None,groups=example.com,resources=databases,verbs=create;update,versions=v1,name=mdatabase.kb.io,admissionReviewVersions=v1

var _ webhook.Defaulter = &Database{}


// Default implements webhook.Defaulter so a webhook will be registered for the type

func (r *Database) Default() {
    // TODO 1: Set default replica count if not specified
    // TODO 2: Set default version if not specified
    // TODO 3: Set default storage size if not specified
    // TODO 4: Set default resource requests if resources specified but requests missing
    // Hint: Use pointer fields to detect when values are not set vs set to zero
}

//+kubebuilder:webhook:path=/mutate-example-com-v1-
//+kubebuilder:webhook:path=/validate-example-com-v1-
database,mutating=false,failurePolicy=fail,sideEffects=None,groups=example.com,resources=databases,verbs=create;update,versions=v1,name=vdatabase.kb.io,admissionReviewVersions=v1

```

```
var _ webhook.Validator = &Database{}
```

// ValidateCreate implements webhook.Validator so a webhook will be registered for the type

```
func (r *Database) ValidateCreate() error {
```

// TODO 1: Validate replica count is reasonable for storage size

// TODO 2: Validate version format and ensure it's a supported version

// TODO 3: Validate backup schedule cron expression if provided

// TODO 4: Validate resource requests are less than or equal to limits

// TODO 5: Validate configuration parameters are valid for the specified version

```
    return r.validateDatabase()
```

```
}
```

// ValidateUpdate implements webhook.Validator so a webhook will be registered for the type

```
func (r *Database) ValidateUpdate(old runtime.Object) error {
```

```
    oldDB := old.(*Database)
```

// TODO 1: Prevent downgrading to older versions

// TODO 2: Prevent reducing storage size (data loss risk)

// TODO 3: Validate replica count changes are within safe limits

// TODO 4: Ensure configuration changes are compatible with current version

```
    return r.validateDatabase()
```

```
}
```

// ValidateDelete implements webhook.Validator so a webhook will be registered for the type

```
func (r *Database) ValidateDelete() error {
```

// TODO 1: Check if database has active connections (would require external check)

// TODO 2: Verify backup is recent if backup schedule is configured

// TODO 3: Implement additional deletion safety checks as needed

```
    return nil
```

```
}
```

```
func (r *Database) validateDatabase() error {
```

```
    var allErrs field.ErrorList
```

// TODO: Add field-specific validation logic here

```

// Use field.ErrorList to collect multiple validation errors

// Return apierrors.NewInvalid() with collected errors

if len(allErrs) == 0 {

    return nil
}

return apierrors.NewInvalid(
    schema.GroupKind{Group: "example.com", Kind: "Database"},

    r.Name,
    allErrs,
)
}

```

E. Language-Specific Hints:

- Use `kubebuilder` markers (`//+kubebuilder:validation:`) for generating OpenAPI schema validation
- Use `controller-gen` to generate CRD manifests from Go struct definitions
- Use `k8s.io/apimachinery/pkg/util/validation/field` for structured validation error reporting
- Use pointer fields (`*string`, `*int32`) to distinguish between unset and zero values
- Use `omitempty` JSON tags for optional fields to keep manifests clean
- Use `sigs.k8s.io/controller-runtime/pkg/webhook` for admission webhook infrastructure

F. Milestone Checkpoint:

After implementing the CRD and webhooks, verify the following behavior:

```
# Generate and apply CRD
make manifests

kubectl apply -f config/crd/bases/

# Test schema validation

kubectl apply -f - <<EOF
apiVersion: example.com/v1
kind: Database
metadata:
  name: test-invalid
spec:
  replicas: 0 # Should fail validation (minimum 1)
  version: "invalid" # Should fail pattern validation
  storageSize: "bad-format" # Should fail pattern validation
EOF

# Test successful creation

kubectl apply -f - <<EOF
apiVersion: example.com/v1
kind: Database
metadata:
  name: test-valid
spec:
  replicas: 3
  version: "5.7.0"
  storageSize: "20Gi"
EOF

# Verify printer columns work

kubectl get databases
```

BASH

Expected output should show validation errors for invalid resource and successful creation with custom columns for valid resource.

G. Debugging Tips:

Symptom	Likely Cause	How to Diagnose	Fix
CRD not found error	CRD manifest not applied or malformed	<code>kubectl get crd databases.example.com -o yaml</code>	Apply CRD manifest with <code>kubectl apply -f</code>
Schema validation not working	Kubebuilder markers incorrect or CRD not regenerated	Check CRD openAPIV3Schema section	Run <code>make manifests</code> and reapply CRD
Webhook not called	Webhook configuration missing or service unreachable	<code>kubectl get validatingwebhookconfiguration,mutatingwebhookconfiguration</code>	Verify webhook configuration and service
Certificate errors	TLS certificates expired or malformed	Check webhook pod logs for certificate errors	Regenerate certificates or check cert-manager
Status updates fail	Status subresource not enabled in CRD	Check CRD for <code>subresources: status:</code>	Add status subresource to CRD spec

Controller and Reconciliation Engine

Milestone(s): Milestone 2 (Controller Setup), Milestone 3 (Reconciliation Loop) - implements the core controller logic that watches custom resources and maintains desired state through continuous reconciliation

The controller is the beating heart of any Kubernetes operator. While custom resources define the API and webhooks validate incoming requests, the controller provides the intelligence that continuously monitors the cluster and takes action to maintain the desired state. This section details how to build a robust controller that can handle the complexities of distributed systems while integrating seamlessly with Kubernetes' declarative model.

Mental Model: The Control System

Think of the Kubernetes controller as an **automatic thermostat for your infrastructure**. Just as a thermostat continuously measures the current temperature, compares it to the desired temperature setting, and adjusts the heating or cooling system to eliminate the difference, a Kubernetes controller continuously observes the actual state of resources, compares it to the desired state declared in custom resources, and takes corrective actions to eliminate any drift.

This analogy reveals several key insights about controller design. First, the controller operates on a **feedback loop** - it measures, compares, acts, and then measures again. The system is inherently reactive rather than predictive. Second, the controller must be **level-triggered rather than edge-triggered** - it doesn't just respond to change events, but continuously ensures the system remains in the correct state even after transient failures or external modifications. Third, like a good thermostat, the controller should be **patient and persistent** - it doesn't panic when the temperature is wrong, but steadily works to correct it over time.

The control theory perspective also highlights the importance of **convergence and stability**. A poorly designed thermostat might oscillate wildly, turning the heat on and off rapidly without ever settling at the target temperature. Similarly, a poorly designed controller might create resources, delete them, recreate them, and never reach a stable state. The reconciliation algorithm must be designed to converge toward the desired state and remain stable once reached.

Unlike a simple thermostat, however, a Kubernetes controller operates in a **distributed system with partial failures**. The "temperature sensor" (informer cache) might be stale, the "heating system" (Kubernetes API) might be temporarily unavailable, or multiple "thermostats" (controller replicas) might be running simultaneously. The controller must be designed to handle these complexities gracefully while still maintaining the core feedback loop behavior.

Informer and Caching Architecture

The informer architecture is the foundation that makes efficient Kubernetes controllers possible. Without informers, controllers would need to constantly poll the API server to check for changes, creating excessive load and introducing latency. The informer pattern solves this by maintaining a **local, eventually consistent cache** of resources that the controller cares about.

Think of the informer as a **dedicated research assistant** who maintains a constantly updated filing system of all relevant documents. Instead of the controller running to the library (API server) every time it needs to check something, the research assistant watches for new publications, updates the local files, and notifies the controller when something important changes. This allows the controller to work efficiently from local information while still staying current with the authoritative source.

The informer architecture consists of several key components working together. The **SharedIndexInformer** is the core component that watches a specific resource type and maintains a local cache. It uses the Kubernetes **list-and-watch** protocol to efficiently stay synchronized with the API server. The informer first performs a list operation to get the current state of all resources, then establishes a watch stream to receive incremental updates.

Informer Component	Purpose	Key Responsibilities
SharedIndexInformer	Core caching and watching	Maintains local cache, watches for changes, triggers event handlers
Reflector	API server communication	Implements list-and-watch protocol, handles connection failures
DeltaFIFO	Event queuing	Queues resource changes, deduplicates events, provides ordering
Store	Local cache storage	Indexes resources for efficient lookup, thread-safe access
ResourceEventHandler	Change notification	Receives add/update/delete events, enqueues work items

The **SharedIndexInformer** uses several sophisticated mechanisms to provide reliable caching. The informer maintains multiple indexes on the cached data, allowing efficient lookups by different criteria. For example, it can quickly find all resources in a specific namespace or all resources owned by a particular parent object. The informer also handles **resync periods** - even without changes, it periodically re-examines all cached resources to ensure the controller processes any items that might have been missed due to transient failures.

Critical Insight: The informer cache is eventually consistent, not strongly consistent. There's always a window where the local cache might not reflect the latest API server state. Controllers must be designed to handle this gracefully by making their operations idempotent and tolerating stale information.

Event handling in the informer architecture follows a specific pattern designed for reliability and performance. When the informer detects a change, it doesn't immediately call the controller's reconciliation logic. Instead, it calls a **ResourceEventHandler** that typically just enqueues a work item into a rate-limited work queue. This decoupling is essential for several reasons: it prevents slow reconciliation logic from blocking the informer's watch stream, it allows for batching and deduplication of rapid changes to the same resource, and it enables sophisticated retry and backoff strategies.

Event Handler Method	When Called	Typical Action
OnAdd	Resource created or informer started	Enqueue resource key for reconciliation
OnUpdate	Resource spec or status changed	Compare old vs new, enqueue if meaningful change
onDelete	Resource deleted	Enqueue resource key for cleanup reconciliation

The work queue used with informers provides several critical capabilities for robust controller operation. The **RateLimitingInterface** prevents controllers from overwhelming the API server with rapid requests when resources are changing frequently. The queue implements **exponential backoff** - if reconciling a resource fails, the queue will wait increasingly longer periods before retrying the same resource. The queue also provides **deduplication** - if the same resource changes multiple times before the controller processes it, only one reconciliation is triggered.

Decision: SharedIndexInformer vs Direct API Calls

- **Context:** Controllers need to read resource state frequently during reconciliation. Direct API calls would create excessive load and latency.
- **Options Considered:** Direct API calls for each read, local caching without watches, SharedIndexInformer pattern
- **Decision:** Use SharedIndexInformer for all resource reads
- **Rationale:** Provides efficient local caching, automatic synchronization, built-in indexing, and integrates with controller-runtime patterns. Reduces API server load by orders of magnitude while providing sub-second change notification.
- **Consequences:** Controller logic must handle eventually consistent data and implement idempotent operations. Memory usage increases with number of cached resources.

The informer lifecycle requires careful management to ensure proper startup and shutdown behavior. During controller startup, the informer must complete its initial **cache synchronization** before the controller begins processing work items. This ensures that the local cache contains a complete view of existing resources before reconciliation begins. The `cache.WaitForCacheSync` function provides this capability, blocking until the informer has completed its initial list operation and is ready to serve requests.

Informer Lifecycle Phase	Actions Required	Common Pitfalls
Startup	Start informer, wait for cache sync	Processing events before sync complete
Running	Process events, handle watch errors	Blocking event handlers with slow operations
Shutdown	Stop informer, drain work queue	Goroutine leaks from improper cleanup

Reconciliation Algorithm

The reconciliation algorithm is where the controller's intelligence lives. This algorithm implements the core feedback loop that compares desired state (from the custom resource spec) with actual state (from the cluster) and takes corrective actions to eliminate any differences. The algorithm must be **idempotent**, **convergent**, and **robust** in the face of partial failures and concurrent modifications.

The reconciliation process follows a well-defined sequence of steps that ensures consistent and reliable behavior. Each step builds upon the previous one, creating a logical progression from observation through decision-making to action. Understanding this sequence is crucial for implementing correct controller behavior.

Step 1: Resource Retrieval and Validation

The reconciliation begins by retrieving the current state of the custom resource that triggered the reconciliation. This involves fetching the resource from the local informer cache (not directly from the API server) and performing basic validation to ensure the resource is in a state that can be reconciled.

During this phase, the controller must handle several edge cases. The resource might have been deleted between the time it was enqueued and when reconciliation begins. The resource might be in the middle of being deleted (having a non-zero deletion timestamp) but still have finalizers that prevent actual removal. The resource might have been modified by another client, changing its generation number since it was last processed.

Step 2: Current State Discovery

Once the target resource is validated, the controller must discover the current state of all resources that should exist to satisfy the desired configuration. This typically involves querying for Deployments, Services, ConfigMaps, and other Kubernetes resources that the operator manages on behalf of the custom resource.

This discovery phase uses **owner references** to identify managed resources. Each resource created by the controller should have an owner reference pointing back to the custom resource that caused its creation. This enables both discovery (finding all owned resources) and automatic garbage collection (when the owner is deleted, Kubernetes automatically removes owned resources).

Resource Discovery Pattern	Query Method	Purpose
List with label selectors	<code>client.List(ctx, &list, client.MatchingLabels{...})</code>	Find resources with specific labels
List with owner references	<code>client.List(ctx, &list, client.MatchingFields{".metadata.ownerReferences.uid": owner.UID})</code>	Find all resources owned by custom resource
Direct name lookup	<code>client.Get(ctx, namespacedName, &resource)</code>	Check existence of expected resource

Step 3: Desired State Computation

With the current state known, the controller must compute what the desired state should look like based on the custom resource's spec. This computation transforms the high-level desired configuration into specific Kubernetes resources with complete specifications.

This step often involves complex logic specific to the application being managed. For a database operator, it might involve determining the appropriate container images, computing resource requirements, generating configuration files, and deciding on networking configuration. The key principle is that this computation should be **deterministic** - given the same input spec, it should always produce the same desired resource configurations.

Step 4: Difference Analysis

The heart of the reconciliation algorithm is comparing the current state with the desired state to identify what changes need to be made. This comparison must be sophisticated enough to distinguish between meaningful differences that require action and insignificant differences that should be ignored.

Some fields in Kubernetes resources are modified by the system after creation and should not be considered during difference analysis. For example, the `resourceVersion` field changes with every update, and many fields have default values set by admission controllers. The controller must know which fields are significant for its purposes and focus only on those during comparison.

Difference Type	Action Required	Example
Missing resource	Create with desired spec	Deployment doesn't exist but should
Extra resource	Delete unwanted resource	Service exists but shouldn't based on current spec
Configuration drift	Update to match desired spec	Container image version differs from spec
Status-only changes	No action (status-only update)	Ready replicas count changed

Step 5: Corrective Actions

Based on the difference analysis, the controller takes corrective actions to move the actual state toward the desired state. These actions typically involve creating, updating, or deleting Kubernetes resources through the client library.

The order of operations matters significantly during this phase. Resources with dependencies should be created in the correct order (ConfigMaps before Deployments that reference them), and deletions should happen in reverse dependency order. The controller should also be prepared for operations to fail partially and leave the system in an intermediate state.

Step 6: Status Update

After attempting corrective actions, the controller updates the custom resource's status to reflect the current state of reconciliation. This provides visibility to users about what the controller has accomplished and whether any problems were encountered.

The status update should include several types of information: overall conditions indicating the health and progress of reconciliation, specific metrics like the number of ready replicas, and detailed messages about any errors encountered. This information is crucial for debugging and monitoring operator behavior.

Status Field	Purpose	Example Values
Conditions	Overall health and progress	Ready=True, Progressing=False
ReadyReplicas	Current operational capacity	3 (out of 3 desired replicas)
ObservedGeneration	Last spec version processed	matches metadata.generation when up-to-date
Phase	High-level lifecycle state	Pending, Ready, Failed
Message	Human-readable status description	"All replicas are ready and serving traffic"

The complete reconciliation algorithm can be expressed as the following decision flow:

1. Retrieve target custom resource from cache

- If not found, log and return (resource was deleted)
- If deletion timestamp set, proceed with cleanup logic
- If generation unchanged and conditions current, return early (optimization)

2. Discovery current owned resources

- Query for all resources with owner references to this custom resource
- Query for resources matching expected names and labels
- Build map of current resource state keyed by resource type and name

3. Compute desired resource specifications

- Transform custom resource spec into concrete Kubernetes resource specs
- Apply naming conventions, labels, owner references
- Generate configurations, secrets, and other derived resources

4. Perform three-way diff analysis

- Compare desired vs current for each resource type
- Identify resources to create, update, or delete
- Filter out insignificant differences (system-managed fields)

5. Execute corrective actions in dependency order

- Create missing dependencies first (ConfigMaps, Secrets)
- Create or update primary resources (Deployments, Services)
- Delete obsolete resources last
- Handle partial failures gracefully

6. Update custom resource status

- Set conditions based on reconciliation success/failure
- Update metrics (ready replicas, observed generation)
- Record detailed error messages for debugging
- Use status subresource to avoid update conflicts

Decision: Reconciliation Frequency Strategy

- **Context:** Controllers need to balance responsiveness with efficiency, ensuring timely updates without overwhelming the API server
- **Options Considered:** Continuous polling, event-driven only, hybrid approach with periodic resync
- **Decision:** Event-driven reconciliation with configurable resync period (default 10 hours)
- **Rationale:** Event-driven provides immediate response to changes while periodic resync catches missed events and handles external modifications. Long resync period minimizes unnecessary work while providing eventual consistency guarantees.
- **Consequences:** Enables efficient resource usage while maintaining strong consistency. Requires careful event handler implementation to avoid missing important changes.

Error Handling and Requeue Strategy

Robust error handling is what separates production-ready controllers from toy implementations. In distributed systems, failures are not exceptional - they are expected and must be handled gracefully. A well-designed controller can distinguish between different types of failures and respond appropriately to each category.

The fundamental insight for controller error handling is that **most failures are transient** and will resolve themselves if the controller simply waits and retries. Network connectivity issues, temporary API server overload, and resource conflicts during concurrent updates all fall into this category. However, some failures are **permanent** and will never succeed no matter how many times they are retried. Attempting to create a resource with an invalid name or referencing a non-existent secret will fail consistently until the underlying problem is fixed.

Transient vs Permanent Error Classification

The controller's error handling strategy begins with correctly classifying errors as transient or permanent. This classification determines whether the controller should retry the operation automatically or surface the error to users for manual intervention.

Error Category	Examples	Handling Strategy
Transient Network	Connection timeout, DNS resolution failure	Exponential backoff retry
Transient API Server	Rate limiting, temporary overload	Exponential backoff with jitter
Transient Resource	Update conflicts, temporary resource constraints	Linear backoff retry
Permanent Configuration	Invalid resource names, malformed specs	Surface error, no retry
Permanent Permission	RBAC denials, forbidden operations	Surface error, log for admin
Permanent Dependency	Non-existent secrets, missing CRDs	Surface error, watch for dependency

Transient errors should trigger **automatic retry with backoff**. The controller enqueues the resource for reprocessing after a delay, with the delay increasing on subsequent failures. This prevents the controller from overwhelming struggling components while still providing timely retry when conditions improve.

Permanent errors should **not trigger automatic retry**. Instead, the controller should update the resource status with a detailed error message, emit a warning event, and wait for external action to resolve the underlying problem. Continuing to retry permanent failures wastes resources and creates excessive log noise.

Exponential Backoff Implementation

The work queue provides sophisticated retry mechanisms that implement exponential backoff automatically. When the controller returns an error from the `Reconcile` function, the controller-runtime framework interprets this as a request to retry the operation after a delay.

The backoff algorithm starts with a base delay (typically 1 second) and doubles the delay after each failure, up to a maximum delay (typically 16 minutes). This exponential growth ensures that transient issues get resolved quickly while preventing permanent failures from consuming excessive resources.

Failure Count	Retry Delay	Cumulative Time	Rationale
1	1 second	1 second	Quick recovery for brief glitches
2	2 seconds	3 seconds	Still responsive for short outages
3	4 seconds	7 seconds	Moderate backing off
4	8 seconds	15 seconds	Reduced load during longer issues
5	16 seconds	31 seconds	Significant backoff
10	16 minutes	~17 minutes	Maximum delay prevents resource waste

Jitter is added to the backoff delays to prevent **thundering herd** problems when many controllers are retrying simultaneously. Without jitter, if the API server becomes unavailable and then recovers, all controllers would retry at exactly the same time, potentially overwhelming the server again. Jitter adds randomness to spread out the retry attempts.

Requeue Strategies

The controller-runtime framework provides several mechanisms for requesting that a resource be reprocessed:

Requeue Method	When to Use	Behavior
<code>return ctrl.Result{}, err</code>	Transient error occurred	Exponential backoff based on error count
<code>return ctrl.Result{RequeueAfter: duration}, nil</code>	Need to check again after specific time	Fixed delay, resets error count
<code>return ctrl.Result{Requeue: true}, nil</code>	Need immediate reprocessing	Immediate enqueue, resets error count
<code>return ctrl.Result{}, nil</code>	Reconciliation successful	No requeue, wait for next event

The choice between these requeue strategies depends on the specific situation encountered during reconciliation. Transient errors should use the error return to trigger exponential backoff. Situations where the controller needs to wait for an external condition (like waiting for a deployment to become ready) should use `RequeueAfter` with an appropriate delay. Immediate requeue should be used sparingly and only when the controller knows that conditions have changed and reconciliation should succeed.

Reconciliation Result Patterns

Different reconciliation outcomes require different result patterns to ensure optimal controller behavior:

```

// Successful reconciliation, no further action needed
return ctrl.Result{}, nil

// Transient error, use exponential backoff
if isTransientError(err) {
    return ctrl.Result{}, err
}

// Permanent error, update status and don't retry
r.updateErrorStatus(ctx, database, err)
return ctrl.Result{}, nil

// Waiting for external condition, check again in 30 seconds
if !allResourcesReady {
    return ctrl.Result{RequeueAfter: REQUEUE_AFTER_DURATION}, nil
}

// Configuration changed during reconciliation, reprocess immediately
if resourceModified {
    return ctrl.Result{Requeue: true}, nil
}

```

Status Condition Management

Error handling integrates closely with status condition management. The controller should maintain conditions that reflect the current state of reconciliation, including any errors encountered. Well-designed conditions provide both programmatic and human-readable information about controller state.

Condition Type	Status	Reason	Meaning
Ready	True	AllResourcesReady	Reconciliation successful, all resources operational
Ready	False	ReconciliationInProgress	Reconciliation in progress, resources not yet ready
Ready	False	ReconciliationFailed	Permanent error prevents successful reconciliation
Progressing	True	CreatingResources	Controller actively working to reach desired state
Progressing	False	ReconciledSuccessfully	Desired state achieved, no work in progress

Decision: Error Classification Strategy

- **Context:** Controllers must distinguish between errors that should trigger retry vs errors that require human intervention
- **Options Considered:** Retry all errors, never retry errors, classify by error type and HTTP status
- **Decision:** Classify errors by HTTP status code and error type, with configurable retry limits
- **Rationale:** HTTP 4xx errors (except 408, 429) are typically permanent while 5xx errors are transient. Specific error types (validation failures) can be classified regardless of HTTP status. Retry limits prevent infinite loops on misclassified errors.
- **Consequences:** Enables efficient resource usage and faster recovery from transient issues. Requires careful error classification logic and monitoring of retry patterns.

Common Pitfalls

⚠️ Pitfall: Processing Events Before Cache Sync

A common mistake is starting the reconciliation loop before the informer cache has completed its initial synchronization. This can lead to creating duplicate resources because the controller doesn't see existing resources in its empty cache. Always wait for `cache.WaitForCacheSync` to return true before starting worker goroutines.

⚠️ Pitfall: Blocking Event Handlers

Event handlers (`OnAdd`, `OnUpdate`, `OnDelete`) run in the informer's goroutine and must not block. Performing expensive operations or API calls directly in event handlers will block the informer's watch stream, preventing it from processing new events. Event handlers should only enqueue work items and return quickly.

Pitfall: Not Implementing Idempotent Operations

Controllers may reconcile the same resource multiple times due to retries, cache inconsistencies, or multiple events. All reconciliation operations must be idempotent - applying the same reconciliation multiple times should produce the same result. This requires checking existing resource state before creating new resources.

Pitfall: Ignoring Owner References

Without proper owner references, the controller cannot reliably discover resources it created, and garbage collection will not work correctly. Every resource created by the controller must have an owner reference pointing to the custom resource that caused its creation. Use `controllerutil.SetControllerReference` to set these correctly.

Pitfall: Status Update Conflicts

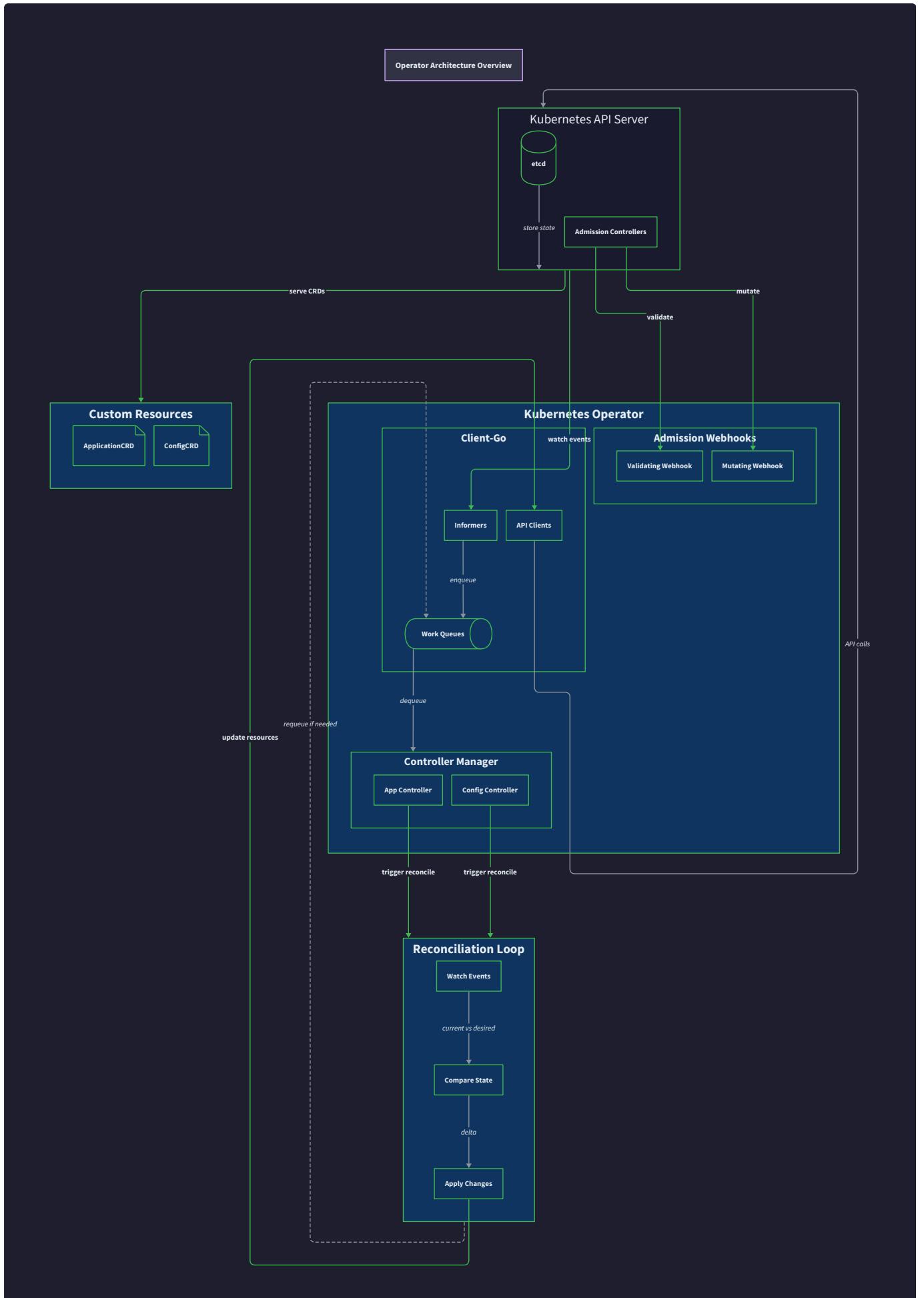
Updating both spec and status fields in the same operation can cause update conflicts with other clients. Use the status subresource (`client.Status().Update()`) to update status independently from spec changes. This allows the controller to update status without conflicting with user spec modifications.

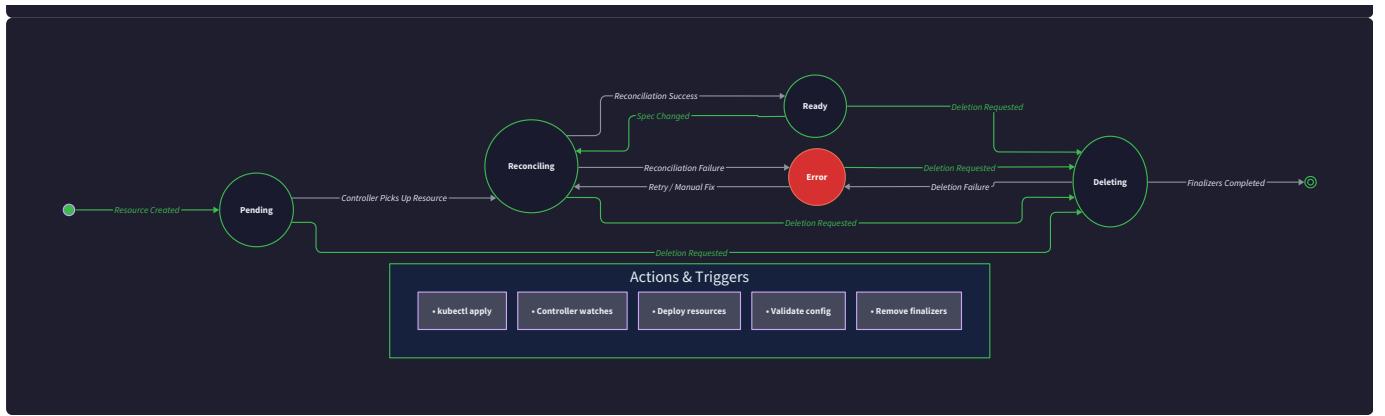
Pitfall: Infinite Reconciliation Loops

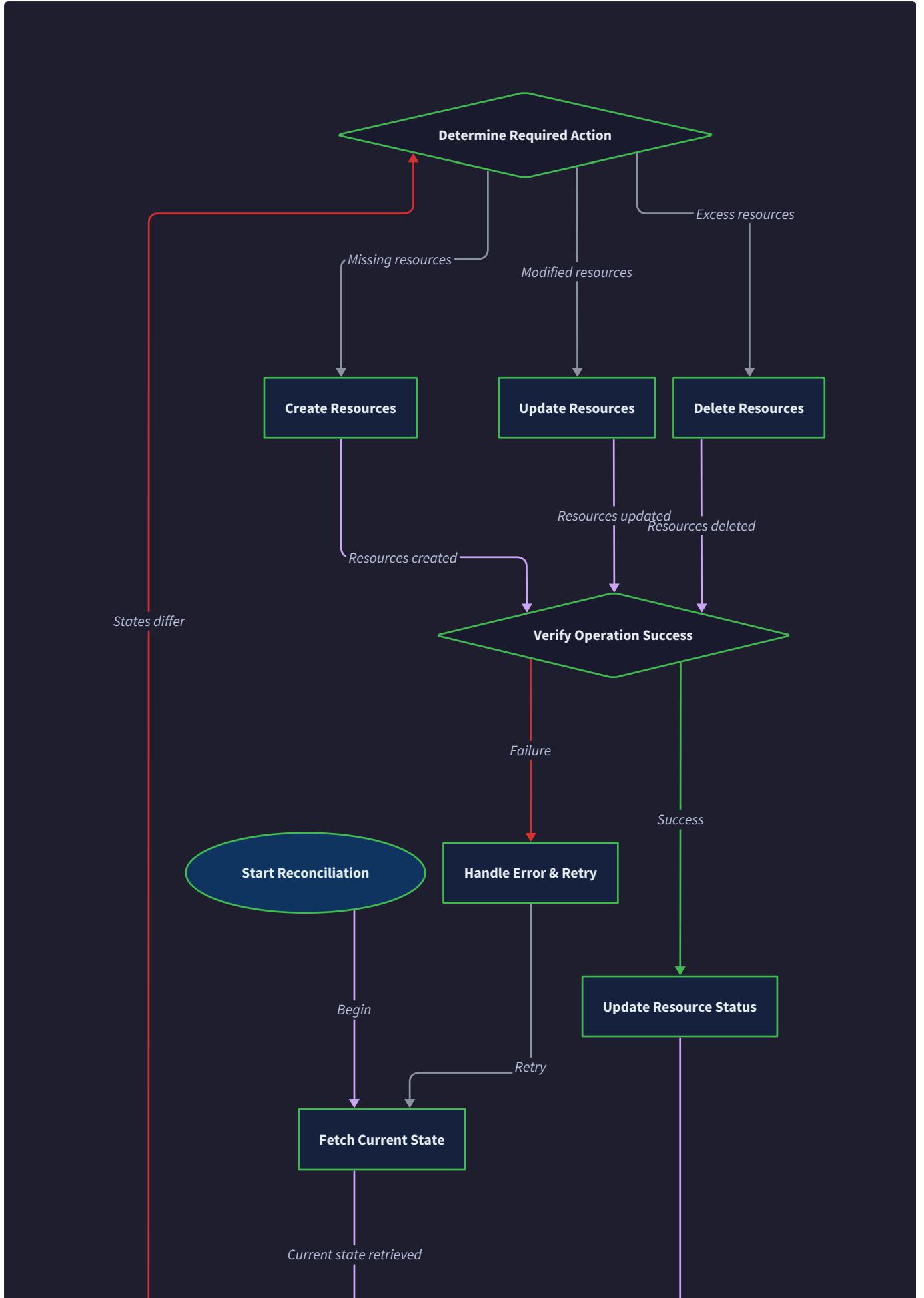
If the controller modifies resources in a way that triggers its own event handlers, it can create infinite reconciliation loops. Be careful when updating resources that the controller watches. Use proper condition checking to avoid unnecessary updates, and ensure that status updates don't trigger spec changes.

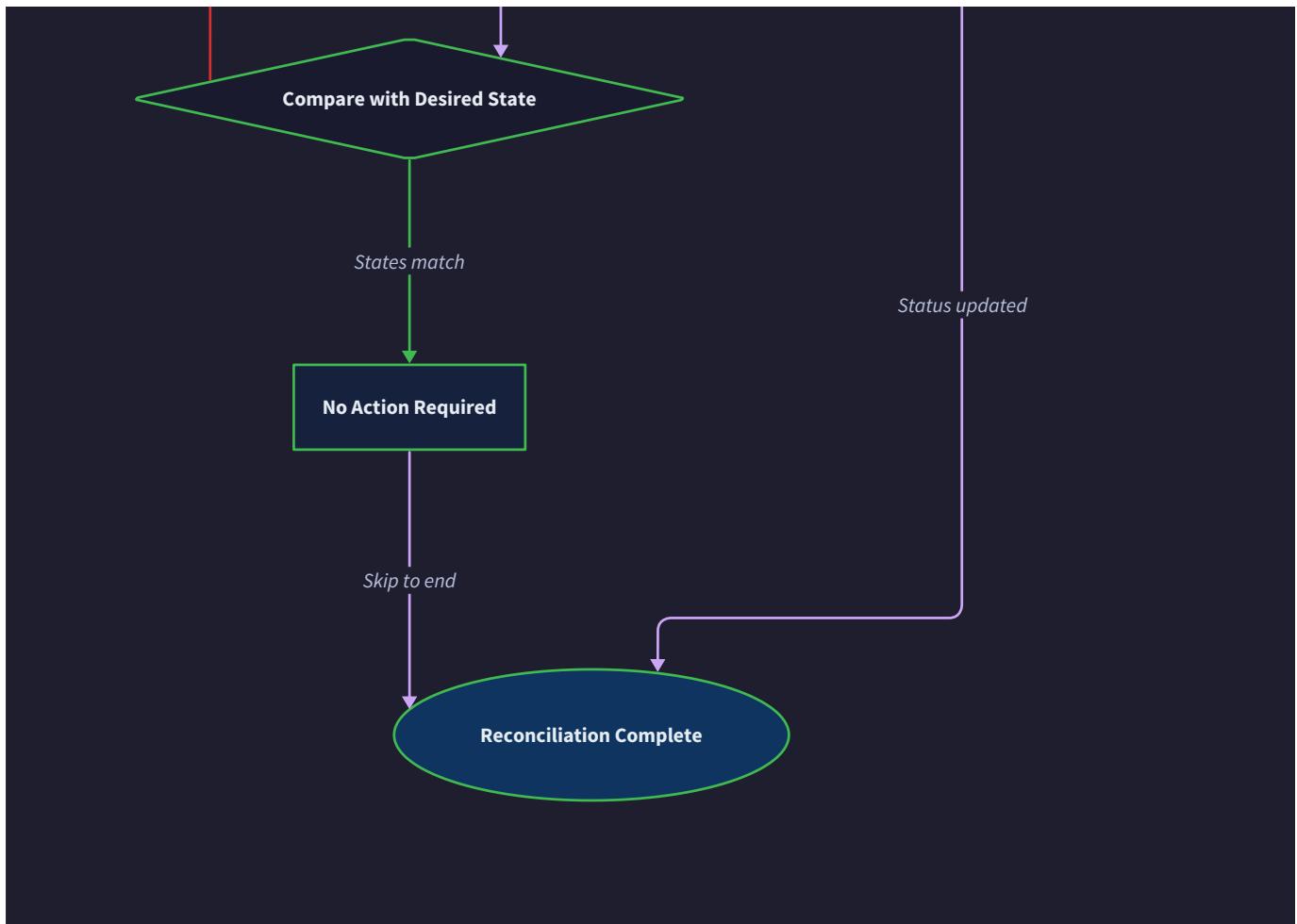
Pitfall: Not Handling Finalizers Properly

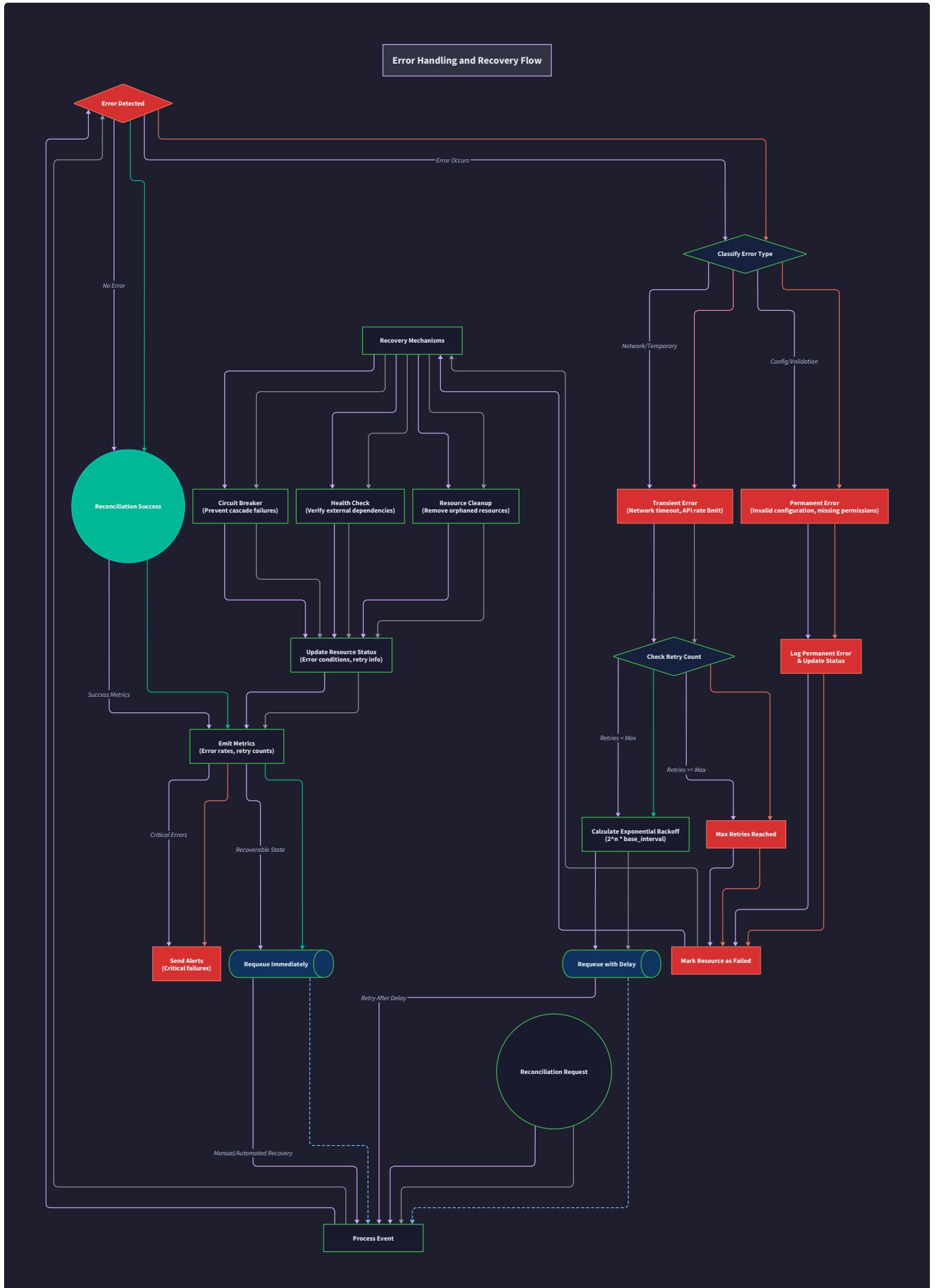
When using finalizers for cleanup, controllers must remove the finalizer after successful cleanup, or the resource will be stuck in deletion state forever. Always implement finalizer removal logic and test deletion scenarios thoroughly.











Implementation Guidance

This section provides concrete code patterns and infrastructure to implement the controller architecture described above. The focus is on providing complete, working infrastructure code while leaving the core business logic for you to implement.

Technology Recommendations

Component	Simple Option	Advanced Option
Controller Framework	controller-runtime (kubebuilder style)	Custom controller with client-go directly
Client Library	controller-runtime client (unified interface)	Separate typed and dynamic clients
Work Queue	controller-runtime managed queues	Manual workqueue.RateLimitingInterface
Informer Setup	controller-runtime managed informers	Manual shared informer factory
Testing Framework	controller-runtime envtest	Custom test cluster setup

Recommended File Structure

```
internal/controller/
    database_controller.go      ← main controller implementation
    database_controller_test.go ← unit tests with fake client
    suite_test.go               ← integration test suite setup
internal/resources/
    deployment.go              ← deployment resource builder
    service.go                 ← service resource builder
    configmap.go               ← configmap resource builder
config/
    samples/                   ← example custom resources
        database_v1_database.yaml
manager/
    manager.yaml               ← controller deployment
```

Infrastructure Starter Code

Complete Controller Structure ([internal/controller/database_controller.go](#)):

```
package controller

import (
    "context"
    "fmt"
    "time"

    appsv1 "k8s.io/api/apps/v1"
    corev1 "k8s.io/api/core/v1"
    "k8s.io/apimachinery/pkg/api/errors"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/apimachinery/pkg/types"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/controller/controllerutil"
    "sigs.k8s.io/controller-runtime/pkg/log"
    databasev1 "github.com/example/database-operator/api/v1"
)

const (
    REQUEUE_AFTER_DURATION = 30 * time.Second
    FINALIZER_NAME         = "database.example.com/finalizer"
)

// DatabaseReconciler reconciles a Database object

type DatabaseReconciler struct {
    client.Client
    Scheme *runtime.Scheme
}

//+kubebuilder:rbac:groups=database.example.com,resources=databases,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups=database.example.com,resources=databases/status,verbs=get;update;patch
//+kubebuilder:rbac:groups=database.example.com,resources=databases/finalizers,verbs=update
//+kubebuilder:rbac:groups=apps,resources=deployments,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups="",resources=services,verbs=get;list;watch;create;update;patch;delete
//+kubebuilder:rbac:groups="",resources=configmaps,verbs=get;list;watch;create;update;patch;delete
```

```
//+kubebuilder:rbac:groups="",resources=events,verbs=create;patch

// SetupWithManager sets up the controller with the Manager.

func (r *DatabaseReconciler) SetupWithManager(mgr ctrl.Manager) error {
    return ctrl.NewControllerManagedBy(mgr).
        For(&databasev1.Database{}).
        Owns(&appsv1.Deployment{}).
        Owns(&corev1.Service{}).
        Owns(&corev1.ConfigMap{}).
        Complete(r)
}
```

Work Queue and Event Handling Infrastructure:

GO

```
// ConditionType represents the type of condition

type ConditionType string

const (
    ConditionTypeReady      ConditionType = "Ready"
    ConditionTypeProgressing ConditionType = "Progressing"
)

// ConditionStatus represents the status of a condition

type ConditionStatus string

const (
    ConditionStatusTrue      ConditionStatus = "True"
    ConditionStatusFalse     ConditionStatus = "False"
    ConditionStatusUnknown   ConditionStatus = "Unknown"
)

// updateCondition updates or adds a condition to the Database status

func (r *DatabaseReconciler) updateCondition(database *databasev1.Database, conditionType ConditionType, status ConditionStatus, reason, message string) {

    condition := databasev1.Condition{

        Type:          string(conditionType),
        Status:        string(status),
        LastTransitionTime: metav1.Now(),
        Reason:        reason,
        Message:       message,
        ObservedGeneration: database.Generation,
    }

    // Find existing condition of same type

    for i, existing := range database.Status.Conditions {

        if existing.Type == condition.Type {

            // Only update if status changed

            if existing.Status != condition.Status {

                database.Status.Conditions[i] = condition
            } else {

                // Update message but keep transition time
            }
        }
    }
}
```

```
        database.Status.Conditions[i].Message = message

        database.Status.Conditions[i].Reason = reason

        database.Status.Conditions[i].ObservedGeneration = database.Generation

    }

    return
}

}

// Add new condition

database.Status.Conditions = append(database.Status.Conditions, condition)

}

// isErrorRetryable determines if an error should trigger exponential backoff retry

func isErrorRetryable(err error) bool {

    if errors.IsTimeout(err) || errors.IsServerTimeout(err) {

        return true

    }

    if errors.IsServiceUnavailable(err) || errors.IsInternalError(err) {

        return true

    }

    if errors.IsTooManyRequests(err) {

        return true

    }

    if errors.IsConflict(err) {

        return true // Resource conflicts are typically transient

    }

    // 4xx errors (except above) are typically permanent

    if errors.NotFound(err) || errors.Exists(err) {

        return false

    }

    if errors.BadRequest(err) || errors.Invalid(err) {

        return false

    }

    if errors.Forbidden(err) || errors.Unauthorized(err) {
```

```
    return false
}

// Default to retryable for unknown errors

return true
}
```

Core Logic Skeleton Code

Main Reconciliation Logic (implement this yourself):

GO

```
// Reconcile is part of the main kubernetes reconciliation loop which aims to
// move the current state of the cluster closer to the desired state.

func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
    log := log.FromContext(ctx)

    log.Info("Starting reconciliation", "database", req.NamespacedName)

    // TODO 1: Fetch the Database instance from the cache

    // Use r.Get(ctx, req.NamespacedName, &database) to retrieve the resource

    // If not found, return success (resource was deleted)

    // Hint: Use errors.NotFound(err) to check for not found errors

    // TODO 2: Handle deletion logic if deletion timestamp is set

    // Check if database.DeletionTimestamp != nil

    // If deleting and has finalizer, perform cleanup and remove finalizer

    // Return after handling deletion

    // Hint: Use controllerutil.ContainsFinalizer() and controllerutil.RemoveFinalizer()

    // TODO 3: Add finalizer if not present

    // Check if finalizer is present using controllerutil.ContainsFinalizer()

    // If missing, add it using controllerutil.AddFinalizer() and update resource

    // Hint: This ensures cleanup logic runs before deletion

    // TODO 4: Discover current state of all owned resources

    // Find existing Deployment, Service, ConfigMap using owner references

    // Build a map of current resources keyed by type and name

    // Hint: Use client.List with owner reference field selector or label selectors

    // TODO 5: Compute desired state based on Database spec

    // Call helper functions to generate desired Deployment, Service, ConfigMap specs

    // Apply naming conventions, labels, and owner references

    // Hint: Use controllerutil.SetControllerReference to set owner references

    // TODO 6: Perform three-way diff and apply changes

    // Compare current vs desired for each resource type

    // Create missing resources, update modified resources, delete obsolete resources

    // Handle partial failures gracefully - update what you can

    // Hint: Use r.createOrUpdateResource() helper for each resource
```

```
// TODO 7: Update Database status based on reconciliation results

// Set conditions (Ready, Progressing) based on success/failure

// Update ReadyReplicas from Deployment status

// Set ObservedGeneration to database.Generation when fully reconciled

// Use r.Status().Update() to persist status changes

// TODO 8: Determine appropriate return value

// Return error for transient failures (triggers exponential backoff)

// Return RequeueAfter for waiting on external conditions

// Return success (empty Result) when reconciliation complete

// Hint: Use isErrorRetryable() to classify errors

return ctrl.Result{}, nil

}
```

Resource Management Helper (implement this yourself):

```
// createOrUpdateResource creates or updates a Kubernetes resource to match desired state

func (r *DatabaseReconciler) createOrUpdateResource(ctx context.Context, desired client.Object, database
*databasev1.Database) error {

    log := log.FromContext(ctx)

    // TODO 1: Check if resource already exists

    // Create a new instance of the same type as desired

    // Use r.Get() to fetch existing resource

    // Handle both found and not found cases

    // Hint: Use errors.NotFound() to distinguish cases

    // TODO 2: If resource doesn't exist, create it

    // Set owner reference to the Database resource

    // Use r.Create() to create the resource

    // Log creation and return any errors

    // Hint: Use controllerutil.SetControllerReference()

    // TODO 3: If resource exists, check if update is needed

    // Compare significant fields between existing and desired

    // Skip fields that are modified by the system (resourceVersion, etc.)

    // Decide whether an update is necessary

    // Hint: Focus on spec fields, ignore status and metadata

    // TODO 4: Update resource if changes detected

    // Copy desired spec into existing resource

    // Preserve system-managed fields (resourceVersion, etc.)

    // Use r.Update() to persist changes

    // Log update and return any errors

    return nil
}
```

Status Management Helper (implement this yourself):

```

// updateDatabaseStatus updates the Database status subresource with current state

func (r *DatabaseReconciler) updateDatabaseStatus(ctx context.Context, database *databasev1.Database) error {

    log := log.FromContext(ctx)

    // TODO 1: Fetch current Deployment status

    // Look up the Deployment owned by this Database

    // Extract ready replicas, total replicas, and conditions

    // Handle case where Deployment doesn't exist yet

    // Hint: Use same naming convention as resource creation

    // TODO 2: Update Database status fields

    // Set ReadyReplicas from Deployment status

    // Set ObservedGeneration to database.Generation when fully processed

    // Determine overall Phase (Pending, Ready, Failed)

    // Update Endpoints if Service exists

    // TODO 3: Update status conditions

    // Set Ready condition based on whether all resources are operational

    // Set Progressing condition based on whether work is in progress

    // Include detailed messages for debugging

    // Hint: Use r.updateCondition() helper

    // TODO 4: Persist status changes

    // Use r.Status().Update() to update only status subresource

    // This avoids conflicts with spec changes

    // Return any update errors

    // Hint: Status updates use different RBAC permissions

    return r.Status().Update(ctx, database)
}

```

Language-Specific Hints

- Use `sigs.k8s.io/controller-runtime/pkg/client` for all API operations - it provides caching and works with both typed and untyped resources
- Use `controllerutil.SetControllerReference()` to establish owner relationships - this enables automatic garbage collection
- Use `errors.NotFound()`, `errors.Conflict()`, etc. for proper error classification rather than string matching
- Use `ctrl.Request.NamespacedName` type for resource keys - it's the standard way to identify resources
- Use `metav1.Now()` for condition timestamps to ensure consistent time format

- Use `runtime.Object` interface for generic resource handling in helper functions

Milestone Checkpoint

After implementing the controller logic:

Test Command: `go test ./internal/controller/... -v`

Expected Behavior:

- Controller starts and waits for cache sync before processing
- Creating a Database resource triggers reconciliation
- Owned Deployment, Service, and ConfigMap are created with proper owner references
- Database status is updated with Ready condition and replica counts
- Deleting Database triggers finalizer cleanup before actual deletion
- Status subresource updates don't conflict with spec changes

Manual Verification:

```
# Create test Database resource
kubectl apply -f config/samples/database_v1_database.yaml

# Verify owned resources were created
kubectl get deployments,services,configmaps -l app.kubernetes.io/name=database-sample

# Check Database status
kubectl get database database-sample -o yaml | grep -A 20 status:

# Verify owner references are set
kubectl get deployment database-sample-deployment -o yaml | grep -A 10 ownerReferences:
```

Signs Something is Wrong:

- "cache not started" errors → informer not properly initialized
- Resources created without owner references → garbage collection won't work
- Status updates causing "conflict" errors → not using status subresource
- Reconciliation not triggered by resource changes → event handlers not configured
- Controller processing before cache sync → missing WaitForCacheSync

Admission Control with Webhooks

Milestone(s): Milestone 4 (Webhooks) - implements admission webhooks for validation and mutation that intercept resource operations before they reach storage

Admission webhooks represent one of the most powerful extensibility mechanisms in Kubernetes, allowing operators to implement custom business logic that executes during the resource lifecycle. These webhooks act as gatekeepers that can inspect, validate, and modify resources before they are persisted to etcd, providing a critical control point for enforcing organizational policies and ensuring data consistency.

Mental Model: The Gatekeeper

Think of admission webhooks as security guards at the entrance to a high-security building. Just as every person entering the building must pass through a checkpoint where guards verify credentials, check for prohibited items, and potentially provide visitor badges, every resource entering the Kubernetes API must pass through admission control where webhooks can validate the request, reject invalid operations, and inject required metadata.

The analogy extends further: just as there might be multiple checkpoints with different responsibilities (one guard checks identification, another performs security screening, a third issues access badges), Kubernetes runs multiple webhook types in a specific order. **Mutating webhooks** act like the badge-issuing station - they can add required information to requests (like default values or labels). **Validating webhooks** act like the final security checkpoint - they perform comprehensive validation and can deny entry entirely if something is wrong.

The critical insight is that these gatekeepers operate synchronously in the request path. When a user runs `kubectl apply`, their request doesn't go directly to storage - it first passes through this entire admission chain. If any webhook rejects the request, the operation fails immediately and nothing gets stored. This provides strong consistency guarantees that manual validation processes cannot match.

Decision: Synchronous vs Asynchronous Webhook Processing

- **Context:** Admission webhooks could potentially operate asynchronously to avoid blocking the API request path, but this would require complex coordination mechanisms
- **Options Considered:** Synchronous processing during admission, asynchronous processing with eventual validation, hybrid approach with fast-path synchronous and slow-path asynchronous
- **Decision:** Synchronous processing during admission control
- **Rationale:** Provides immediate feedback to users, ensures invalid resources never reach storage, and maintains strong consistency without complex reconciliation logic for validation failures
- **Consequences:** Webhook latency directly impacts API response time, requires careful timeout configuration, but provides superior user experience and data integrity guarantees

Validating Webhook Implementation

Validating admission webhooks implement business rule validation that goes beyond what OpenAPI schema validation can express in the Custom Resource Definition. While CRD schemas excel at structural validation (ensuring fields have correct types, required fields are present), validating webhooks handle cross-field validation, external dependency checks, and complex business logic that requires computational evaluation.

The webhook receives an `AdmissionReview` request containing the proposed resource and must return an `AdmissionResponse` indicating whether to allow or deny the operation. The webhook has access to both the new resource (for create operations) and the old resource (for update operations), enabling sophisticated change validation that can prevent invalid state transitions.

AdmissionReview Request Structure:

Field	Type	Description
<code>kind</code>	string	Always "AdmissionReview"
<code>apiVersion</code>	string	API version of the admission review (admission.k8s.io/v1)
<code>request.uid</code>	string	Unique identifier for this admission request - must be included in response
<code>request.kind</code>	GroupVersionKind	The resource type being validated (e.g., Database)
<code>request.namespace</code>	string	Namespace of the resource being created/updated
<code>request.operation</code>	string	CREATE, UPDATE, or DELETE
<code>request.object</code>	runtime.RawExtension	The new resource being created or updated (JSON-encoded)
<code>request.oldObject</code>	runtime.RawExtension	The existing resource (only present for UPDATE operations)
<code>request.userInfo</code>	UserInfo	Information about the user making the request
<code>request.dryRun</code>	bool	True if this is a dry-run operation (kubectl apply --dry-run)

AdmissionResponse Structure:

Field	Type	Description
<code>uid</code>	string	Must match the request UID
<code>allowed</code>	bool	Whether the operation should be permitted
<code>result.status</code>	string	"Success" for allowed operations, "Failure" for denied
<code>result.message</code>	string	Human-readable explanation for denial
<code>result.reason</code>	string	Machine-readable reason code
<code>result.code</code>	int32	HTTP status code (400 for validation errors, 403 for policy violations)

The validation logic should be **idempotent** and **stateless** - the same resource should always produce the same validation result regardless of when the webhook is called. This is critical because Kubernetes may call the webhook multiple times for the same resource during retries or internal processing.

Database Validation Rules Implementation:

Consider implementing these validation rules for the `Database` custom resource:

- 1. Version Compatibility Validation:** Ensure the requested database version is supported and that version downgrades follow safe migration paths
- 2. Resource Constraint Validation:** Verify that CPU and memory requests are within organizational limits and that storage size is appropriate for the database version
- 3. Backup Schedule Validation:** Validate cron expressions in `BackupSchedule` and ensure backup frequency meets compliance requirements
- 4. Configuration Validation:** Check that database configuration parameters in the `Config` map are valid for the specified version
- 5. Update Path Validation:** For UPDATE operations, ensure changes don't violate constraints (e.g., storage size cannot be decreased, replica count changes follow safe scaling patterns)

Decision: Client-Side vs Server-Side Validation Balance

- **Context:** Validation can occur in multiple places - CRD schema, admission webhooks, and reconciliation logic
- **Options Considered:** All validation in webhooks, all validation in reconciliation, layered validation with different concerns at each level
- **Decision:** Layered validation with schema for structure, webhooks for business rules, reconciliation for runtime constraints
- **Rationale:** Schema validation provides immediate feedback with minimal latency, webhooks handle complex business logic synchronously, reconciliation handles constraints that depend on cluster state
- **Consequences:** Multiple validation layers require careful coordination, but provides optimal user experience and strong consistency guarantees

Common Validation Patterns:

Validation Type	Implementation Approach	Example
Cross-field validation	Compare multiple fields in same resource	Validate that memory limits are greater than memory requests
Enumeration validation	Check field values against allowed lists	Ensure database version is in supported versions list
External dependency validation	Query external systems or APIs	Verify that specified storage class exists and is available
Policy compliance validation	Apply organizational rules and constraints	Check that resource requests don't exceed namespace quotas
Change safety validation	Compare old and new resources for updates	Prevent scaling down below minimum replica count

⚠ Pitfall: Webhook Validation Loops

A common mistake is implementing validation logic that can change its decision based on external factors like current time or cluster state. For example, validating that "backups can only be scheduled during maintenance windows" might reject the same resource at different times of day. This creates confusing user experiences where `kubectl apply` succeeds at one time but fails at another for the same YAML.

Instead, validate the structural correctness of the schedule but let the controller's reconciliation logic handle timing-based decisions. The webhook should ensure the cron expression is valid, but the controller should skip backups that fall outside maintenance windows.

Error Response Best Practices:

When rejecting a resource, provide actionable error messages that help users fix the problem:

```
✗ Bad: "Invalid configuration"
✓ Good: "spec.config.max_connections: value 10000 exceeds maximum allowed value of 1000 for database version 14.x"

✗ Bad: "Update not allowed"
✓ Good: "spec.replicas: cannot decrease replica count from 5 to 2 in a single operation. Scale down gradually by reducing count by at most 1 replica per operation."
```

Mutating Webhook Implementation

Mutating admission webhooks modify incoming resources by injecting default values, adding required metadata, and transforming fields to ensure consistency across the cluster. These webhooks run before validating webhooks in the admission chain, allowing them to normalize resources before validation occurs.

The mutation process uses **JSON Patch** operations to describe changes to the incoming resource. Rather than returning a modified copy of the entire resource, the webhook returns a list of patch operations that the API server applies to the original request. This approach is more efficient and provides better conflict resolution when multiple mutating webhooks modify the same resource.

JSON Patch Operations:

Operation	Purpose	Example
add	Insert new field or array element	Add default storage size
replace	Change value of existing field	Normalize version string format
remove	Delete field or array element	Remove deprecated fields
move	Relocate field within document	Reorganize configuration structure
copy	Duplicate field value to another location	Copy name to display name
test	Verify field has expected value	Ensure field hasn't changed since mutation started

Default Value Injection Strategy:

The `Default` method should implement smart defaulting that considers the environment and context:

1. **Environment-Aware Defaults:** Set different default resource requests for development vs production namespaces
2. **Version-Appropriate Defaults:** Apply defaults that are appropriate for the specified database version
3. **Organizational Policy Defaults:** Inject standard labels, annotations, and configuration that align with organizational policies
4. **Resource-Dependent Defaults:** Set memory requests based on storage size, adjust replica counts based on resource tier

Database Default Values Implementation:

Field	Default Logic	Rationale
<code>spec.replicas</code>	1 for development namespaces, 3 for production	Ensures high availability in production while keeping development lightweight
<code>spec.version</code>	Latest stable version from supported versions list	Provides security updates while avoiding bleeding-edge releases
<code>spec.storageSize</code>	"20Gi" for single replica, "100Gi" for multi-replica	Scales storage allocation with expected load
<code>spec.backupSchedule</code>	"0 2 * * *" (daily at 2 AM) if not specified	Ensures data protection without manual configuration
<code>spec.resources.requests.memory</code>	Calculated as <code>storageSize / 10</code>	Provides reasonable memory allocation based on database size
<code>spec.resources.requests.cpu</code>	"500m" for development, "2" for production	Balances performance with resource efficiency

The mutation logic should be **deterministic** - given the same input resource and cluster state, it should always produce the same mutations. This ensures that repeated applications of the same YAML file don't cause unnecessary updates.

Decision: Default Value Source Hierarchy

- **Context:** Default values can come from multiple sources: CRD schema defaults, mutating webhooks, controller logic, or external configuration
- **Options Considered:** CRD schema defaults only, webhook-based defaulting, external configuration service, hybrid approach
- **Decision:** Hybrid approach with CRD schema for simple defaults, webhook for context-aware defaults
- **Rationale:** CRD schema defaults work well for static values but cannot consider namespace, user context, or organizational policies. Webhooks can implement sophisticated defaulting logic.
- **Consequences:** More complex implementation but enables smart defaults that improve user experience and enforce organizational standards

Mutation Ordering Considerations:

When multiple mutating webhooks are registered for the same resource type, Kubernetes calls them in alphabetical order by webhook name. This ordering is critical because later webhooks see the mutations applied by earlier webhooks. Design your webhook names carefully to ensure proper ordering:

- `000-baseline-defaults` : Sets fundamental default values
- `500-policy-injection` : Adds organizational policy requirements
- `900-finalization` : Performs final transformations and cleanup

Field Transformation Patterns:

Transformation Type	Use Case	Implementation
Normalization	Convert user input to canonical format	Transform "postgres:14" to "postgresql-14.6"
Policy Injection	Add required organizational metadata	Inject cost center labels, security scanning annotations
Resource Calculation	Derive resource requirements from high-level specs	Calculate memory requests from connection pool size
Configuration Generation	Generate complex config from simple parameters	Build database configuration map from performance tier
Cross-Resource References	Establish relationships between resources	Set owner references, add finalizers

⚠ Pitfall: Mutation Conflicts and Overwrites

When implementing mutating webhooks, be careful not to overwrite values that users explicitly set. A common mistake is unconditionally setting default values without checking if the field is already populated:

 Bad: Always set `spec.replicas = 3` for production
 Good: Set `spec.replicas = 3` only if `spec.replicas` is `nil` and `namespace` has `production` label

This ensures that users can override defaults when needed while still getting sensible default behavior for most cases.

TLS Certificate Management

Admission webhooks require HTTPS communication between the Kubernetes API server and the webhook endpoint. This requirement exists because webhooks handle sensitive resource data and admission decisions that affect cluster security. The API server validates webhook certificates against a trusted CA certificate, making certificate management a critical operational concern.

The webhook server must present a valid TLS certificate that matches its service name and is signed by a Certificate Authority that the API server trusts. Kubernetes provides several approaches for managing these certificates, each with different complexity and automation characteristics.

Certificate Management Approaches:

Approach	Pros	Cons	Best For
cert-manager	Automatic renewal, industry standard, handles CA distribution	Additional dependency, requires cert-manager installation	Production deployments
Self-signed certificates	No external dependencies, simple for development	Manual renewal, complex CA distribution, security concerns	Development and testing
External CA integration	Integrates with organizational PKI, auditable	Complex setup, requires external CA access	Enterprise environments
Kubernetes CSR API	Built into Kubernetes, no external dependencies	Manual approval workflow, limited automation	Controlled environments

Decision: Certificate Management Strategy

- Context:** Webhook TLS certificates need provisioning, renewal, and distribution to API server configuration
- Options Considered:** Manual certificate management, cert-manager automation, Kubernetes CSR API, external CA integration
- Decision:** cert-manager for production, self-signed for development
- Rationale:** cert-manager provides robust automation with industry-standard ACME protocol support, while self-signed certificates reduce dependencies during development
- Consequences:** Production deployments require cert-manager installation, but operational overhead is minimal once configured

cert-manager Integration:

cert-manager automates the certificate lifecycle by creating `Certificate` resources that define the desired certificate properties. The cert-manager controller provisions certificates from various issuers (Let's Encrypt, HashiCorp Vault, self-signed CA) and automatically renews them before expiration.

Certificate Resource Configuration:

Field	Value	Purpose
<code>spec.secretName</code>	<code>database-webhook-certs</code>	Secret where certificate and key are stored
<code>spec.issuerRef.name</code>	<code>selfsigned-issuer</code>	Certificate issuer (CA) to use for signing
<code>spec.issuerRef.kind</code>	<code>ClusterIssuer</code>	Cluster-wide or namespace-scoped issuer
<code>spec.dnsNames</code>	<code>[database-webhook.system.svc, database-webhook.system.svc.cluster.local]</code>	DNS names that certificate should be valid for
<code>spec.duration</code>	<code>8760h (1 year)</code>	Certificate validity period
<code>spec.renewBefore</code>	<code>720h (30 days)</code>	How long before expiry to renew

Webhook Configuration TLS Setup:

The `ValidatingAdmissionWebhook` and `MutatingAdmissionWebhook` configurations must reference the CA certificate that signed the webhook's TLS certificate. cert-manager can automatically inject this CA bundle using the `cert-manager.io/inject-ca-from` annotation:

```

metadata:
  annotations:
    cert-manager.io/inject-ca-from: system/selfsigned-issuer

```

YAML

This annotation instructs cert-manager to populate the `caBundle` field in the webhook configuration with the appropriate CA certificate, ensuring the API server can validate the webhook's TLS certificate.

Certificate Rotation and High Availability:

Certificate rotation happens automatically when cert-manager renews certificates. The controller updates the secret containing the certificate and private key, and the webhook server must reload these certificates without downtime. Implement certificate reloading by:

1. **Watching the certificate secret** for changes using a Kubernetes informer
2. **Gracefully reloading** the TLS configuration without dropping existing connections
3. **Health checking** the new certificate before switching to ensure validity
4. **Logging certificate events** for operational visibility and debugging

For high availability deployments with multiple webhook replicas, ensure that certificate updates propagate to all instances consistently. Use a shared secret for certificate storage and implement proper leader election if coordination is needed during certificate transitions.

Development Certificate Setup:

For development environments, self-signed certificates provide a simpler alternative that avoids external dependencies:

1. **Generate a CA key and certificate** for signing webhook certificates
2. **Create a webhook certificate** signed by the development CA with appropriate DNS names
3. **Store certificates in Kubernetes secrets** using the same naming conventions as production
4. **Configure webhook registration** with the CA bundle for validation

⚠ Pitfall: Certificate DNS Name Mismatches

A common certificate error occurs when the webhook certificate doesn't include all the DNS names that the API server might use to connect to the webhook service. The certificate must include:

- `service-name.namespace.svc` (standard service DNS name)
- `service-name.namespace.svc.cluster.local` (fully qualified service DNS name)
- Any additional DNS names if using custom networking or ingress

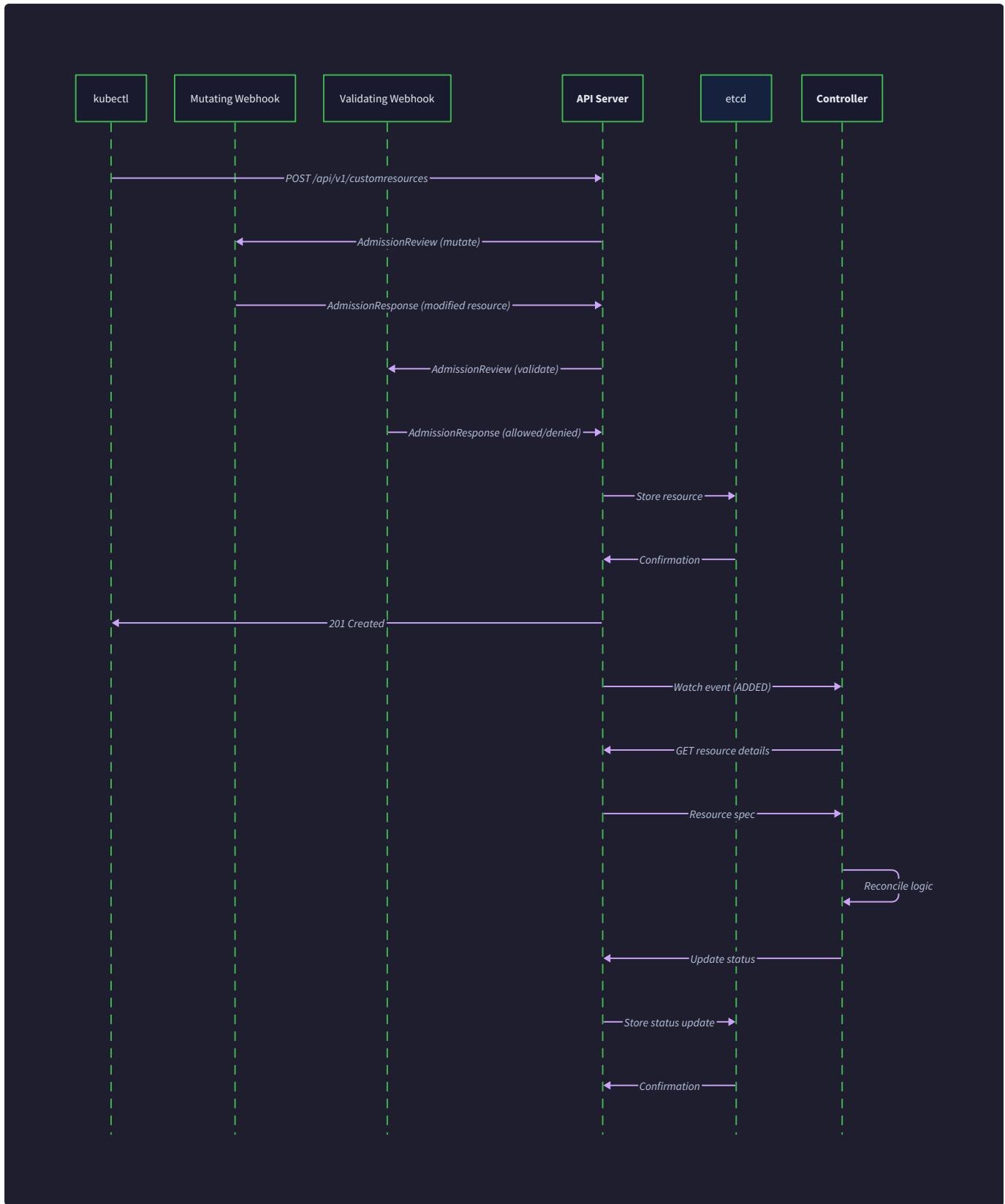
Missing DNS names in the certificate cause connection failures with TLS verification errors that can be difficult to diagnose.

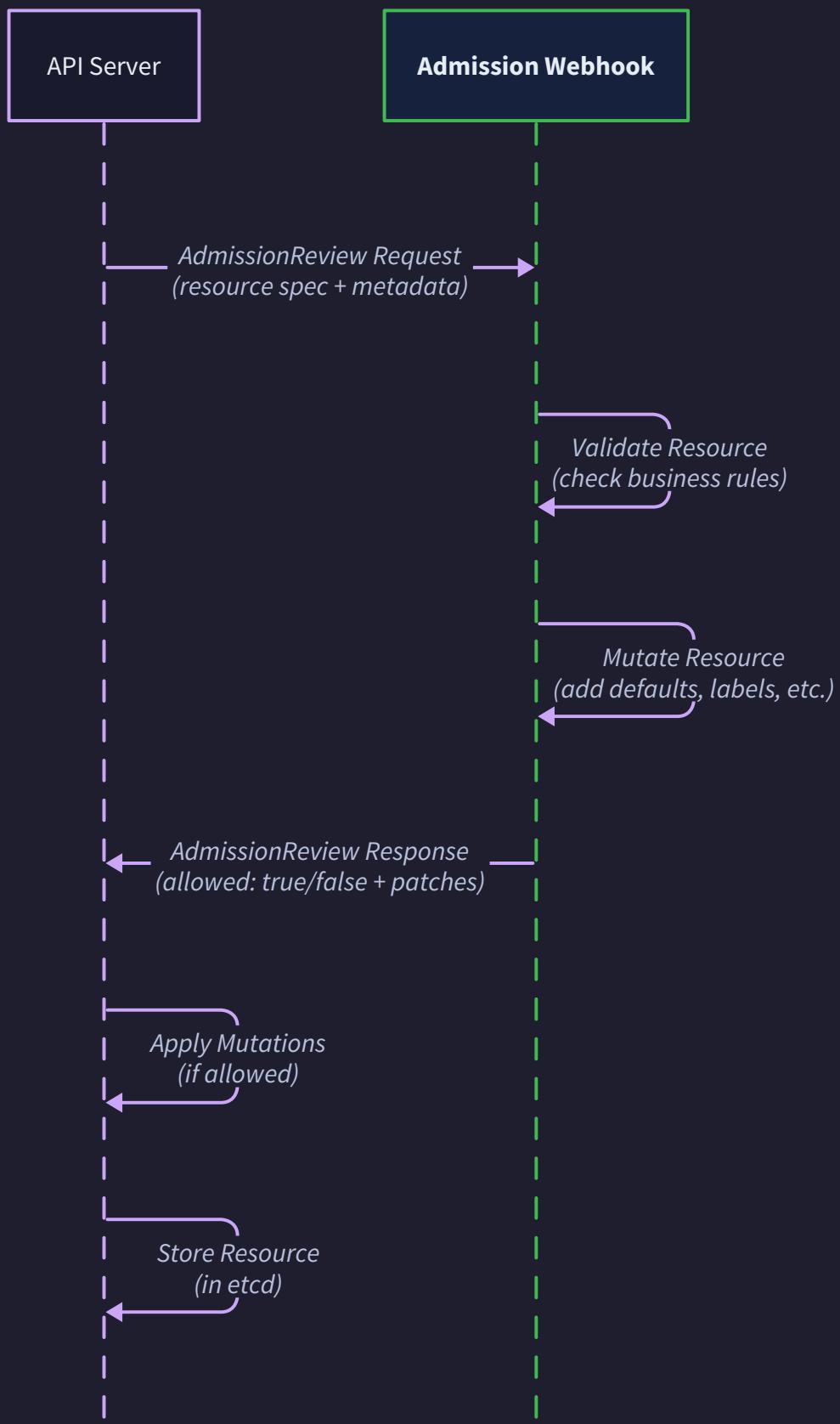
Certificate Monitoring and Observability:

Implement monitoring for certificate health to prevent webhook failures due to certificate expiration:

Metric	Purpose	Alert Threshold
Certificate expiry time	Track time until certificate expiration	Alert 7 days before expiry
Certificate renewal success	Monitor automatic renewal operations	Alert on renewal failures
Webhook TLS handshake errors	Track TLS connection failures	Alert on increased error rate
cert-manager controller health	Ensure cert-manager is running properly	Alert on controller downtime

Certificate-related webhook failures often manifest as cryptic TLS errors in API server logs. Establish clear runbooks for certificate troubleshooting that cover common scenarios like expired certificates, CA bundle mismatches, and renewal failures.





Implementation Guidance

Building admission webhooks requires careful attention to HTTP server implementation, JSON processing, and certificate management. This section provides complete infrastructure code and implementation skeletons that handle the operational complexities while leaving the core validation and mutation logic for you to implement.

Technology Recommendations:

Component	Simple Option	Advanced Option
HTTP Server	net/http with standard library	controller-runtime webhook server
JSON Processing	encoding/json for AdmissionReview	controller-runtime admission package
Certificate Management	Self-signed certificates with openssl	cert-manager with Certificate resources
TLS Configuration	tls.LoadX509KeyPair for static certificates	Dynamic certificate reloading with fsnotify
Testing Framework	httptest for webhook endpoint testing	envtest with real admission controller

Recommended File Structure:

```
internal/webhooks/
    admission.go          ← AdmissionReview request/response handling
    database_defaulter.go ← Mutating webhook implementation
    database_validator.go ← Validating webhook implementation
    webhook_server.go     ← HTTP server and certificate management
    admission_test.go     ← Unit tests for webhook logic
config/webhook/
    manifests.yaml        ← Webhook registration YAML
    certificate.yaml      ← cert-manager Certificate resource
deploy/certificates/
    ca.crt               ← Development CA certificate
    tls.crt              ← Development webhook certificate
    tls.key              ← Development webhook private key
```

Complete Admission Infrastructure (admission.go):

```
package webhooks

import (
    "context"
    "encoding/json"
    "fmt"
    "net/http"

    admissionv1 "k8s.io/api/admission/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/apimachinery/pkg/runtime/serializer"
    utilruntime "k8s.io/apimachinery/pkg/util/runtime"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/log"
    "sigs.k8s.io/controller-runtime/pkg/webhook/admission"
)

var (
    scheme = runtime.NewScheme()
    codecs = serializer.NewCodecFactory(scheme)
)

func init() {
    utilruntime.Must(admissionv1.AddToScheme(scheme))
}

// AdmissionHandler wraps validation and mutation logic in HTTP handler

type AdmissionHandler struct {
    decoder *admission.Decoder
}

// NewAdmissionHandler creates an admission handler with proper decoding setup

func NewAdmissionHandler() *AdmissionHandler {
    return &AdmissionHandler{
        decoder: admission.NewDecoder(scheme),
```

```
}

}

// ServeHTTP handles incoming admission review requests

func (h *AdmissionHandler) ServeHTTP(w http.ResponseWriter, r *http.Request) {
    logger := log.FromContext(r.Context())

    // Parse AdmissionReview from request body

    body := make([]byte, r.ContentLength)

    if _, err := r.Body.Read(body); err != nil {

        logger.Error(err, "Failed to read request body")

        http.Error(w, err.Error(), http.StatusBadRequest)

        return
    }

    var admissionReview admissionv1.AdmissionReview

    if err := json.Unmarshal(body, &admissionReview); err != nil {

        logger.Error(err, "Failed to unmarshal AdmissionReview")

        http.Error(w, err.Error(), http.StatusBadRequest)

        return
    }

    // Process the admission request

    response := h.processAdmissionReview(r.Context(), &admissionReview)

    // Send AdmissionReview response

    responseBytes, err := json.Marshal(response)

    if err != nil {

        logger.Error(err, "Failed to marshal AdmissionResponse")

        http.Error(w, err.Error(), http.StatusInternalServerError)

        return
    }

    w.Header().Set("Content-Type", "application/json")
```

```

w.Write(responseBytes)

}

func (h *AdmissionHandler) processAdmissionReview(ctx context.Context,
    review *admissionv1.AdmissionReview) *admissionv1.AdmissionReview {

    req := review.Request

    response := &admissionv1.AdmissionResponse{
        UID: req.UID,
        Allowed: true,
    }

    // Route to appropriate handler based on operation and resource
    switch {
    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Create:
        response = h.handleDatabaseCreate(ctx, req)
    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Update:
        response = h.handleDatabaseUpdate(ctx, req)
    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Delete:
        response = h.handleDatabaseDelete(ctx, req)
    default:
        // Allow operations on unknown resources
        response.Allowed = true
    }

    return &admissionv1.AdmissionReview{
        TypeMeta: metav1.TypeMeta{
            APIVersion: "admission.k8s.io/v1",
            Kind:       "AdmissionReview",
        },
        Response: response,
    }
}

```

Core Validation Logic Skeleton (`database_validator.go`):

```
package webhooks

import (
    "context"
    "encoding/json"
    "fmt"
    "strconv"
    "strings"

    admissionv1 "k8s.io/api/admission/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"

    databasev1 "github.com/example/database-operator/api/v1"
)

// handleDatabaseCreate validates new Database resources

func (h *AdmissionHandler) handleDatabaseCreate(ctx context.Context,
    req *admissionv1.AdmissionRequest) *admissionv1.AdmissionResponse {

    var database databasev1.Database

    if err := json.Unmarshal(req.Object.Raw, &database); err != nil {
        return &admissionv1.AdmissionResponse{
            UID: req.UID,
            Allowed: false,
            Result: &metav1.Status{
                Message: fmt.Sprintf("Failed to decode Database: %v", err),
            },
        }
    }

    // TODO 1: Validate spec.version is in supported versions list
    // Hint: Define supportedVersions = []string{"14.6", "15.1", "15.2"}

    // TODO 2: Validate spec.storageSize format and minimum size
    // Hint: Use resource.ParseQuantity() and compare against minimum (1Gi)
```

```
// TODO 3: Validate spec.replicas is within allowed range (1-10)

// Hint: Check database.Spec.Replicas >= 1 && database.Spec.Replicas <= 10


// TODO 4: Validate spec.backupSchedule cron expression if provided

// Hint: Use "github.com/robfig/cron/v3" parser.Parse()


// TODO 5: Validate spec.config map contains only allowed keys

// Hint: Define allowedConfigKeys and check all keys in database.Spec.Config


return &admissionv1.AdmissionResponse{

    UID: req.UID,

    Allowed: true,

}

}

// handleDatabaseUpdate validates changes to existing Database resources

func (h *AdmissionHandler) handleDatabaseUpdate(ctx context.Context,
    req *admissionv1.AdmissionRequest) *admissionv1.AdmissionResponse {

    var oldDatabase, newDatabase databasev1.Database

    if err := json.Unmarshal(req.OldObject.Raw, &oldDatabase); err != nil {

        return &admissionv1.AdmissionResponse{

            UID: req.UID,

            Allowed: false,

            Result: &metav1.Status{

                Message: fmt.Sprintf("Failed to decode old Database: %v", err),
            },
        }
    }

    if err := json.Unmarshal(req.Object.Raw, &newDatabase); err != nil {

        return &admissionv1.AdmissionResponse{

            UID: req.UID,
```

```
    Allowed: false,
    Result: &metav1.Status{
        Message: fmt.Sprintf("Failed to decode new Database: %v", err),
    },
}

}

// TODO 1: Validate storage size is not decreased
// Hint: Parse both sizes and compare: newSize >= oldSize

// TODO 2: Validate version changes are upgrade-only (no downgrades)
// Hint: Compare version numbers, reject if newVersion < oldVersion

// TODO 3: Validate replica scaling follows safe patterns
// Hint: Allow scaling up freely, but limit scaling down to max 1 replica per operation

// TODO 4: Validate immutable fields haven't changed
// Hint: Check that critical fields like database name/type remain unchanged

return &admissionv1.AdmissionResponse{
    UID: req.UID,
    Allowed: true,
}

}

// handleDatabaseDelete validates Database deletion requests

func (h *AdmissionHandler) handleDatabaseDelete(ctx context.Context,
    req *admissionv1.AdmissionRequest) *admissionv1.AdmissionResponse {

    // TODO 1: Check if database has dependent resources that need cleanup
    // Hint: Look for finalizers that indicate cleanup is required

    // TODO 2: Validate deletion is allowed based on database state
    // Hint: Prevent deletion if database is in "backup-in-progress" state
}
```

```
// TODO 3: Check for protection annotations that prevent deletion  
  
// Hint: Look for "database.example.com/protect: true" annotation  
  
  
return &admissionv1.AdmissionResponse{  
  
    UID: req.UID,  
  
    Allowed: true,  
  
}  
  
}
```

Core Mutation Logic Skeleton (`database_defaulter.go`):

```
package webhooks

import (
    "context"
    "encoding/json"
    "fmt"

    admissionv1 "k8s.io/api/admission/v1"
    corev1 "k8s.io/api/core/v1"
    "k8s.io/apimachinery/pkg/api/resource"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"

    databasev1 "github.com/example/database-operator/api/v1"
)

// handleDatabaseDefaulting applies default values to new Database resources

func (h *AdmissionHandler) handleDatabaseDefaulting(ctx context.Context,
    req *admissionv1.AdmissionRequest) *admissionv1.AdmissionResponse {

    var database databasev1.Database

    if err := json.Unmarshal(req.Object.Raw, &database); err != nil {
        return &admissionv1.AdmissionResponse{
            UID: req.UID,
            Allowed: false,
            Result: &metav1.Status{
                Message: fmt.Sprintf("Failed to decode Database: %v", err),
            },
        }
    }

    var patches []JSONPatch

    // TODO 1: Set default replicas based on namespace environment
    // Hint: Check namespace labels for "environment=production" → 3 replicas, else 1
```

```
// TODO 2: Set default version to latest stable if not specified

// Hint: if database.Spec.Version == "" { patches = append(...) }

// TODO 3: Set default storage size based on replica count

// Hint: Single replica → 20Gi, Multiple replicas → 100Gi

// TODO 4: Generate default backup schedule if not provided

// Hint: Use "0 2 * * *" for daily backups at 2 AM

// TODO 5: Calculate default resource requests based on storage size

// Hint: Memory = StorageSize / 10, CPU = "500m" for dev, "2" for prod

// TODO 6: Inject organizational labels and annotations

// Hint: Add cost-center, team, and monitoring labels

if len(patches) == 0 {

    return &admissionv1.AdmissionResponse{
        UID: req.UID,
        Allowed: true,
    }
}

patchBytes, err := json.Marshal(patches)

if err != nil {

    return &admissionv1.AdmissionResponse{
        UID: req.UID,
        Allowed: false,
        Result: &metav1.Status{
            Message: fmt.Sprintf("Failed to marshal patches: %v", err),
        },
    }
}
```

```

patchType := admissionv1.PatchTypeJSONPatch

return &admissionv1.AdmissionResponse{
    UID: req.UID,
    Allowed: true,
    Patch: patchBytes,
    PatchType: &patchType,
}

}

// JSONPatch represents a JSON patch operation

type JSONPatch struct {
    Op      string      `json:"op"`
    Path   string      `json:"path"`
    Value interface{} `json:"value,omitempty"`
}

// Helper function to create an "add" patch operation

func createAddPatch(path string, value interface{}) JSONPatch {
    return JSONPatch{
        Op:      "add",
        Path:   path,
        Value: value,
    }
}

// Helper function to create a "replace" patch operation

func createReplacePatch(path string, value interface{}) JSONPatch {
    return JSONPatch{
        Op:      "replace",
        Path:   path,
        Value: value,
    }
}

```

Complete Webhook Server Setup (`webhook_server.go`):

```
package webhooks
```

GO

```
import (
    "context"
    "crypto/tls"
    "fmt"
    "net/http"
    "path/filepath"
    "time"

    "k8s.io/apimachinery/pkg/util/wait"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/healthz"
    "sigs.k8s.io/controller-runtime/pkg/log"
    "sigs.k8s.io/controller-runtime/pkg/manager"
    "sigs.k8s.io/controller-runtime/pkg/webhook"
)
```

```
const (
```

```
    WebhookPort = 9443
    CertDir     = "/tmp/k8s-webhook-server/serving-certs"
    CertName    = "tls.crt"
    KeyName     = "tls.key"
)
```

```
// SetupWebhookServer configures the webhook server with TLS and health checks
```

```
func SetupWebhookServer(mgr ctrl.Manager) error {
```

```
    hookServer := mgr.GetWebhookServer()
```

```
// Configure TLS certificate paths
```

```
    hookServer.Port = WebhookPort
```

```
    hookServer.CertDir = CertDir
```

```
    hookServer.CertName = CertName
```

```
    hookServer.KeyName = KeyName
```

```
// Register admission handlers

admissionHandler := NewAdmissionHandler()

hookServer.Register("/mutate", &webhook.Admission{Handler: admissionHandler})

hookServer.Register("/validate", &webhook.Admission{Handler: admissionHandler})


// Add health check endpoints

if err := mgr.AddHealthzCheck("webhook", healthz.Ping); err != nil {

    return fmt.Errorf("failed to add health check: %w", err)
}

if err := mgr.AddReadyzCheck("webhook", healthz.Ping); err != nil {

    return fmt.Errorf("failed to add ready check: %w", err)
}

return nil
}

// CertificateReloader handles dynamic certificate reloading

type CertificateReloader struct {

    certPath string

    keyPath string

    cert     *tls.Certificate
}

// NewCertificateReloader creates a certificate reloader for the given paths

func NewCertificateReloader(certPath, keyPath string) (*CertificateReloader, error) {

    reloader := &CertificateReloader{

        certPath: certPath,

        keyPath: keyPath,
    }

    if err := reloader.loadCertificate(); err != nil {

        return nil, fmt.Errorf("failed to load initial certificate: %w", err)
    }
}
```

```

    return reloader, nil
}

func (r *CertificateVerifier) Verify() error {
    cert, err := tls.LoadX509KeyPair(r.certPath, r.keyPath)

    if err != nil {
        return err
    }

    r.cert = &cert

    return nil
}

// GetCertificate returns the current certificate (implements tls.Config.GetCertificate)
func (r *CertificateVerifier) GetCertificate(*tls.ClientHelloInfo) (*tls.Certificate, error) {
    return r.cert, nil
}

// Start begins watching for certificate changes and reloading automatically
func (r *CertificateVerifier) Start(ctx context.Context) error {
    logger := log.FromContext(ctx)

    return wait.PollImmediateUntil(30*time.Second, func() (bool, error) {
        if err := r.loadCertificate(); err != nil {
            logger.Error(err, "Failed to reload certificate")
            return false, nil // Continue polling
        }

        logger.Info("Certificate reloaded successfully")
        return false, nil // Continue polling
    }, ctx.Done())
}

```

Language-Specific Implementation Hints:

- **JSON Patch Operations:** Use the `JSONPatch` struct to build mutation operations. The `path` field uses JSON Pointer syntax (e.g., `/spec/replicas`, `/metadata/labels/environment`)
- **AdmissionReview Handling:** Always copy the request UID to the response UID - this is required for proper request tracking

- **Error Response Formatting:** Use `metav1.Status` for structured error responses with appropriate HTTP status codes (400 for validation errors, 500 for internal errors)
- **TLS Certificate Loading:** Use `tls.LoadX509KeyPair()` for loading certificates and implement `GetCertificate` callback for dynamic reloading
- **Resource Quantity Parsing:** Use `resource.ParseQuantity()` for validating storage sizes and resource limits

Milestone Checkpoint:

After implementing the webhook logic, verify the following behavior:

1. **Webhook Registration:** `kubectl get validatingadmissionwebhooks` and `kubectl get mutatingadmissionwebhooks` show your webhooks
2. **Certificate Validation:** Check that webhook pod logs show successful TLS certificate loading
3. **Mutation Behavior:** Create a `Database` resource with minimal spec and verify defaults are applied
4. **Validation Behavior:** Try creating invalid `Database` resources and verify they are rejected with clear error messages
5. **Update Validation:** Try invalid updates (like decreasing storage size) and verify they are properly blocked

Common Debugging Scenarios:

Symptom	Likely Cause	How to Diagnose	Fix
"connection refused" errors	Webhook server not running or wrong port	Check pod logs and service endpoints	Verify webhook server is listening on correct port
"certificate verification failed"	CA bundle mismatch or expired certificate	Check webhook configuration caBundle vs actual cert	Update caBundle or renew certificate
"webhook timeout" errors	Webhook processing too slow or hanging	Check webhook response time and goroutine leaks	Optimize validation logic, add timeouts
"admission webhook denied" without clear reason	Poor error message formatting	Check <code>AdmissionResponse.Result.Message</code> field	Improve error message content and structure
Mutations not being applied	Webhook called but patches malformed	Check JSON patch syntax and field paths	Validate patch operations against JSON Pointer spec

Component Interactions and Data Flow

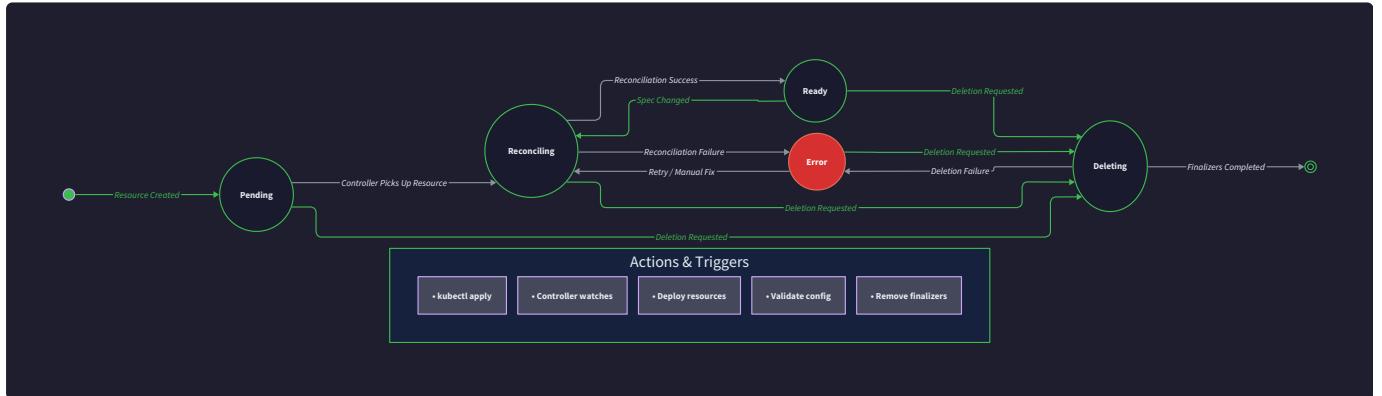
Milestone(s): Milestone 3 (Reconciliation Loop), Milestone 4 (Webhooks), Milestone 5 (Testing & Deployment) - demonstrates the orchestrated behavior between controller, webhooks, and API server during complete resource lifecycle operations

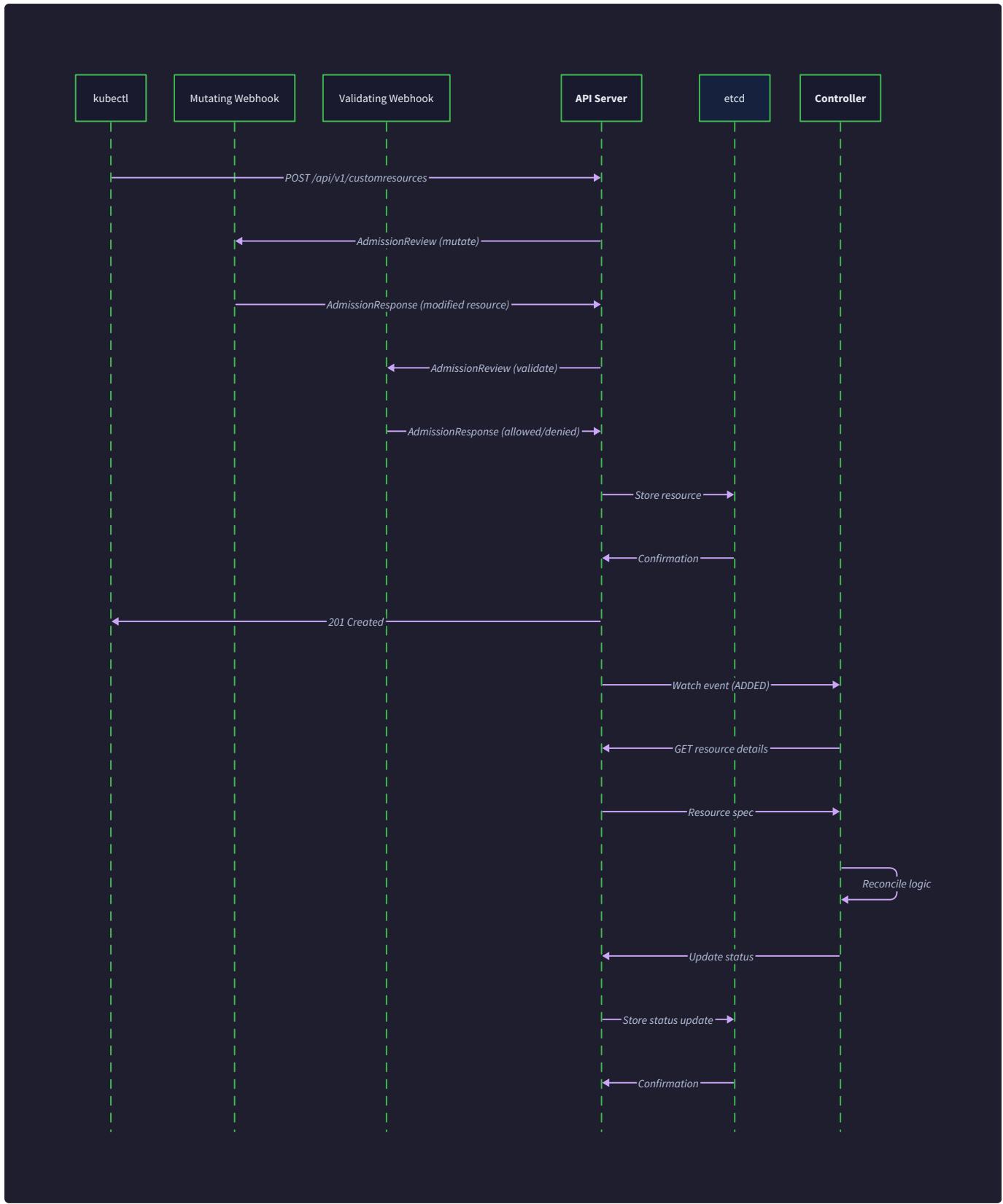
The interaction between the operator's components creates a sophisticated orchestration that maintains system consistency while handling the complexities of distributed operations. Understanding these interaction patterns is crucial for building operators that behave predictably under both normal conditions and failure scenarios.

Mental Model: The Orchestra Conductor

Think of the Kubernetes API server as an orchestra conductor coordinating multiple musicians (the controller and webhooks). When someone in the audience (a user) requests a song (creates a resource), the conductor first asks the first violin section (mutating webhooks) to tune their instruments (apply defaults), then asks the brass section (validating webhooks) to confirm they can play the piece (validate the request). Only after everyone agrees does the conductor store the sheet music (persist the resource) and signal the orchestra to begin playing (trigger controller reconciliation). Throughout the performance, musicians continuously watch the conductor for tempo changes (status updates) and adjust their playing accordingly.

This mental model helps explain why webhook failures can block resource creation (the brass section refusing to play stops the entire orchestra) and why controller reconciliation happens asynchronously (the orchestra continues playing even if some audience members leave).





Resource Creation Flow

The resource creation flow represents the most complex interaction pattern in the operator ecosystem, involving multiple validation and processing stages that must complete successfully before the controller begins reconciliation. This flow demonstrates how Kubernetes' admission control system provides strong consistency guarantees while allowing operators to extend the platform's behavior.

Pre-Storage Validation Pipeline

When a user submits a `Database` resource through `kubectl apply`, the API server initiates a multi-stage validation pipeline before the resource reaches persistent storage. This pipeline ensures that only valid, properly configured resources enter the cluster state, preventing downstream controller errors and maintaining system integrity.

The API server first performs schema validation against the CRD's OpenAPI v3 specification, rejecting resources that violate basic type constraints or required field rules. This validation catches structural errors like missing required fields, incorrect field types, or values that exceed defined string length limits. Schema validation provides the first line of defense against malformed resources, ensuring that subsequent processing stages can assume basic structural correctness.

Validation Stage	Purpose	Failure Impact	Recovery Action
Schema Validation	Type safety and required fields	Immediate rejection with schema error	Fix resource definition and reapply
Mutating Webhook	Default value injection	Webhook timeout or denial	Check webhook pod logs and certificate
Validating Webhook	Business rule enforcement	Resource creation blocked	Address validation error and resubmit
Storage Layer	Optimistic concurrency control	Conflict on resource version	Retry with fresh resource version

After schema validation succeeds, the API server invokes any configured mutating webhooks in the order specified by their admission registration. The mutating webhook receives an `AdmissionReview` containing the proposed resource and can modify field values before storage occurs. Common mutations include injecting default values for optional fields, adding labels or annotations for operational tracking, or transforming user-friendly field values into internal representations.

Design Insight: Mutation Order Matters

Multiple mutating webhooks execute in a defined sequence, with each webhook receiving the potentially modified resource from previous webhooks. This creates a pipeline where early webhooks can set defaults that later webhooks can validate or further modify. Operators must consider this ordering when designing webhook logic to avoid conflicts or unexpected interactions between different mutation policies.

Mutating Webhook Processing

The mutating webhook implementation receives admission review requests containing the raw resource JSON along with metadata about the operation context. The webhook deserializes this JSON into the `Database` struct, applies default values for any unspecified optional fields, and returns a JSON patch describing the required modifications.

Mutating Webhook Request Processing:

1. API server serializes Database resource to JSON
2. API server constructs AdmissionReview with resource JSON and operation metadata
3. API server sends HTTPS POST to webhook endpoint with AdmissionReview
4. Webhook deserializes AdmissionReview and extracts Database resource
5. Webhook applies defaulting logic to populate missing optional fields
6. Webhook generates JSON patch operations for modified fields
7. Webhook returns AdmissionResponse with patch and admission decision
8. API server applies patch to original resource JSON
9. Modified resource proceeds to validating webhook stage

The JSON patch format provides precise control over field modifications, supporting operations like adding missing fields, replacing existing values, or removing unwanted elements. This precision ensures that webhooks only modify their intended fields without accidentally affecting other resource properties.

Default Field	Condition	Applied Value	Justification
<code>StorageSize</code>	If empty or unspecified	"10Gi"	Reasonable default for development databases
<code>Replicas</code>	If zero or negative	1	Single replica provides basic functionality
<code>Version</code>	If empty string	"14.5"	Current stable PostgreSQL release
<code>Resources.Limits</code>	If nil	CPU: "500m", Memory: "512Mi"	Prevent resource exhaustion
<code>Config</code>	If nil	<code>map[string]string{"shared_buffers": "256MB"}</code>	Basic performance tuning

Validating Webhook Processing

After successful mutation, the API server invokes validating webhooks to enforce business rules and cross-field validation constraints that cannot be expressed in the OpenAPI schema. Unlike mutating webhooks, validating webhooks cannot modify the resource but can reject the entire operation by returning a denial response.

The validating webhook performs comprehensive business logic validation, checking constraints like database version compatibility, resource requirement feasibility, and configuration parameter validity. These validations often require external knowledge or complex logic that exceeds the capabilities of declarative schema validation.

Common validation scenarios include verifying that the requested `StorageSize` is larger than any existing persistent volume claims, ensuring that the specified `Version` is compatible with upgrade paths from currently running versions, and validating that custom configuration parameters in the `Config` map use valid PostgreSQL settings with appropriate values.

Decision: Webhook vs Schema Validation Split

- **Context:** Complex validation logic can be implemented in webhooks or expressed in CRD schemas using CEL expressions
- **Options Considered:**
 1. Implement all validation in admission webhooks with full programming flexibility
 2. Use CEL expressions in CRD schema for all validation logic
 3. Hybrid approach with schema validation for simple rules and webhooks for complex logic
- **Decision:** Hybrid approach with simple validations in schema and complex business rules in webhooks
- **Rationale:** Schema validation provides faster feedback and better error messages for simple constraints, while webhooks handle complex logic requiring external state or computations
- **Consequences:** Reduces webhook load for simple validation failures while maintaining flexibility for sophisticated business rules

Storage and Initial Status

Once all webhook validations pass, the API server commits the resource to etcd storage and initializes the status subresource with default condition values. The newly stored resource receives a unique `resourceVersion` that enables optimistic concurrency control for subsequent updates.

The initial status includes a `Progressing` condition set to `True` with a reason of "ResourceCreated", indicating that the controller should begin reconciliation processing. The `ObservedGeneration` field initializes to zero, signaling that no controller reconciliation has occurred yet.

Initial Status Field	Value	Purpose
Conditions	[{"Type": "Progressing", "Status": "True", "Reason": "ResourceCreated"}]	Indicates pending reconciliation
ReadyReplicas	0	No workloads created yet
ObservedGeneration	0	No reconciliation completed
Phase	"Pending"	Resource awaiting controller processing
Message	"Resource created, awaiting reconciliation"	Human-readable status

Update and Reconciliation Flow

The update and reconciliation flow demonstrates how the operator maintains system consistency when users modify existing resources or when external conditions change the actual cluster state. This flow showcases the declarative nature of Kubernetes operators, where controllers continuously work to align actual state with desired state regardless of what triggered the divergence.

Spec Change Detection

When users modify a `Database` resource specification, the API server increments the resource's `generation` field and processes the update through the same webhook pipeline used during creation. However, validating webhooks receive both the old and new resource versions, enabling validation logic that considers the implications of specific changes.

The controller's informer cache detects the specification change through its watch connection to the API server. The shared informer compares the new resource version against its cached copy and triggers an update event that enqueues the resource key for reconciliation processing.

- ```
Update Detection and Queuing:
1. User submits resource update via kubectl apply or API call
2. API server processes update through mutating and validating webhooks
3. API server increments generation field and stores modified resource
4. Controller informer receives watch event with updated resource
5. Event handler compares new generation against cached generation
6. If generation changed, handler enqueues resource key in work queue
7. Controller worker goroutine dequeues key and begins reconciliation
8. Reconciler fetches fresh resource version from cache
9. Reconciler compares generation vs observedGeneration to detect changes
```

The controller uses generation comparison to distinguish between user-initiated spec changes and controller-initiated status updates. Only when `metadata.generation` exceeds `status.observedGeneration` does the controller recognize that new reconciliation work is required.

### Incremental Reconciliation Strategy

Rather than recreating all owned resources during every reconciliation, the controller implements incremental updates that modify only the specific resources affected by the spec change. This approach minimizes disruption to running workloads and reduces the risk of cascading failures during updates.

The reconciler fetches the current state of all owned resources and compares their configuration against the desired state derived from the updated spec. For each owned resource, the controller determines whether creation, update, or deletion is required to align with the new desired state.

| Change Type            | Affected Resources     | Reconciliation Action             | Disruption Level   |
|------------------------|------------------------|-----------------------------------|--------------------|
| Replica count increase | StatefulSet, Service   | Update StatefulSet replicas field | Rolling scale-up   |
| Replica count decrease | StatefulSet, Service   | Update StatefulSet replicas field | Rolling scale-down |
| Storage size increase  | PersistentVolumeClaim  | Create new PVC, migrate data      | High disruption    |
| Version upgrade        | StatefulSet, ConfigMap | Update image and configuration    | Rolling restart    |
| Configuration change   | ConfigMap, StatefulSet | Update ConfigMap, restart pods    | Controlled restart |
| Resource limits change | StatefulSet            | Update pod template resources     | Rolling restart    |

### Critical Insight: Storage Size Immutability

Kubernetes PersistentVolumeClaims cannot be resized in all storage classes, and decreasing storage size is never supported. The controller must detect storage size changes and implement appropriate strategies like creating new volumes and orchestrating data migration rather than attempting in-place PVC modifications that will fail.

### Status Update Propagation

After completing reconciliation actions, the controller updates the `Database` status subresource to reflect the current state of owned resources and the progress of the reconciliation operation. Status updates use a separate API call to the status subresource, preventing race conditions between controller status updates and user spec modifications.

The controller aggregates status information from all owned resources, computing values like `ReadyReplicas` from the StatefulSet status and updating condition states based on the overall health of the database deployment. The `ObservedGeneration` field advances to match the current `metadata.generation`, indicating that reconciliation has processed the latest spec changes.

#### Status Update Process:

1. Controller completes reconciliation actions for owned resources
2. Controller queries current status of all owned resources
3. Controller aggregates health information from StatefulSet, Services, PVCs
4. Controller computes overall readiness state and condition values
5. Controller constructs updated DatabaseStatus with current state
6. Controller calls status subresource API to persist status changes
7. Controller sets `ObservedGeneration` to current `metadata.generation`
8. Controller updates `LastTransitionTime` for any changed conditions
9. Updated status becomes visible to users via `kubectl get database`

The status update includes human-readable messages that help users understand the current state and any issues that might prevent successful reconciliation. These messages appear in `kubectl describe` output and provide valuable debugging information when troubleshooting deployment problems.

### Deletion and Cleanup Flow

The deletion and cleanup flow represents the most complex interaction pattern because it must coordinate resource cleanup across multiple Kubernetes objects while respecting dependency relationships and handling potential failures during the cleanup process. This flow demonstrates how finalizers provide deletion safety in distributed systems where cleanup operations may take significant time or encounter transient failures.

#### Finalizer-Protected Deletion

When users delete a `Database` resource, Kubernetes does not immediately remove the resource from storage if finalizers are present in the `metadata.finalizers` array. Instead, the API server sets the `metadata.deletionTimestamp` field and waits for controllers to remove their finalizers after completing cleanup operations.

The `DatabaseReconciler` adds the `FINALIZER_NAME` to newly created resources during their first reconciliation, ensuring that deletion cannot proceed without controller involvement. This finalizer acts as a safety mechanism that prevents orphaned resources when users delete the parent `Database` object.

| Finalizer                                          | Purpose                | Cleanup Responsibility                            | Removal Condition                        |
|----------------------------------------------------|------------------------|---------------------------------------------------|------------------------------------------|
| <code>database.example.com/finalizer</code>        | Owned resource cleanup | Delete StatefulSet, Services, ConfigMaps, PVCs    | All owned resources successfully deleted |
| <code>database.example.com/backup-finalizer</code> | Backup completion      | Ensure backup job completes or fails              | Backup job reaches terminal state        |
| <code>external-secrets.io/finalizer</code>         | Secret cleanup         | Remove database credentials from external systems | External secret deletion confirmed       |

The controller detects deletion requests by checking for the presence of `metadata.deletionTimestamp` during reconciliation. When this timestamp is set, the controller switches from normal reconciliation logic to deletion processing, focusing exclusively on cleanup operations rather than maintaining desired state.

### Ordered Resource Cleanup

The cleanup process follows a specific order designed to minimize data loss risk and prevent dependency violations. The controller begins by stopping new connections to the database, then drains existing connections, performs final backups, and finally removes storage resources.

```
Deletion Processing Order:
1. Controller detects deletionTimestamp is set on Database resource
2. Controller updates status phase to "Terminating"
3. Controller initiates graceful shutdown of database pods via StatefulSet
4. Controller waits for active connections to drain (with timeout)
5. Controller triggers final backup job if backup schedule configured
6. Controller waits for backup completion or timeout
7. Controller deletes Service resources to prevent new connections
8. Controller deletes StatefulSet and waits for pod termination
9. Controller deletes ConfigMaps and Secrets
10. Controller deletes PersistentVolumeClaims (if policy allows)
11. Controller removes finalizer from Database resource
12. Kubernetes completes resource deletion automatically
```

Each cleanup step includes timeout handling to prevent indefinite blocking on failing operations. The controller maintains detailed status messages throughout the deletion process, helping users understand cleanup progress and identify any issues that require manual intervention.

#### Decision: PVC Deletion Policy

- **Context:** Database PersistentVolumeClaims contain critical data that users might want to preserve after operator deletion
- **Options Considered:**
  1. Always delete PVCs automatically during cleanup
  2. Never delete PVCs, leaving cleanup to administrators
  3. Configurable deletion policy per Database instance
- **Decision:** Configurable deletion policy with annotation-based control
- **Rationale:** Provides flexibility for both development environments (auto-cleanup) and production environments (manual cleanup)
- **Consequences:** Requires clear documentation about PVC lifecycle and potential storage costs from orphaned volumes

## Backup and Data Protection

Before deleting storage resources, the controller creates a final backup job to preserve database contents. This backup operation runs asynchronously while the controller monitors job status to determine when cleanup can proceed safely.

The backup job creation follows the same owner reference pattern used for other owned resources, ensuring that backup jobs are cleaned up automatically if the Database deletion process fails and requires restart. The controller includes backup metadata in the Database status, providing users with information about backup location and completion status.

| Backup Stage         | Duration Estimate   | Failure Handling                         | User Visibility                        |
|----------------------|---------------------|------------------------------------------|----------------------------------------|
| Job Creation         | < 1 second          | Retry with exponential backoff           | Status: "Creating backup job"          |
| Backup Execution     | Varies by data size | Monitor job status, timeout after 1 hour | Status: "Backing up data"              |
| Backup Verification  | < 30 seconds        | Check backup file existence and size     | Status: "Verifying backup"             |
| Cleanup Continuation | < 10 seconds        | Proceed with resource deletion           | Status: "Backup complete, cleaning up" |

## Cleanup Failure Recovery

The deletion flow includes comprehensive error handling for scenarios where cleanup operations fail due to temporary issues like network connectivity problems, insufficient permissions, or resource conflicts. The controller uses exponential backoff retry logic to handle transient failures while providing escape mechanisms for permanent failure conditions.

When cleanup operations encounter errors, the controller updates the Database status with detailed error information and requeues the reconciliation request for retry. The exponential backoff algorithm prevents overwhelming failing systems while ensuring that temporary issues resolve automatically once underlying conditions improve.

Common cleanup failures include PersistentVolumeClaims that cannot be deleted due to active pod mounts, backup jobs that fail due to storage connectivity issues, and StatefulSets that hang during termination due to pod disruption budgets or resource constraints.

### **⚠ Pitfall: Incomplete Finalizer Removal**

Controllers must remove their finalizers only after ALL cleanup operations complete successfully. Removing finalizers prematurely allows Kubernetes to complete resource deletion, potentially orphaning external resources or causing data loss. Always implement comprehensive error checking before finalizer removal, and prefer leaving finalizers in place during uncertain failure conditions rather than risking incomplete cleanup.

## Graceful Shutdown Coordination

The StatefulSet deletion process includes graceful shutdown coordination to ensure that database processes have sufficient time to flush pending writes, close connections cleanly, and perform other shutdown procedures. The controller configures appropriate `terminationGracePeriodSeconds` values based on the database type and expected shutdown duration.

During StatefulSet deletion, the controller monitors pod termination status and can extend grace periods if shutdown operations require additional time. This monitoring prevents forceful pod termination that could corrupt database files or lose pending transactions.

### Graceful Shutdown Monitoring:

1. Controller initiates StatefulSet deletion with configured grace period
2. Controller watches for pod termination events from informer
3. Controller checks database process status via readiness probes
4. If shutdown exceeds grace period, controller logs warning
5. Controller waits for all pods to reach Terminated phase
6. Controller verifies no orphaned persistent connections remain
7. Controller proceeds with next cleanup phase

The graceful shutdown process includes health checks that verify database processes have stopped completely before proceeding to delete storage resources. These checks prevent data corruption scenarios where storage cleanup occurs while database processes are still writing data files.

## Common Pitfalls

### ⚠️ Pitfall: Webhook Timeout During Resource Creation

Admission webhooks have strict timeout limits (typically 10-30 seconds) enforced by the API server. Complex validation logic or external API calls in webhooks can exceed these timeouts, causing resource creation to fail with cryptic timeout errors. Implement webhook logic with fast execution paths and avoid synchronous external dependencies. Use caching for expensive validations and prefer eventual consistency models where immediate validation is not critical.

### ⚠️ Pitfall: Status Update Race Conditions

Controllers that update status frequently can encounter race conditions where multiple status updates attempt to modify the same resource version simultaneously. This causes "conflict" errors that appear as reconciliation failures. Always fetch fresh resource versions before status updates, implement retry logic for conflict errors, and consider batching multiple status changes into single update operations.

### ⚠️ Pitfall: Finalizer Deadlock Scenarios

Multiple controllers adding finalizers to the same resource can create deadlock situations where each controller waits for others to complete cleanup before removing their own finalizers. Design finalizer cleanup logic to be independent of other finalizers and implement timeout mechanisms that allow emergency finalizer removal during administrative intervention.

### ⚠️ Pitfall: Informer Cache Staleness

Controllers making decisions based on informer cache data may operate on stale information that doesn't reflect recent cluster changes. Critical decisions like resource creation or deletion should fetch fresh data from the API server rather than relying solely on cache contents. Use `client.Get()` calls for authoritative state checks before irreversible operations.

### ⚠️ Pitfall: Webhook Certificate Expiration

Admission webhooks require valid TLS certificates for HTTPS communication with the API server. Certificate expiration causes immediate webhook failures that block all resource operations. Implement automated certificate rotation using cert-manager or similar tools, monitor certificate expiration dates, and establish emergency procedures for manual certificate renewal during outages.

## Implementation Guidance

### Technology Recommendations

| Component              | Simple Option                                                     | Advanced Option                                                            |
|------------------------|-------------------------------------------------------------------|----------------------------------------------------------------------------|
| HTTP Client            | <code>net/http</code> with retry logic                            | <code>controller-runtime/pkg/client</code> with built-in retries           |
| JSON Processing        | <code>encoding/json</code> with struct tags                       | <code>sigs.k8s.io/controller-runtime/pkg/webhook/admission</code> decoder  |
| Status Updates         | Direct client calls to status subresource                         | <code>controller-runtime/pkg/client</code> with <code>SubResource()</code> |
| Event Recording        | <code>client-go/tools/record</code><br><code>EventRecorder</code> | Controller-runtime manager event recorder                                  |
| Certificate Management | Manual certificate creation                                       | <code>cert-manager.io</code> automatic certificate lifecycle               |

## Recommended File Structure

```
internal/controller/
 database_controller.go ← Main reconciliation logic
 database_controller_test.go ← Reconciliation unit tests
 status_updater.go ← Status management utilities
 finalizer_manager.go ← Finalizer addition/removal logic

internal/webhook/
 database_webhook.go ← Admission webhook handlers
 database_webhook_test.go ← Webhook validation tests
 certificate_manager.go ← TLS certificate management
 admission_handler.go ← Webhook HTTP server setup

internal/resources/
 statefulset.go ← StatefulSet creation and updates
 service.go ← Service resource management
 configmap.go ← ConfigMap generation logic
 pvc.go ← PersistentVolumeClaim handling

config/webhook/
 manifests.yaml ← ValidatingAdmissionWebhook configuration
 kustomization.yaml ← Webhook resource overlay
 certificate.yaml ← Certificate resource for cert-manager
```

## Infrastructure Starter Code

Complete HTTP Server for Admission Webhooks:

```
package webhook
```

GO

```
import (
 "context"
 "crypto/tls"
 "fmt"
 "net/http"
 "time"

 admissionv1 "k8s.io/api/admission/v1"
 metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
 "k8s.io/apimachinery/pkg/runtime"
 "sigs.k8s.io/controller-runtime/pkg/log"
 "sigs.k8s.io/controller-runtime/pkg/webhook/admission"
)

// AdmissionHandler processes admission review requests for Database resources

type AdmissionHandler struct {
 decoder *admission.Decoder
 scheme *runtime.Scheme
}

func NewAdmissionHandler(decoder *admission.Decoder, scheme *runtime.Scheme) *AdmissionHandler {
 return &AdmissionHandler{
 decoder: decoder,
 scheme: scheme,
 }
}

func (h *AdmissionHandler) ServeHTTP(w http.ResponseWriter, r *http.Request) {
 logger := log.FromContext(r.Context())

 var review admissionv1.AdmissionReview
 if err := json.NewDecoder(r.Body).Decode(&review); err != nil {
 logger.Error(err, "Failed to decode admission review request")
 http.Error(w, "Invalid admission review", http.StatusBadRequest)
 }
}
```

```

 return

}

response := h.processAdmissionReview(r.Context(), &review)

review.Response = response

review.Response.UID = review.Request.UID

w.Header().Set("Content-Type", "application/json")

if err := json.NewEncoder(w).Encode(review); err != nil {

 logger.Error(err, "Failed to encode admission review response")

 http.Error(w, "Internal server error", http.StatusInternalServerError)

}

}

func (h *AdmissionHandler) processAdmissionReview(ctx context.Context, review *admissionv1.AdmissionReview) *admissionv1.AdmissionResponse {

 req := review.Request

 switch req.Kind.Kind {

 case "Database":

 switch req.Operation {

 case admissionv1.Create:

 return h.handleDatabaseCreate(ctx, req)

 case admissionv1.Update:

 return h.handleDatabaseUpdate(ctx, req)

 case admissionv1.Delete:

 return h.handleDatabaseDelete(ctx, req)

 }

 }

 return &admissionv1.AdmissionResponse{

 Allowed: true,

 Result: &metav1.Status{

 Code: http.StatusOK,

```

```
 },
}
}
```

Complete TLS Certificate Manager:

```
package webhook
```

GO

```
import (
 "context"
 "crypto/tls"
 "fmt"
 "os"
 "sync"
 "time"
)

// CertificateReloader manages TLS certificate lifecycle for webhook servers

type CertificateReloader struct {
 certPath string
 keyPath string
 cert *tls.Certificate
 mutex sync.RWMutex
}

func NewCertificateReloader(certPath, keyPath string) *CertificateReloader {
 return &CertificateReloader{
 certPath: certPath,
 keyPath: keyPath,
 }
}

func (cr *CertificateReloader) GetCertificate(*tls.ClientHelloInfo) (*tls.Certificate, error) {
 cr.mutex.RLock()
 defer cr.mutex.RUnlock()

 if cr.cert == nil {
 return nil, fmt.Errorf("certificate not loaded")
 }

 return cr.cert, nil
}
```

```
func (cr *CertificateReloader) loadCertificate() error {
 cert, err := tls.LoadX509KeyPair(cr.certPath, cr.keyPath)

 if err != nil {
 return fmt.Errorf("failed to load certificate: %w", err)
 }

 cr.mutex.Lock()
 cr.cert = &cert
 cr.mutex.Unlock()

 return nil
}

func (cr *CertificateReloader) Start(ctx context.Context) error {
 // Initial certificate load
 if err := cr.loadCertificate(); err != nil {
 return err
 }

 // Watch for certificate file changes
 ticker := time.NewTicker(30 * time.Second)
 defer ticker.Stop()

 for {
 select {
 case <-ctx.Done():
 return ctx.Err()
 case <-ticker.C:
 if err := cr.loadCertificate(); err != nil {
 // Log error but continue watching
 continue
 }
 }
 }
}
```

```
 }
}
```

### Core Logic Skeleton Code

Main Reconciliation Flow:

```
func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
 logger := log.FromContext(ctx)

 // TODO 1: Fetch the Database resource from the API server

 // Use r.Get() to retrieve the current resource state

 // Handle NotFound errors by returning successful result (resource deleted)

 // TODO 2: Check if resource is being deleted (deletionTimestamp set)

 // If deletion detected, call handleDeletion() and return

 // Deletion flow bypasses normal reconciliation logic

 // TODO 3: Add finalizer if not present

 // Check if FINALIZER_NAME exists in metadata.finalizers

 // Add finalizer and update resource if missing

 // TODO 4: Check if reconciliation is needed

 // Compare metadata.generation with status.observedGeneration

 // Skip reconciliation if generations match (no spec changes)

 // TODO 5: Fetch current state of owned resources

 // Get StatefulSet, Service, ConfigMap, PVC owned by this Database

 // Use owner reference labels for efficient resource discovery

 // TODO 6: Compute desired state from spec

 // Generate desired StatefulSet, Service, ConfigMap based on Database spec

 // Apply configuration templates with values from spec fields

 // TODO 7: Reconcile each owned resource type

 // Call createOrUpdateResource() for StatefulSet, Service, ConfigMap

 // Handle creation, updates, and error scenarios

 // TODO 8: Wait for StatefulSet readiness

 // Check StatefulSet status.readyReplicas vs spec.replicas

 // Return requeue if StatefulSet not ready yet
```

```
// TODO 9: Update Database status

// Set conditions, readyReplicas, phase based on owned resource status

// Update observedGeneration to current metadata.generation

// TODO 10: Return appropriate result

// Return ctrl.Result{} for success, ctrl.Result{RequeueAfter: duration} for retry

// Return error for permanent failures that need exponential backoff

}
```

**Database Validation Logic:**

```
func (h *AdmissionHandler) handleDatabaseCreate(ctx context.Context, req *admissionv1.AdmissionRequest) *admissionv1.AdmissionResponse {
 var db Database

 // TODO 1: Decode the Database resource from admission request
 // Use h.decoder.DecodeRaw() to convert req.Object into Database struct
 // Handle JSON decode errors with appropriate error response

 // TODO 2: Validate replica count constraints
 // Check if db.Spec.Replicas is between 1 and maximum allowed (e.g., 5)
 // Return denial response for invalid replica counts

 // TODO 3: Validate storage size format and minimum
 // Parse db.Spec.StorageSize using resource.ParseQuantity()
 // Ensure storage size meets minimum requirements (e.g., 1Gi)

 // TODO 4: Validate PostgreSQL version format
 // Check db.Spec.Version matches expected version pattern (major.minor)
 // Verify version is in supported versions list

 // TODO 5: Validate resource requirements
 // Check db.Spec.Resources.Limits and Requests for reasonable values
 // Ensure memory limits are sufficient for database operation

 // TODO 6: Validate configuration parameters
 // Iterate through db.Spec.Config map entries
 // Check parameter names against allowed PostgreSQL settings
 // Validate parameter values are within acceptable ranges

 // TODO 7: Return admission response
 // Return &admissionv1.AdmissionResponse{Allowed: true} for valid resources
 // Return denial response with detailed error message for validation failures
}
```

#### Status Update Logic:

```
func (r *DatabaseReconciler) updateDatabaseStatus(ctx context.Context, db *Database, sts *appsv1.StatefulSet) error {
 // TODO 1: Create updated status object

 // Initialize new DatabaseStatus struct with current field values

 // Preserve existing status fields that should not change

 // TODO 2: Update ready replica count

 // Set status.ReadyReplicas from sts.Status.ReadyReplicas

 // Handle case where StatefulSet is nil (not created yet)

 // TODO 3: Determine current phase

 // Set status.Phase based on StatefulSet status and readiness

 // Use phases: "Pending", "Progressing", "Ready", "Failed"

 // TODO 4: Update Ready condition

 // Set Ready condition to True if readyReplicas == spec.replicas

 // Set Ready condition to False if StatefulSet has errors

 // Set Ready condition to Unknown during initial deployment

 // TODO 5: Update Progressing condition

 // Set Progressing to True during rollouts or scaling operations

 // Set Progressing to False when deployment reaches stable state

 // Include reason and message describing current operation

 // TODO 6: Set observed generation

 // Update status.ObservedGeneration to current metadata.generation

 // This signals that reconciliation has processed latest spec changes

 // TODO 7: Update last reconciliation timestamp

 // Set status.LastReconcileTime to current time

 // Provides visibility into reconciliation activity

 // TODO 8: Persist status update

 // Use r.Status().Update() to save status changes
```

```
// Handle conflict errors with retry logic
}
```

## Milestone Checkpoint

### After Milestone 3 Implementation:

Run the following commands to verify complete resource lifecycle:

```
Test resource creation flow

kubectl apply -f config/samples/database_v1_database.yaml

kubectl get database sample-database -o yaml

Verify webhook processing

kubectl describe database sample-database

Should show defaulted values and validation success

Test update and reconciliation flow

kubectl patch database sample-database --type='merge' -p='{"spec":{"replicas":3}}'

kubectl get database sample-database -o jsonpath='{.status.readyReplicas}'

Test deletion and cleanup flow

kubectl delete database sample-database

kubectl get statefulset,service,pvc

Should show progressive cleanup of owned resources
```

### Expected Behavior:

- Resource creation triggers mutating webhook defaulting followed by validating webhook checks
- Controller reconciliation creates StatefulSet, Service, and ConfigMap with owner references
- Status updates show progression from Pending → Progressing → Ready phases
- Replica count changes trigger StatefulSet updates and status reflects new ready count
- Resource deletion initiates finalizer-controlled cleanup sequence
- All owned resources are removed before Database deletion completes

### Signs of Issues:

- Webhook timeouts indicate certificate or networking problems
- Reconciliation stuck in Progressing phase suggests StatefulSet creation failures
- Status never updates indicates RBAC permission issues
- Finalizer not removed indicates cleanup operation failures

## Debugging Tips

| Symptom                                  | Likely Cause                       | How to Diagnose                                          | Fix                                                         |
|------------------------------------------|------------------------------------|----------------------------------------------------------|-------------------------------------------------------------|
| Resource creation hangs                  | Webhook timeout                    | Check webhook pod logs and certificate validity          | Restart webhook pods, renew certificates                    |
| Validation rejected with unclear message | Webhook validation logic error     | Examine admission response details in kubectl output     | Review webhook validation code, add detailed error messages |
| Controller reconciliation never triggers | Informer cache sync failure        | Check controller manager logs for watch errors           | Verify RBAC permissions, restart controller                 |
| Status updates fail with conflict errors | Concurrent status modifications    | Enable verbose logging to see resource version conflicts | Implement retry logic with fresh resource fetches           |
| Finalizer cleanup stuck                  | External resource deletion failure | Check Database status message and owned resource states  | Manually clean up stuck resources, remove finalizer         |
| Owned resources not deleted              | Missing owner references           | Verify owner reference metadata on StatefulSet, Service  | Add proper owner references during resource creation        |

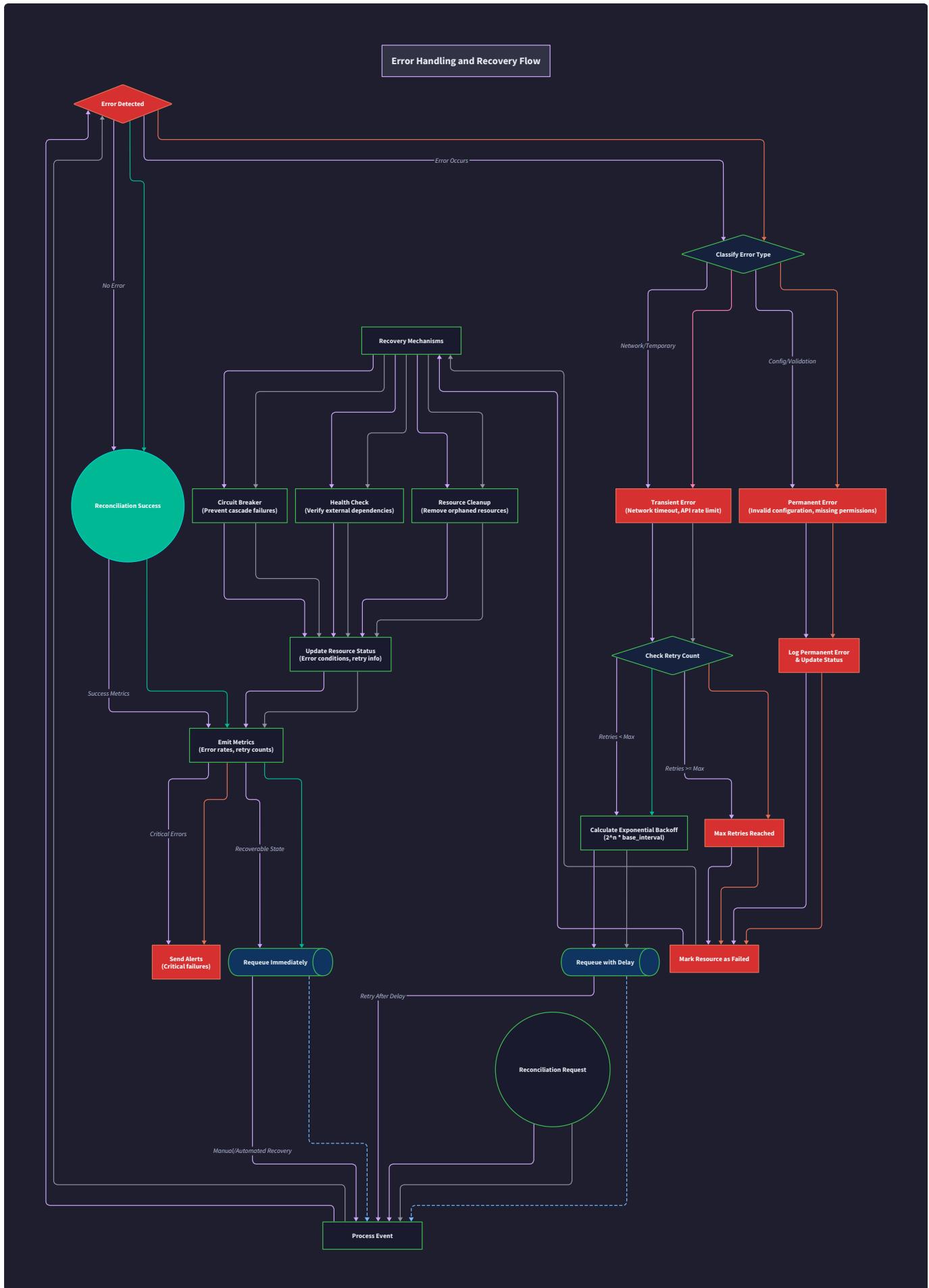
## Error Handling and Edge Cases

**Milestone(s):** Milestone 3 (Reconciliation Loop), Milestone 4 (Webhooks), Milestone 5 (Testing & Deployment) - addresses robust error handling across controller reconciliation, webhook processing, and production deployment scenarios

### Mental Model: The Resilient Orchestra Conductor

Think of a Kubernetes operator as an orchestra conductor who must keep the music playing even when instruments malfunction, sheet music gets lost, or the concert hall loses power. A resilient conductor doesn't panic when the violin section misses an entrance—they calmly adjust the tempo, signal the missed cue again, and ensure the overall performance continues. Similarly, our operator must gracefully handle API server outages, webhook timeouts, and resource conflicts while maintaining the declarative promises made to users.

The key insight is that distributed systems failures are not exceptional—they are normal operating conditions. Just as a professional conductor expects and prepares for missed cues, a well-designed operator anticipates network partitions, transient errors, and partial failures as part of its standard repertoire.



## Common Failure Modes

Understanding failure modes is crucial for building resilient operators. Each failure mode requires different detection, classification, and recovery strategies. The operator must distinguish between transient errors that will resolve themselves with retry and permanent errors that require human intervention or alternative approaches.

### API Server Connectivity Failures

API server connectivity issues are among the most common failures in Kubernetes environments. These can range from brief network hiccups to extended outages during cluster maintenance or infrastructure failures.

#### Connection Loss During Informer Operations

When the informer loses its connection to the API server, it stops receiving resource events and its local cache becomes stale. The controller might continue processing items from its work queue, but it cannot fetch current resource state or update resource status. This creates a dangerous situation where the controller operates on outdated information.

| Failure Scenario          | Symptoms                        | Detection Method                         | Recovery Strategy                                   |
|---------------------------|---------------------------------|------------------------------------------|-----------------------------------------------------|
| Network partition         | Informer stops receiving events | Connection error logs, cache staleness   | Exponential backoff reconnection with jitter        |
| API server restart        | HTTP 503 responses              | Failed API calls with server unavailable | Wait for server readiness, rebuild cache            |
| Certificate expiration    | TLS handshake failures          | Certificate validation errors            | Automatic certificate reload or rotation            |
| Resource quota exhaustion | API calls rejected              | HTTP 429 or quota exceeded errors        | Back off and alert, may require manual intervention |

The informer's shared index cache provides some protection during brief outages, allowing read operations to continue using cached data. However, the controller must be aware that cached data may be stale and avoid making destructive decisions based on potentially outdated information.

#### Decision: Informer Cache Staleness Handling

- **Context:** When API server connectivity is lost, the informer cache becomes stale but remains readable
- **Options Considered:**
  1. Stop all reconciliation during connectivity loss
  2. Continue reconciliation with stale cache data
  3. Limit reconciliation to read-only operations during outages
- **Decision:** Continue reconciliation with stale cache but avoid destructive operations
- **Rationale:** Stopping completely would prevent recovery when connectivity returns, but destructive operations on stale data could cause cascading failures
- **Consequences:** Requires careful classification of operations as safe/unsafe with stale data

### API Call Failures During Reconciliation

Individual API calls within the reconciliation loop can fail for various reasons, from temporary network issues to resource conflicts with concurrent operations. The controller must handle these failures gracefully without corrupting the overall reconciliation state.

#### Reconciliation Error Classification:

1. Check if error indicates stale cache (NotFound for recently created resource)
2. Classify as transient (network timeout) or permanent (validation failure)
3. For transient errors: requeue with exponential backoff
4. For permanent errors: update status with error condition, stop retrying
5. For cache staleness: force cache refresh and retry once

## Webhook Processing Failures

Admission webhooks introduce additional failure points in the resource creation and update pipeline. Webhook failures can block all resource operations, making them particularly critical to handle correctly.

### Webhook Timeout and Unavailability

The Kubernetes API server has a default 10-second timeout for webhook calls. If the webhook service is unavailable or overloaded, these timeouts can block all operations on the affected resource types.

| Timeout Scenario       | Impact                       | Detection                             | Mitigation                                         |
|------------------------|------------------------------|---------------------------------------|----------------------------------------------------|
| Webhook pod not ready  | All resource operations fail | Pod status, readiness probes          | Failure policy configuration, multiple replicas    |
| Network latency spikes | Intermittent timeouts        | Response time monitoring              | Reduce webhook processing time, increase timeout   |
| Webhook overloaded     | Cascading timeout failures   | Queue depth metrics, CPU/memory usage | Horizontal scaling, request rate limiting          |
| Certificate issues     | TLS handshake failures       | Certificate validation logs           | Automated certificate management with cert-manager |

The webhook configuration's `failurePolicy` setting determines whether operations should be allowed or denied when the webhook is unreachable. This represents a fundamental trade-off between availability and security enforcement.

### Decision: Webhook Failure Policy Configuration

- **Context:** Webhook unavailability can block all resource operations
- **Options Considered:**
  1. `failurePolicy: Fail` - deny operations when webhook unavailable (secure but less available)
  2. `failurePolicy: Ignore` - allow operations when webhook unavailable (available but less secure)
  3. Hybrid approach with different policies for validation vs mutation
- **Decision:** Use `Fail` for critical validations, `Ignore` for optional mutations
- **Rationale:** Critical business rules must be enforced even at the cost of availability, but optional enhancements shouldn't block operations
- **Consequences:** Requires careful classification of webhook logic as critical vs optional

## Webhook Processing Errors

Even when the webhook service is reachable, processing errors within the webhook logic can cause admission failures. These errors must be handled gracefully with appropriate error messages for users.

The webhook's error handling affects user experience significantly. A webhook that returns a generic "internal error" message leaves users unable to diagnose or fix their resource definitions, while clear, actionable error messages enable self-service problem resolution.

## Controller Process Failures

Controller crashes and restarts are inevitable in distributed environments. The controller must be designed to recover gracefully from failures without losing track of resources or corrupting reconciliation state.

### Controller Crash Recovery

When a controller process crashes and restarts, it loses all in-memory state including work queue contents and ongoing reconciliation progress. The controller must reconstruct its understanding of the world from the persistent state stored in the Kubernetes API server.

| Recovery Aspect         | Challenge                              | Solution                                          |
|-------------------------|----------------------------------------|---------------------------------------------------|
| Work queue state        | In-memory queue contents lost on crash | Rebuild from informer full sync on startup        |
| Ongoing reconciliations | Partial operations may be incomplete   | Implement idempotent reconciliation logic         |
| Finalizer cleanup       | Resources may be stuck in deletion     | Check for stale finalizers on startup scan        |
| Status updates          | Status may not reflect current state   | Recompute and update status during reconciliation |

The informer's resync mechanism provides a safety net by periodically re-enqueueing all resources, ensuring that any missed events or incomplete reconciliations are eventually processed.

### Leader Election Failures

In high-availability deployments with multiple controller replicas, leader election ensures only one replica actively reconciles resources. Leader election failures can result in split-brain scenarios or periods with no active leader.

Leader election uses Kubernetes lease resources with time-based renewal. The leader must continuously renew its lease, and followers must monitor for lease expiration to detect when they should attempt to acquire leadership.

## Resource Conflict Scenarios

Multiple controllers or external actors can modify the same Kubernetes resources concurrently, leading to conflicts that must be resolved without data loss or inconsistent state.

### Optimistic Locking Conflicts

Kubernetes uses optimistic locking through resource versions to detect concurrent modifications. When a controller attempts to update a resource that has been modified since it was last read, the API server returns a conflict error.

- ```
Optimistic Locking Conflict Resolution:  
1. Attempt resource update with current resource version  
2. If conflict error occurs, re-read resource from API server  
3. Merge local changes with current state (if possible)  
4. Retry update with new resource version  
5. If conflicts persist, implement exponential backoff  
6. Alert if conflicts indicate systematic issues
```

The controller's reconciliation logic must be designed to handle resource version conflicts gracefully by re-reading current state and recomputing desired changes.

Concurrent Controller Modifications

When multiple controllers manage different aspects of the same resource, coordination becomes essential to prevent conflicts and ensure coherent behavior.

Conflict Type	Example	Resolution Strategy
Field ownership	Two controllers updating same field	Use Server-Side Apply with field management
Resource creation	Multiple controllers creating same resource	Use owner references and naming conventions
Status updates	Concurrent status field modifications	Use status subresource with separate resource versions
Finalizer management	Multiple controllers adding/removing finalizers	Coordinate finalizer names and removal order

Recovery and Self-Healing

The operator's recovery mechanisms must be automatic, reliable, and comprehensive enough to handle the full spectrum of failure scenarios without human intervention for common cases.

Automatic Retry Logic

The retry strategy forms the foundation of the operator's resilience. Different types of errors require different retry approaches, and the strategy must prevent overwhelming failing systems while ensuring rapid recovery when conditions improve.

Exponential Backoff Implementation

Exponential backoff prevents the controller from overwhelming a struggling API server or webhook service with repeated failed requests. The backoff algorithm increases delay between retries exponentially, with jitter to prevent thundering herd problems.

Retry Attempt	Base Delay	Jitter Range	Total Delay Range
1	1s	±0.5s	0.5s - 1.5s
2	2s	±1s	1s - 3s
3	4s	±2s	2s - 6s
4	8s	±4s	4s - 12s
5	16s	±8s	8s - 24s
Max	300s	±150s	150s - 450s

The controller-runtime's rate-limiting work queue provides built-in exponential backoff with configurable parameters. The `isErrorRetryable` function classifies errors to determine whether retry is appropriate.

Error Classification for Retry Decisions

Not all errors should trigger automatic retries. Permanent errors like validation failures will not resolve through retry and should be reported to users immediately rather than consuming retry capacity.

```
Error Classification Algorithm:
1. Check error type and HTTP status code
2. Transient errors (network timeouts, 5xx responses): retry with backoff
3. Retryable client errors (409 conflicts, 429 rate limits): retry with backoff
4. Permanent client errors (400 validation, 403 permission): report immediately
5. Resource not found errors: check if resource was recently deleted, may retry once
6. Unknown errors: default to transient classification for safety
```

The classification function must be conservative, erring on the side of retrying potentially recoverable errors rather than giving up prematurely.

Leader Election Failover

High-availability deployments require seamless failover when the current leader becomes unavailable. The leader election mechanism must detect failures quickly while preventing split-brain scenarios.

Leader Health Monitoring

The active leader must continuously prove its health by renewing its lease before expiration. The lease renewal interval must be shorter than the lease duration to provide a safety margin for network delays and processing time.

Configuration Parameter	Value	Rationale
Lease Duration	15s	Long enough to survive brief network hiccups
Renew Deadline	10s	Provides 5s safety margin before lease expiration
Retry Period	2s	Allows multiple renewal attempts within deadline

Followers monitor the leader's lease and attempt to acquire leadership when the lease expires without renewal. The acquisition process uses atomic compare-and-swap operations to prevent multiple followers from becoming leader simultaneously.

Graceful Leadership Transition

When leadership changes, the new leader must reconstruct the controller state and resume reconciliation without missing resources or duplicating operations. This requires careful coordination between the outgoing and incoming leaders.

The incoming leader performs a full informer resync to rebuild its local cache and work queue, ensuring it processes all resources regardless of what the previous leader may have missed. This approach trades some performance for reliability and simplicity.

State Reconstruction from Cluster State

After failures or leadership changes, the controller must reconstruct its understanding of the world from the authoritative state stored in the Kubernetes API server. This reconstruction must be complete and accurate to prevent inconsistencies.

Full Reconciliation Sweep

The informer's initial cache sync triggers reconciliation for all existing custom resources, allowing the controller to verify that actual state matches desired state and correct any drift that may have occurred during downtime.

State Reconstruction Process:

1. Start informer and wait for initial cache sync completion
2. Controller receives reconcile requests for all existing resources
3. Each reconciliation compares actual state with desired state
4. Controller corrects any drift by creating, updating, or deleting owned resources
5. Status is updated to reflect current reconciliation state
6. Normal event-driven reconciliation resumes

This approach ensures that the controller can recover from extended downtime, configuration changes, or external modifications to owned resources.

Finalizer Cleanup Verification

During startup, the controller must verify that all finalizers under its management are still needed and remove any stale finalizers that may prevent resource deletion. This prevents resources from becoming stuck in deletion state due to controller failures or configuration changes.

The controller scans all resources with its finalizer and verifies that cleanup operations are complete. If cleanup has already been performed (perhaps by a previous controller instance), the finalizer is removed to allow deletion to proceed.

Consistency and Conflict Resolution

Distributed systems like Kubernetes provide eventual consistency rather than strong consistency, requiring operators to handle concurrent modifications and conflicting updates gracefully.

Handling Concurrent Updates

Multiple actors may attempt to modify the same Kubernetes resource simultaneously, including the operator itself, other controllers, and human administrators using `kubectl`. The operator must coordinate these updates without losing data or creating inconsistent state.

Resource Version Conflicts

Every Kubernetes resource includes a `resourceVersion` field that changes with each modification. When updating a resource, the client must provide the current resource version, and the API server rejects updates that specify an outdated version.

Conflict Resolution Step	Action	Purpose
1. Read current resource	<code>GET /api/v1/namespaces/{ns}/databases/{name}</code>	Obtain latest resourceVersion
2. Compute desired changes	Compare spec with actual state	Determine required modifications
3. Apply changes	<code>PUT</code> or <code>PATCH</code> with resourceVersion	Attempt atomic update
4. Handle conflict	If 409 Conflict, retry from step 1	Resolve concurrent modifications
5. Update work queue	Requeue with exponential backoff if needed	Prevent overwhelming API server

The controller's reconciliation logic must be prepared to restart the reconciliation process when conflicts occur, ensuring that it operates on current resource state rather than stale cached data.

Server-Side Apply for Field Management

Server-Side Apply provides a more sophisticated approach to handling concurrent updates by tracking field ownership and allowing controllers to declare which fields they manage. This prevents accidental overwrites of fields managed by other controllers or users.

Each controller declares its identity as a field manager and specifies exactly which fields it wants to control. The API server automatically resolves conflicts by preserving fields owned by other managers while applying only the requesting controller's managed fields.

Decision: Update Strategy Selection

- **Context:** Controllers must handle concurrent updates to shared resources
- **Options Considered:**
 1. Traditional Update (`PUT`) with conflict retry - simple but can overwrite other controllers' changes
 2. Strategic Merge Patch - better than `PUT` but still can conflict
 3. JSON Patch - precise but fragile to concurrent changes
 4. Server-Side Apply - most sophisticated but requires field management design
- **Decision:** Use Server-Side Apply for spec updates, traditional Update for status
- **Rationale:** Spec updates often conflict with user or other controller changes, while status is typically owned by single controller
- **Consequences:** Requires defining field management policies and handling field ownership conflicts

Optimistic Locking Patterns

Optimistic locking assumes that conflicts are rare and allows multiple readers to access data concurrently, detecting conflicts only when writes occur. This provides better performance than pessimistic locking but requires careful conflict resolution logic.

Read-Modify-Write Cycles

The controller's typical update pattern follows a read-modify-write cycle where it reads the current resource state, computes required changes, and attempts to write the updated resource back to the API server.

```
Optimistic Locking Pattern:  
1. Read resource and note resourceVersion (e.g., "12345")  
2. Compute desired changes based on current state  
3. Modify resource object with intended changes  
4. Attempt update with resourceVersion: "12345"  
5. If successful: operation complete  
6. If conflict (409): another actor modified resource, restart from step 1  
7. If other error: classify and handle according to retry policy
```

The key insight is that the controller must be prepared to discard its computed changes and restart the reconciliation process with fresh data when conflicts occur.

Idempotent Reconciliation Design

Idempotent reconciliation ensures that applying the same desired state multiple times produces the same result, making conflict resolution and retry logic safe and predictable.

Idempotency Principle	Example	Implementation
Create operations are idempotent	Creating a Deployment	Check if resource exists before creating, or use <code>kubectl apply</code> semantics
Update operations are idempotent	Scaling replicas from 2 to 3	Set absolute values rather than relative changes
Delete operations are idempotent	Removing a Service	Check if resource exists before deletion, handle NotFound gracefully
Status updates are idempotent	Setting Ready condition	Compare current status, update only if different

Idempotent design allows the controller to safely retry operations after conflicts without worrying about double-application of changes or accumulating errors.

Eventual Consistency Expectations

Kubernetes provides eventual consistency, meaning that changes propagate through the system over time but may not be immediately visible to all observers. The operator must account for this consistency model in its design and expectations.

Informer Cache Lag

The informer's local cache may lag behind the authoritative state in the API server, particularly during periods of high change volume or network issues. The controller must handle cases where its cache doesn't yet reflect recent changes made by other actors.

A common scenario occurs when the controller creates a resource (like a Deployment) and immediately queries its informer cache to check the resource's status. The cache may not yet contain the newly created resource, leading to incorrect decisions if the controller assumes the resource doesn't exist.

The critical insight here is that informer caches provide **eventually consistent** views of cluster state, not **immediately consistent** views. Controllers must design their logic to handle cache lag gracefully rather than assuming instant consistency.

Cross-Resource Dependency Timing

When the operator manages multiple related resources (like a Deployment, Service, and ConfigMap), the timing of their creation and readiness can vary. The controller must handle scenarios where dependencies aren't yet available or haven't reached their desired state.

Dependency Handling Strategy:

1. Create all owned resources in dependency order (ConfigMap → Deployment → Service)
2. Check readiness of dependencies before proceeding with dependent resources
3. Use owner references to ensure garbage collection if parent resource is deleted
4. Implement timeout and retry logic for dependency readiness checks
5. Update status with detailed progress information for troubleshooting

The `updateCondition` function maintains detailed status information about dependency progress, allowing users and debugging tools to understand why reconciliation may be incomplete.

Common Pitfalls

⚠️ Pitfall: Assuming Immediate Consistency After Creates Many developers assume that after successfully creating a resource via the API server, it will immediately be visible in the informer cache. This leads to race conditions where the controller makes incorrect decisions based on stale cache data. The fix is to either use direct API server queries for recently created resources or implement retry logic that handles temporary cache inconsistency.

⚠️ Pitfall: Not Implementing Proper Resource Version Handling Some controllers ignore resource version conflicts and keep retrying with the original resource version, leading to permanent update failures. The correct approach is to re-read the resource after each conflict and recompute changes based on the current state rather than the original state.

⚠️ Pitfall: Blocking Reconciliation on Transient Errors Controllers that fail fast on any error, including transient network issues, provide poor user experience. The operator should distinguish between permanent errors (validation failures) that require user intervention and transient errors (network timeouts) that should be retried automatically with exponential backoff.

⚠️ Pitfall: Missing Finalizer Cleanup in Error Scenarios If the controller adds a finalizer but crashes before completing cleanup, the resource becomes stuck in deletion state. The controller must implement startup logic that scans for stale finalizers and completes any pending cleanup operations to prevent resources from being permanently stuck.

⚠️ Pitfall: Not Handling Webhook Unavailability Gracefully Webhook services that crash or become unavailable can block all operations on their resource types if configured with `failurePolicy: Fail`. While this prevents bypassing validation, it can make the entire system unusable. The solution is to implement webhook high availability with multiple replicas, health checks, and careful consideration of failure policy settings.

Implementation Guidance

Understanding error handling concepts is only half the battle—implementing robust error handling requires careful attention to Go's error handling patterns, controller-runtime's retry mechanisms, and Kubernetes API conventions.

Technology Recommendations

Component	Simple Option	Advanced Option
Error Classification	Basic type assertions and string matching	Custom error types with structured classification
Retry Logic	controller-runtime's default rate limiter	Custom rate limiter with circuit breaker patterns
Monitoring	Standard Go logging with structured fields	Prometheus metrics with custom error counters
Leader Election	controller-runtime's built-in leader election	Custom leader election with health checks

Recommended File Structure

```
internal/controller/
    database_controller.go      ← main reconciliation logic with error handling
    database_controller_test.go ← tests including error scenarios
    error_handling.go          ← error classification and retry logic
    finalizer_cleanup.go       ← startup cleanup and recovery logic
    status_manager.go          ← status update logic with conflict handling
internal/webhook/
    admission_handler.go       ← Webhook error handling and timeout management
    tls_manager.go              ← certificate management and rotation
internal/metrics/
    error_metrics.go           ← error tracking and alerting
```

Error Classification Infrastructure

This complete error classification system provides the foundation for intelligent retry decisions throughout the operator:

```
package controller

import (
    "errors"
    "net"
    "time"

    apierrors "k8s.io/apimachinery/pkg/api/errors"
    "sigs.k8s.io/controller-runtime/pkg/reconcile"
)

// ErrorClassification represents different types of errors and their handling

type ErrorClassification int

const (
    ErrorTransient ErrorClassification = iota // Retry with exponential backoff
    ErrorRetryable                      // Retry immediately or with short delay
    ErrorPermanent                       // Don't retry, update status with error
    ErrorConflict                        // Special handling for optimistic locking
)

// ClassifyError determines how an error should be handled during reconciliation

func (r *DatabaseReconciler) ClassifyError(err error) ErrorClassification {
    if err == nil {
        return ErrorTransient // Should not happen, but safe default
    }

    // Kubernetes API errors

    if apierrors.IsConflict(err) {
        return ErrorConflict
    }

    if apierrors.NotFound(err) {
        return ErrorRetryable // May be cache lag
    }

    if apierrors.IsServerTimeout(err) || apierrors.IsTimeout(err) {
        return ErrorTransient
    }
}
```

GO

```
}

if apierrors.IsTooManyRequests(err) {

    return ErrorTransient

}

if apierrors.IsInternalError(err) || apierrors.IsServiceUnavailable(err) {

    return ErrorTransient

}

if apierrors.IsBadRequest(err) || apierrors.IsValid(err) {

    return ErrorPermanent

}

if apierrors.IsForbidden(err) || apierrors.IsUnauthorized(err) {

    return ErrorPermanent

}

// Network errors

var netErr net.Error

if errors.As(err, &netErr) {

    if netErr.Timeout() {

        return ErrorTransient

    }

    return ErrorTransient // Most network errors are transient

}

// Default to transient for safety

return ErrorTransient

}

// HandleReconcileError processes errors according to their classification

func (r *DatabaseReconciler) HandleReconcileError(ctx context.Context,

    db *Database, err error) (reconcile.Result, error) {

    classification := r.ClassifyError(err)

    switch classification {
```

```

case ErrorTransient:

    // Let controller-runtime handle exponential backoff

    r.updateCondition(db, ConditionTypeReady, "False", "TransientError", err.Error())

    return reconcile.Result{}, err


case ErrorRetryable:

    // Quick retry, might be cache staleness

    r.updateCondition(db, ConditionTypeReady, "False", "RetryableError", err.Error())

    return reconcile.Result{RequeueAfter: 5 * time.Second}, nil


case ErrorConflict:

    // Resource version conflict - retry immediately

    return reconcile.Result{Requeue: true}, nil


case ErrorPermanent:

    // Don't retry, update status with error

    r.updateCondition(db, ConditionTypeReady, "False", "PermanentError", err.Error())

    return reconcile.Result{}, nil // Don't return error to avoid retry


default:

    return reconcile.Result{}, err
}
}

```

Robust Reconciliation Loop Skeleton

This skeleton implements the complete error handling patterns discussed in the design section:

```
func (r *DatabaseReconciler) Reconcile(ctx context.Context, req reconcile.Request) (reconcile.Result, error) {    GO

    logger := r.Log.WithValues("database", req.NamespacedName)

    // TODO 1: Fetch the Database resource with error handling

    // Hint: Handle NotFound errors gracefully - resource may have been deleted

    var db Database

    if err := r.Get(ctx, req.NamespacedName, &db); err != nil {

        if apierrors.IsNotFound(err) {

            // Resource deleted, stop reconciliation

            return reconcile.Result{}, nil

        }

        return r.HandleReconcileError(ctx, &db, err)
    }

    // TODO 2: Handle finalizer logic for deletion

    // Hint: Check if deletion timestamp is set, add/remove finalizers accordingly

    if !db.DeletionTimestamp.IsZero() {

        return r.handleDeletion(ctx, &db)
    }

    if !controllerutil.ContainsFinalizer(&db, FINALIZER_NAME) {

        controllerutil.AddFinalizer(&db, FINALIZER_NAME)

        return reconcile.Result{Requeue: true}, r.Update(ctx, &db)
    }

    // TODO 3: Get current state of owned resources with conflict handling

    // Hint: Use Get() calls with proper error classification

    currentState, err := r.getCurrentState(ctx, &db)

    if err != nil {

        return r.HandleReconcileError(ctx, &db, err)
    }

    // TODO 4: Compare desired state from spec with current state

    // Hint: Implement comparison logic that's resilient to cache staleness

    desiredState := r.computeDesiredState(&db)
```

```
changes := r.computeChanges(currentState, desiredState)

// TODO 5: Apply changes with proper error handling and conflict resolution

// Hint: Use Server-Side Apply for spec updates, handle conflicts gracefully

for _, change := range changes {

    if err := r.applyChange(ctx, change); err != nil {

        if r.ClassifyError(err) == ErrorConflict {

            // Immediate retry for conflicts

            return reconcile.Result{Requeue: true}, nil
        }
    }

    return r.HandleReconcileError(ctx, &db, err)
}

// TODO 6: Update status with current conditions and state

// Hint: Use status subresource, handle conflicts by re-reading and retrying

return r.updateDatabaseStatus(ctx, &db, currentState)
}
```

Webhook Error Handling Infrastructure

Complete webhook error handling that addresses timeout and certificate issues:

```
package webhook
```

GO

```
import (
    "context"
    "crypto/tls"
    "fmt"
    "net/http"
    "time"

    admissionv1 "k8s.io/api/admission/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/runtime"
)

type AdmissionHandler struct {
    decoder *admission.Decoder
    certReloader *CertificateReloader
}

// ServeHTTP handles admission review requests with comprehensive error handling
func (h *AdmissionHandler) ServeHTTP(w http.ResponseWriter, r *http.Request) {
    // TODO 1: Set reasonable timeout for webhook processing
    // Hint: Kubernetes has 10s default timeout, leave buffer for response
    ctx, cancel := context.WithTimeout(r.Context(), 8*time.Second)
    defer cancel()

    // TODO 2: Parse admission review request with error handling
    // Hint: Handle malformed requests gracefully with proper HTTP status codes
    review, err := h.parseAdmissionReview(r)
    if err != nil {
        h.sendErrorResponse(w, http.StatusBadRequest,
            fmt.Sprintf("Failed to parse admission review: %v", err))
        return
    }
}
```

```

// TODO 3: Process admission review with timeout handling

// Hint: Check context deadline and return partial response if needed

response := h.processAdmissionReview(ctx, review)

// TODO 4: Send response with proper error handling

// Hint: Set appropriate HTTP status codes and content types

h.sendAdmissionResponse(w, response)

}

// processAdmissionReview routes requests to appropriate handlers with error recovery

func (h *AdmissionHandler) processAdmissionReview(ctx context.Context,
    review *admissionv1.AdmissionReview) *admissionv1.AdmissionReview {

    // TODO 1: Check context deadline before processing

    // Hint: Return admission denial if insufficient time remains

    if deadline, ok := ctx.Deadline(); ok && time.Until(deadline) < 1*time.Second {

        return h.timeoutResponse(review, "Webhook processing timeout")

    }

    // TODO 2: Route to appropriate handler based on operation and kind

    // Hint: Use type switches and handle unknown operations gracefully

    req := review.Request

    var response *admissionv1.AdmissionResponse

    switch {

    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Create:

        response = h.handleDatabaseCreate(ctx, req)

    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Update:

        response = h.handleDatabaseUpdate(ctx, req)

    case req.Kind.Kind == "Database" && req.Operation == admissionv1.Delete:

        response = h.handleDatabaseDelete(ctx, req)

    default:

        response = h.allowResponse(req.UID, "Operation not handled by this webhook")

    }
}

```

```
// TODO 3: Create admission review response with proper metadata

// Hint: Preserve request UID and API version information

return &admissionv1.AdmissionReview{

    TypeMeta: metav1.TypeMeta{

        APIVersion: "admission.k8s.io/v1",

        Kind:      "AdmissionReview",

    },

    Response: response,

}

}
```

TLS Certificate Management

Robust certificate handling that prevents webhook outages due to certificate issues:

```
type CertificateReloader struct {

    certPath string

    keyPath string

    cert     *tls.Certificate

    certTime time.Time

    keyTime  time.Time

}

func NewCertificateReloader(certPath, keyPath string) *CertificateReloader {

    return &CertificateReloader{

        certPath: certPath,

        keyPath: keyPath,

    }

}

// GetCertificate returns current certificate, reloading if files have changed

func (cr *CertificateReloader) GetCertificate(*tls.ClientHelloInfo) (*tls.Certificate, error) {

    // TODO 1: Check if certificate files have been modified

    // Hint: Use os.Stat to check modification times, reload if changed

    // TODO 2: Load certificate with proper error handling

    // Hint: Don't update stored cert if loading fails, return previous cert

    // TODO 3: Validate certificate expiration

    // Hint: Log warnings if certificate expires soon, return error if expired

    return cr.cert, nil

}
```

GO

Milestone Checkpoints

Checkpoint 1: Error Classification Testing Run this test to verify error classification works correctly:

```
go test ./internal/controller -run TestErrorClassification
```

BASH

Expected behavior: All Kubernetes API error types are correctly classified as transient, retryable, or permanent. Network errors are classified as transient. Unknown errors default to transient.

Checkpoint 2: Reconciliation Error Handling Create a Database resource with invalid configuration and verify the controller updates status with appropriate error condition instead of retrying indefinitely.

Checkpoint 3: Webhook Timeout Simulation Add artificial delays to webhook processing and verify that requests near the timeout deadline are handled gracefully with appropriate admission responses.

Debugging Tips

Symptom	Likely Cause	How to Diagnose	Fix
Resources stuck in pending state	Controller crashing on reconciliation errors	Check controller logs for panic stacktraces	Add error handling and classification to reconciliation loop
Webhook blocking all operations	Certificate expired or webhook pod not ready	Check webhook service endpoints and certificate validity	Implement certificate reloading and webhook health checks
Status never updates after errors	Status update conflicts not handled	Look for 409 Conflict errors in logs during status updates	Implement conflict retry in status update logic
Controller consuming excessive CPU	Infinite reconciliation loop from missing error handling	Monitor reconciliation frequency for specific resources	Add proper requeue delays and error classification
Resources stuck in deletion	Finalizer not removed due to cleanup errors	Check for resources with finalizers but no owning controller	Implement finalizer cleanup verification on controller startup

The error handling and recovery mechanisms form the foundation of operator reliability. By implementing comprehensive error classification, intelligent retry logic, and robust conflict resolution, the operator can provide consistent behavior even in the face of network partitions, API server outages, and concurrent modifications that are inherent in distributed Kubernetes environments.

Testing Strategy and Validation

Milestone(s): Milestone 5 (Testing & Deployment) - provides comprehensive testing approaches for unit testing, integration testing, and milestone validation checkpoints that ensure the operator functions correctly at each development stage

Testing a Kubernetes operator presents unique challenges that differ significantly from testing traditional applications. The operator must interact correctly with the Kubernetes API server, handle distributed system failures gracefully, and maintain consistent behavior across different cluster states. A robust testing strategy requires multiple layers of validation, from isolated unit tests that verify individual functions to full integration tests that exercise the complete controller lifecycle against a real API server.

Mental Model: The Quality Assurance Laboratory

Think of operator testing like a medical laboratory that validates treatments at multiple levels. Unit tests are like testing individual chemical reactions in isolation - they verify that specific functions produce expected outputs when given controlled inputs. Integration tests resemble clinical trials that test the complete treatment protocol on live subjects - they validate that the entire operator behaves correctly when interacting with a real Kubernetes environment. Finally, milestone validation checkpoints are like regulatory approval processes that ensure each development phase meets safety and efficacy standards before proceeding to the next phase.

Just as medical testing requires both controlled laboratory conditions and real-world clinical validation, operator testing demands both isolated unit tests with predictable fake clients and comprehensive integration tests against actual Kubernetes API servers. Each testing layer catches different categories of defects and provides different types of confidence in the system's correctness.

Unit Testing with Fake Clients

Unit testing forms the foundation of operator validation by testing individual controller functions in complete isolation from external dependencies. The controller-runtime library provides sophisticated fake client implementations that simulate Kubernetes API behavior

without requiring an actual cluster. These fake clients maintain in-memory representations of cluster state and respond to client operations with realistic behavior, including resource versioning, conflict detection, and validation.

Fake clients excel at testing reconciliation logic because they provide deterministic, controllable environments where specific error conditions can be simulated reliably. Unlike integration tests that may exhibit timing-dependent behavior, unit tests with fake clients produce identical results on every execution, making them ideal for testing edge cases and error scenarios that would be difficult to reproduce consistently in real clusters.

The fake client architecture maintains separate storage for each resource type, implementing the same optimistic locking semantics as the real Kubernetes API server. When a test creates a `Database` resource with the fake client, subsequent `Get` operations return the exact object with proper resource versions and generation fields. Updates that specify outdated resource versions fail with conflict errors, allowing tests to verify that the controller handles concurrent modifications correctly.

Core Testing Components

The unit testing infrastructure relies on several key components that work together to provide comprehensive test coverage:

Component	Type	Purpose
<code>fake.NewClientBuilder()</code>	Client Builder	Creates configurable fake Kubernetes client instances
<code>scheme.Scheme</code>	Runtime Scheme	Defines known resource types for client serialization
<code>envtest.Environment</code>	Test Environment	Provides real API server for integration tests
<code>reconcile.Request</code>	Reconciliation Request	Represents controller work queue items
<code>ctrl.Result</code>	Reconciliation Result	Indicates requeue behavior and timing

Testing Reconciliation Logic

The core of unit testing focuses on validating the `Reconcile` method behavior under various scenarios. Each test case should establish a known cluster state using the fake client, trigger reconciliation for a specific resource, and verify both the returned result and the resulting cluster state changes.

Testing successful reconciliation requires creating a `Database` resource with a valid specification, calling the reconciler, and verifying that owned resources are created with correct configurations. The test must check that the `Database` status is updated with appropriate conditions and that the reconciliation result indicates successful completion without requeue requirements.

Testing error conditions involves simulating various failure modes and verifying that the controller responds appropriately. This includes testing scenarios where owned resources cannot be created due to API server errors, where resource specifications contain invalid configurations, and where external dependencies are unavailable.

Testing idempotent behavior ensures that repeated reconciliation of the same resource produces identical results. This critical property guarantees that controller restarts or duplicate events do not cause unintended side effects or resource duplication.

Fake Client Configuration and Setup

Proper fake client configuration requires careful attention to scheme registration and resource type definitions. The fake client must understand all custom resource types that the controller manages, requiring explicit registration of CRDs in the runtime scheme.

```
// Example scheme registration (Implementation Guidance will show complete code)

scheme.AddToScheme(DatabaseCRD)

scheme.AddToScheme(SecretCRD)

fakeClient := fake.NewClientBuilder().WithScheme(scheme).Build()
```

GO

The fake client supports advanced features like resource generation tracking, owner reference validation, and finalizer processing. Tests can configure the fake client to simulate specific API server behaviors, such as temporary unavailability or resource quota enforcement.

Initial state setup involves creating prerequisite resources that the controller expects to exist. This might include `ConfigMap` resources containing default configurations, `Secret` resources with authentication credentials, or `Namespace` resources that own the custom resources being tested.

Status Update Validation

Testing status updates requires careful verification that the controller correctly writes status information without affecting the resource specification. The fake client maintains separate tracking for spec and status subresources, allowing tests to verify that status updates do not inadvertently modify the desired state.

Status Field	Test Verification	Purpose
<code>Conditions</code>	Contains expected condition types with correct status values	Communicates current reconciliation state
<code>ReadyReplicas</code>	Matches number of ready owned resources	Indicates operational capacity
<code>ObservedGeneration</code>	Equals metadata generation of reconciled resource	Shows which spec version was processed
<code>Phase</code>	Reflects current lifecycle stage	Provides high-level status summary
<code>LastBackupTime</code>	Updated when backup operations complete	Tracks operational activities
<code>Endpoints</code>	Lists accessible service endpoints	Provides connectivity information

Condition testing focuses on verifying that the controller creates appropriate condition entries for different scenarios. Ready conditions should indicate when all owned resources are operational, while progressing conditions should appear during active reconciliation activities. Error conditions must provide clear, actionable messages that help users understand and resolve issues.

Error Handling and Requeue Testing

Unit tests must thoroughly validate the controller's error handling behavior, ensuring that different error types trigger appropriate responses. Transient errors should result in exponential backoff requeue behavior, while permanent errors should update the resource status without scheduling automatic retries.

The fake client can be configured to return specific error types for different operations, allowing tests to verify controller behavior in response to various failure modes:

Error Scenario	Expected Behavior	Test Validation
API server unavailable	Return error with immediate requeue	Verify <code>Result.Requeue</code> is true
Resource already exists	Continue reconciliation normally	Verify no error returned
Invalid resource configuration	Update status with error condition	Verify condition type and message
Insufficient permissions	Return error without requeue	Verify <code>Result.RequeueAfter</code> is zero
Resource conflict	Return error with short requeue delay	Verify <code>Result.RequeueAfter</code> < 1 minute

Requeue timing validation ensures that the controller requests appropriate delays between reconciliation attempts. Immediate requeues should be reserved for situations where external state changes are expected quickly, while longer delays should be used for scenarios requiring user intervention or external system recovery.

Testing Finalizer Logic

Finalizer processing represents one of the most critical and error-prone aspects of operator behavior. Unit tests must verify that finalizers are added during resource creation, that cleanup operations execute correctly during deletion, and that finalizers are removed only after

successful cleanup completion.

Testing finalizer addition involves creating a new `Database` resource and verifying that the controller adds the appropriate finalizer during initial reconciliation. The test should confirm that the finalizer appears in the resource metadata and that subsequent reconciliation recognizes the finalizer's presence.

Testing deletion and cleanup requires setting a deletion timestamp on a resource with finalizers and verifying that cleanup operations execute in the correct order. The controller should delete all owned resources, wait for deletion completion, and only then remove the finalizer to allow garbage collection.

Testing cleanup failures ensures that incomplete cleanup operations do not remove finalizers prematurely. If owned resource deletion fails, the controller should leave finalizers in place and report error conditions in the resource status.

Common Unit Testing Pitfalls

Pitfall: Not Configuring Fake Client Scheme Properly

Tests frequently fail because the fake client doesn't understand custom resource types that the controller attempts to manipulate. This manifests as serialization errors or "no kind is registered" panics during test execution. The fix requires explicitly adding all relevant CRDs and built-in Kubernetes types to the fake client's scheme during test setup.

Pitfall: Testing Only Happy Path Scenarios

Many unit test suites focus exclusively on successful reconciliation scenarios while ignoring error conditions and edge cases. This creates false confidence in controller reliability because real-world environments frequently encounter API server unavailability, resource conflicts, and configuration errors. Comprehensive test suites should include more error scenarios than success scenarios.

Pitfall: Not Verifying Status Updates Separately

Tests often check that reconciliation completes successfully but fail to verify that status updates contain correct information. Since status subresources update independently from spec changes, tests must explicitly fetch and validate status fields after reconciliation completion. Missing status validation can hide bugs where controllers perform correct actions but provide incorrect feedback to users.

Pitfall: Ignoring Resource Version and Generation Fields

Fake clients properly implement optimistic locking through resource versions, but tests that don't account for these fields may pass artificially. Real controllers must handle resource version conflicts and generation tracking correctly, so unit tests should verify that these mechanisms work properly by attempting concurrent modifications and checking generation-based condition logic.

Integration Testing with Envtest

Integration testing validates operator behavior against real Kubernetes API servers, providing confidence that the controller works correctly within actual cluster environments. The `envtest` framework creates lightweight Kubernetes control planes that include `etcd` and `kube-apiserver` components without the overhead of full cluster deployments. These test environments support all standard Kubernetes API operations while running quickly enough for automated testing pipelines.

Envtest bridges the gap between unit tests and full cluster testing by providing real API server behavior while maintaining test isolation and repeatability. Unlike fake clients that simulate API responses, envtest environments process requests through actual Kubernetes controllers, including garbage collection, owner reference processing, and admission controller execution.

The envtest environment supports advanced Kubernetes features that fake clients cannot replicate, such as server-side apply, strategic merge patches, and admission webhook processing. This comprehensive API compatibility ensures that operators tested with envtest behave identically in production clusters, reducing the likelihood of integration issues during deployment.

Envtest Environment Setup

Envtest environments require careful configuration to match production cluster characteristics while maintaining test performance. The environment setup process involves downloading Kubernetes binaries, configuring API server options, and establishing client connections with appropriate authentication credentials.

Configuration Option	Purpose	Typical Value
KubernetesVersion	Specifies API server version	"1.28.0"
CRDInstallOptions.Paths	Locations of CRD YAML files	["./config/crd/bases"]
UseExistingCluster	Connect to running cluster instead of creating new one	false for CI, true for debugging
AttachControlPlaneOutput	Show API server logs during tests	true for debugging
ControlPlaneStartTimeout	Maximum time to wait for startup	60 * time.Second
ControlPlaneStopTimeout	Maximum time to wait for shutdown	60 * time.Second

CRD installation happens automatically during environment startup when CRD paths are configured correctly. The envtest framework reads YAML files from specified directories and creates custom resource definitions before starting the test API server. This process ensures that custom resources are available immediately when tests begin execution.

Client configuration uses the same controller-runtime managers and clients that production operators employ. This consistency ensures that integration tests exercise identical code paths to production environments, including client-side caching, informer behavior, and work queue processing.

Controller Manager Integration

Integration tests run complete controller managers that include informers, work queues, and reconciliation loops. This comprehensive testing approach validates not only individual reconciliation logic but also the interactions between different controller components and their behavior under concurrent load.

The controller manager setup process involves registering custom resource schemes, configuring controller options, and starting background processing goroutines. Integration tests must carefully coordinate controller startup with test execution to ensure that informer caches are populated and event processing is active before triggering test scenarios.

```
// Example controller manager setup (complete implementation in Implementation Guidance) GO
mgr := ctrl.NewManager(cfg, ctrl.Options{Scheme: scheme})

reconciler := &DatabaseReconciler{Client: mgr.GetClient(), Scheme: scheme}

reconciler.SetupWithManager(mgr)

ctx, cancel := context.WithCancel(context.Background())

go mgr.Start(ctx)
```

Manager lifecycle coordination requires proper startup sequencing and graceful shutdown handling. Tests must wait for manager readiness before executing test scenarios and must trigger clean shutdown to prevent goroutine leaks between test cases.

Informer cache synchronization introduces timing considerations that don't exist in unit tests. Integration tests must wait for informer caches to populate with initial resource state before verifying controller behavior, as controllers operating with empty caches may exhibit different behavior than those with fully synchronized state.

End-to-End Reconciliation Testing

Integration tests validate complete reconciliation cycles by creating custom resources, waiting for controller processing, and verifying that both owned resources and status updates reflect the expected final state. These tests provide confidence that all components work together correctly under realistic conditions.

Resource creation testing involves applying `Database` resources to the envtest cluster and monitoring both the reconciliation process and its outcomes. Tests must account for asynchronous processing by polling resource status until reconciliation completes or timeout periods expire.

Owned resource verification ensures that controllers create, configure, and manage owned resources correctly. Integration tests should verify not only that owned resources exist but also that they contain correct specifications, labels, annotations, and owner references that enable proper garbage collection.

Status progression tracking monitors how resource status evolves during reconciliation cycles. Tests should verify that progressing conditions appear during active reconciliation, that ready conditions indicate successful completion, and that error conditions provide actionable diagnostic information.

Testing Concurrent Reconciliation

Integration tests can validate controller behavior under concurrent load by creating multiple custom resources simultaneously and verifying that reconciliation proceeds correctly for all resources. This testing approach reveals race conditions, resource contention issues, and queue processing problems that single-resource tests cannot detect.

Multiple resource testing creates several `Database` resources with different configurations and verifies that each resource reconciles independently without interference. These tests should confirm that controller scaling works correctly and that individual resource failures do not affect processing of other resources.

Resource update testing modifies existing resources while reconciliation is active, simulating real-world scenarios where users update specifications while controllers are processing previous changes. Tests should verify that controllers handle spec changes gracefully and that status updates reflect the most recent reconciliation results.

Concurrent modification testing attempts to modify resources from multiple clients simultaneously, verifying that controllers handle resource version conflicts and retry operations appropriately. These tests ensure that operators behave correctly in multi-user environments where administrative operations may overlap with controller activities.

Webhook Integration Testing

Integration tests with envtest can validate admission webhook behavior by configuring webhook endpoints and processing admission review requests through the test API server. This comprehensive testing approach ensures that webhooks integrate correctly with the Kubernetes admission control system.

Webhook server setup involves starting HTTPS servers with appropriate TLS certificates and registering admission webhook configurations with the test API server. The envtest framework supports webhook testing through proper certificate management and admission controller configuration.

Validation webhook testing creates resources with invalid specifications and verifies that admission webhooks reject requests with appropriate error messages. Tests should confirm that validation logic correctly identifies invalid field combinations and provides clear diagnostic information to users.

Mutation webhook testing creates resources with minimal specifications and verifies that mutating webhooks inject appropriate default values. Tests should confirm that mutation logic produces consistent results and that mutated resources pass subsequent validation checks.

Performance and Load Testing

Integration tests can evaluate operator performance characteristics by creating large numbers of resources and measuring reconciliation throughput, memory usage, and response times. These measurements provide baseline performance expectations and help identify scalability limitations.

Throughput testing creates batches of resources and measures how quickly the controller processes them to completion. These tests help establish capacity planning guidelines and identify performance bottlenecks in reconciliation logic.

Memory usage testing monitors controller memory consumption while processing various workloads to identify memory leaks or excessive resource usage. Long-running tests can reveal gradual memory growth that would not be apparent in short-duration scenarios.

Response time testing measures the delay between resource creation and reconciliation completion, providing insights into user experience and system responsiveness. These measurements help establish appropriate timeout values and requeue intervals.

Common Integration Testing Pitfalls

⚠️ Pitfall: Not Waiting for Controller Readiness

Integration tests frequently fail because they attempt to validate controller behavior before informer caches are fully synchronized or before controllers have registered their watches. This manifests as tests that pass when run individually but fail when executed as part of larger test suites. The fix requires implementing proper readiness checks that wait for manager startup completion and informer cache synchronization.

⚠️ Pitfall: Ignoring Asynchronous Processing

Many integration tests assume that reconciliation completes immediately after resource creation, leading to race conditions where tests check results before processing finishes. Real controllers process events asynchronously through work queues, requiring tests to poll for expected conditions with appropriate timeout values rather than checking results immediately.

⚠️ Pitfall: Not Cleaning Up Resources Between Tests

Integration tests that don't properly clean up created resources can exhibit order-dependent failures where previous test state affects subsequent test execution. Since envtest environments persist resources across individual test cases within a suite, tests must explicitly delete created resources or use unique namespaces to ensure isolation.

⚠️ Pitfall: Insufficient Timeout Values

Integration tests often use timeout values that work during local development but fail in resource-constrained CI environments. Envtest startup, controller initialization, and reconciliation processing all require sufficient time to complete, particularly when running on shared infrastructure with limited CPU resources.

Milestone Validation Checkpoints

Each development milestone requires specific validation checkpoints that verify correct implementation before proceeding to subsequent phases. These checkpoints combine automated testing with manual verification to ensure that foundational components work correctly before adding additional complexity.

Milestone validation provides confidence gates that prevent cascading failures where issues in early milestones cause confusing symptoms in later development phases. By establishing clear success criteria for each milestone, developers can focus on one set of challenges at a time while building incrementally toward a complete operator implementation.

The validation checkpoints include both positive testing that verifies expected behavior and negative testing that confirms appropriate error handling. Each checkpoint also includes performance expectations and operational characteristics that should be observable during testing.

Milestone 1: Custom Resource Definition Validation

The first milestone focuses on validating that custom resource definitions are properly configured with appropriate schema validation, status subresources, and versioning support. These foundational elements must work correctly before controller implementation begins.

Validation Check	Expected Behavior	Verification Method
CRD Installation	CRD appears in cluster and accepts valid resources	<code>kubectl get crd databases.example.com</code>
Schema Validation	Invalid resources are rejected with clear error messages	<code>kubectl apply</code> with invalid YAML
Status Subresource	Status updates don't affect spec fields	Separate <code>kubectl patch</code> operations
Printer Columns	Custom columns appear in <code>kubectl get</code> output	<code>kubectl get databases</code> shows custom fields
OpenAPI Integration	API discovery includes CRD schema information	<code>kubectl explain database.spec</code>
Multiple Versions	Both v1alpha1 and v1beta1 versions are accepted	Create resources with different apiVersions

Schema validation testing involves creating `Database` resources with various invalid configurations and verifying that the API server rejects them with helpful error messages. Tests should cover missing required fields, invalid enum values, and constraint violations like negative replica counts.

Status subresource testing confirms that status updates operate independently from spec modifications. This involves patching status fields directly and verifying that the resource generation field doesn't increment, indicating that spec changes and status updates are properly isolated.

Version conversion testing creates resources using older API versions and verifies that they can be retrieved using newer versions, with appropriate field mapping and default value injection occurring during conversion.

Milestone 2: Controller Setup Validation

The second milestone validates that the controller infrastructure is properly configured with working client connections, informer caches, and work queue processing. These components must function correctly before implementing reconciliation logic.

Validation Check	Expected Behavior	Verification Method
Client Connection	Controller can list and watch custom resources	Log messages show successful connection
Informer Startup	Informer cache synchronizes with API server state	Ready condition in manager status
Event Processing	Resource changes trigger work queue entries	Debug logs show enqueued reconciliation requests
Worker Goroutines	Multiple reconciliation requests process concurrently	Concurrent processing visible in logs
Leader Election	Only one controller replica processes events	Multiple replicas with one active leader
RBAC Permissions	Controller can read, create, update, and delete required resources	No permission denied errors

Client functionality testing creates, updates, and deletes resources using the controller's Kubernetes client to verify that authentication and authorization work correctly. This includes testing both cached reads through informers and direct API calls for write operations.

Informer behavior testing creates resources outside the controller and verifies that informer caches reflect changes within reasonable time periods. Tests should also verify that informer resync operations occur at configured intervals.

Work queue testing triggers rapid resource changes and verifies that work queue processing handles the load without dropping events or creating excessive goroutines. Rate limiting behavior should be observable during high-frequency updates.

Milestone 3: Reconciliation Loop Validation

The third milestone validates that reconciliation logic correctly compares desired versus actual state and takes appropriate corrective actions. This represents the core operator functionality and requires comprehensive testing across various scenarios.

Validation Check	Expected Behavior	Verification Method
Resource Creation	New Database resources trigger owned resource creation	Deployment and Service objects appear
Resource Updates	Spec changes propagate to owned resources	Modified owned resources match new spec
Resource Deletion	Finalizer processing cleans up owned resources	Owned resources deleted before Database removal
Status Updates	Status reflects current reconciliation state	Conditions indicate progress and completion
Error Handling	Reconciliation failures update status appropriately	Error conditions with diagnostic messages
Idempotent Operations	Repeated reconciliation produces identical results	Multiple reconcile calls show no changes

State comparison testing modifies owned resources directly and verifies that the controller detects differences and restores desired state. This includes testing configuration drift scenarios where external processes modify controller-managed resources.

Reconciliation timing testing measures how quickly controllers detect changes and complete reconciliation cycles. Response times should be consistent and appropriate for the complexity of managed resources.

Error recovery testing introduces various failure conditions and verifies that controllers handle them gracefully with appropriate status reporting and retry behavior. This includes testing API server unavailability, resource conflicts, and invalid configurations.

Milestone 4: Webhook Validation

The fourth milestone validates that admission webhooks correctly validate and mutate resources according to business rules and policy requirements. Webhook behavior must be thoroughly tested because admission failures can prevent resource creation entirely.

Validation Check	Expected Behavior	Verification Method
Webhook Registration	Admission webhook configurations appear in cluster	<code>kubectl get validatingwebhookconfiguration</code>
TLS Certificate Setup	Webhook server accepts HTTPS connections	Certificate validation in server logs
Validation Logic	Invalid resources are rejected with clear messages	<code>kubectl apply</code> with policy violations
Mutation Logic	Default values are injected into new resources	Compare applied vs original resource content
Webhook Availability	Admission processing continues during webhook updates	No service disruption during deployment
Error Responses	Admission failures include actionable error messages	Clear feedback about validation failures

Validation testing attempts to create resources that violate business rules and verifies that webhooks reject them with informative error messages. Tests should cover both simple field validation and complex cross-field business logic.

Mutation testing creates resources with minimal specifications and verifies that webhooks inject appropriate default values and computed fields. The resulting resources should be valid and functionally complete.

Certificate management testing verifies that webhook certificates are properly provisioned, renewed before expiration, and updated without service disruption. This includes testing both cert-manager integration and manual certificate management approaches.

Milestone 5: Testing and Deployment Validation

The final milestone validates that the complete operator functions correctly in production-like environments with proper security configuration, monitoring integration, and operational procedures. This comprehensive validation ensures that the operator is ready for production deployment.

Validation Check	Expected Behavior	Verification Method
Unit Test Coverage	All reconciliation paths have comprehensive test coverage	Code coverage reports show >80% coverage
Integration Tests	End-to-end scenarios pass consistently	Automated test suite passes in CI
Container Image Build	Operator image builds successfully with proper labeling	<code>docker build</code> produces tagged image
RBAC Configuration	Deployed operator has minimal required permissions	Security audit of ClusterRole permissions
High Availability	Multiple operator replicas with leader election work correctly	Pod failures don't interrupt reconciliation
Monitoring Integration	Metrics and logs provide operational visibility	Prometheus metrics and structured logs

End-to-end testing deploys the operator to test clusters and validates that complete user workflows function correctly from resource creation through deletion. These tests should simulate real user interactions and verify that all features work together seamlessly.

Production readiness testing evaluates the operator under production-like conditions including resource constraints, network partitions, and concurrent user activity. These tests identify potential issues that might not appear under ideal testing conditions.

Operational testing verifies that administrative procedures like operator upgrades, certificate rotation, and backup/restore operations work correctly without disrupting managed applications.

Common Milestone Validation Pitfalls

⚠ Pitfall: Skipping Negative Testing

Many validation checkpoints focus only on verifying that expected behaviors work correctly while ignoring error conditions and edge cases. Production environments frequently encounter unexpected inputs, network failures, and resource constraints that can expose untested code paths. Comprehensive validation requires testing both success and failure scenarios at each milestone.

⚠ Pitfall: Not Testing Real-World Conditions

Validation that works in clean, isolated test environments may fail when deployed to production clusters with existing workloads, resource quotas, and security policies. Effective milestone validation should include testing under conditions that approximate production environments, including limited resources and competing workloads.

⚠ Pitfall: Insufficient Soak Testing

Many operators that pass short-duration tests exhibit issues like memory leaks, performance degradation, or state corruption over longer periods. Milestone validation should include extended testing periods that allow time-dependent issues to manifest, particularly for controllers that process high event volumes.

⚠ Pitfall: Not Validating User Experience

Technical validation often focuses on internal correctness while ignoring the user experience aspects like error message clarity, documentation accuracy, and operational complexity. Production operators must provide excellent user experiences, requiring validation that includes usability testing from the perspective of actual users.

Implementation Guidance

This subsection provides practical implementation details for building comprehensive testing infrastructure that supports effective operator development and validation.

Technology Recommendations

Testing Component	Simple Option	Advanced Option
Unit Testing Framework	Go's built-in <code>testing</code> package	Ginkgo + Gomega for BDD-style testing
Integration Testing	controller-runtime's <code>envtest</code> package	Kind clusters for full Kubernetes features
Mock Generation	Manual fake implementations	<code>go generate</code> with <code>counterfeiter</code> or <code>gomock</code>
Test Data Management	Hard-coded YAML strings	External YAML files with template processing
Assertion Library	Basic Go comparisons	<code>testify/assert</code> for rich assertion methods
Coverage Analysis	<code>go test -cover</code> built-in coverage	<code>gocov + gocov-html</code> for detailed reports

Recommended File Structure

```
project-root/
  controllers/
    database_controller.go      ← main controller implementation
    database_controller_test.go ← unit tests with fake client
    suite_test.go               ← integration test suite setup
  webhooks/
    database_webhook.go         ← admission webhook implementation
    database_webhook_test.go    ← unit tests for webhook logic
    webhook_suite_test.go       ← integration tests with envtest
  testdata/
    valid-database.yaml        ← test resource manifests
    invalid-database.yaml      ← resources for validation testing
    webhook-certificates/
      tls.crt
      tls.key
  hack/
    setup-envtest.sh           ← script to install envtest binaries
    run-integration-tests.sh   ← script for CI integration testing
```

Complete Unit Testing Infrastructure

This infrastructure provides a complete foundation for unit testing with fake clients:

```
// controllers/suite_test.go - Complete unit test setup
```

package controllers

```
import (
    "context"
    "path/filepath"
    "testing"
    "time"

    . "github.com/onsi/ginkgo/v2"
    . "github.com/onsi/gomega"

    corev1 "k8s.io/api/core/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/client-go/kubernetes/scheme"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/client/fake"
    "sigs.k8s.io/controller-runtime/pkg/envtest"
    "sigs.k8s.io/controller-runtime/pkg/log/zap"
    ctrl "sigs.k8s.io/controller-runtime"

    examplev1 "github.com/example/database-operator/api/v1"
)
```

var (

```
    k8sClient client.Client
    testEnv   *envtest.Environment
    ctx       context.Context
    cancel    context.CancelFunc
)
```

```
func TestControllers(t *testing.T) {
    RegisterFailHandler(Fail)
    RunSpecs(t, "Controller Suite")
}
```

```

var _ = BeforeSuite(func() {
    ctrl.SetLogger(zap.New(zap.WriteTo(GinkgoWriter), zap.UseDevMode(true)))

    ctx, cancel = context.WithCancel(context.Background())

    By("bootstrapping test environment")

    testEnv = &envtest.Environment{
        CRDDirectoryPaths: []string{filepath.Join("../", "config", "crd", "bases")},
        ErrorIfCRDPathMissing: true,
    }

    cfg, err := testEnv.Start()

    Expect(err).NotTo(HaveOccurred())
    Expect(cfg).NotTo(BeNil())

    err = examplev1.AddToScheme(scheme.Scheme)

    Expect(err).NotTo(HaveOccurred())

    k8sClient, err = client.New(cfg, client.Options{Scheme: scheme.Scheme})
    Expect(err).NotTo(HaveOccurred())
    Expect(k8sClient).NotTo(BeNil())
})

var _ = AfterSuite(func() {
    cancel()

    By("tearing down the test environment")

    err := testEnv.Stop()
    Expect(err).NotTo(HaveOccurred())
})
}

// Test helper functions

func CreateTestNamespace() *corev1.Namespace {
    ns := &corev1.Namespace{
        ObjectMeta: metav1.ObjectMeta{
            GenerateName: "test-",
        },
}

```

```
}

Expect(k8sClient.Create(ctx, ns)).To(Succeed())

return ns

}

func CreateTestDatabase(namespace, name string) *examplev1.Database {

db := &examplev1.Database{
ObjectMeta: metav1.ObjectMeta{
Name:      name,
Namespace: namespace,
},
Spec: examplev1.DatabaseSpec{
Replicas:   3,
Version:    "13.0",
StorageSize: "10Gi",
},
}

Expect(k8sClient.Create(ctx, db)).To(Succeed())

return db
}
```

Core Reconciliation Test Skeleton

This skeleton provides the structure for comprehensive reconciliation testing:

```
// controllers/database_controller_test.go - Reconciliation test template
GO

package controllers

import (
    . "github.com/onsi/ginkgo/v2"
    . "github.com/onsi/gomega"

    corev1 "k8s.io/api/core/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/types"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client/fake"

    examplev1 "github.com/example/database-operator/api/v1"
)

var _ = Describe("Database Controller", func() {
    var (
        reconciler *DatabaseReconciler
        namespace  *corev1.Namespace
        database   *examplev1.Database
    )

    BeforeEach(func() {
        namespace = CreateTestNamespace()
        reconciler = &DatabaseReconciler{
            Client: k8sClient,
            Scheme: scheme.Scheme,
        }
    })

    Context("When creating a new Database", func() {
        BeforeEach(func() {
            database = CreateTestDatabase(namespace.Name, "test-db")
        })
    })
})
```

```

It("Should create owned resources successfully", func() {
    // TODO: Call reconciler.Reconcile() with appropriate request
    // TODO: Verify that Deployment and Service resources are created
    // TODO: Check that owner references are set correctly
    // TODO: Validate that finalizers are added to the Database resource
    // TODO: Confirm that status conditions indicate successful creation
})

It("Should handle missing namespace gracefully", func() {
    // TODO: Delete the namespace before reconciliation
    // TODO: Call reconciler.Reconcile() and expect appropriate error handling
    // TODO: Verify that status conditions indicate the error state
    // TODO: Confirm that no owned resources are created in invalid state
})

Context("When updating an existing Database", func() {
    BeforeEach(func() {
        database = CreateTestDatabase(namespace.Name, "test-db")
        // TODO: Wait for initial reconciliation to complete
        // TODO: Verify that owned resources exist and are ready
    })

    It("Should propagate spec changes to owned resources", func() {
        // TODO: Update database spec (e.g., change replica count)
        // TODO: Call reconciler.Reconcile() with update event
        // TODO: Verify that owned Deployment reflects new replica count
        // TODO: Check that status indicates progressing state during update
        // TODO: Confirm that ready condition updates after completion
    })

    It("Should detect and correct configuration drift", func() {
        // TODO: Manually modify an owned resource (e.g., change Deployment image)
        // TODO: Call reconciler.Reconcile() to trigger drift detection
        // TODO: Verify that controller restores correct configuration
    })
})

```

Integration Testing with Controller Manager

This template demonstrates full integration testing with running controllers:

```
// controllers/integration_test.go - Full integration test template

package controllers

import (
    "time"

    . "github.com/onsi/ginkgo/v2"
    . "github.com/onsi/gomega"

    appsv1 "k8s.io/api/apps/v1"
    corev1 "k8s.io/api/core/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/types"
    ctrl "sigs.k8s.io/controller-runtime"

    examplev1 "github.com/example/database-operator/api/v1"
)

var _ = Describe("Database Integration", func() {

    var (
        mgr      ctrl.Manager
        namespace *corev1.Namespace
        database *examplev1.Database
    )

    BeforeEach(func() {

        var err error

        mgr, err = ctrl.NewManager(cfg, ctrl.Options{
            Scheme: scheme.Scheme,
            Port:   9443, // webhook port
        })
        Expect(err).ToNot(HaveOccurred())

        reconciler := &DatabaseReconciler{
            Client: mgr.GetClient(),
            Scheme: mgr.GetScheme(),
        }
    })
})
```

GO


```

    }, deployment)

    }, time.Minute, time.Second).Should(Succeed())

    By("Verifying Deployment configuration")

    Expect(deployment.Spec.Replicas).To(Equal(int32Ptr(2)))

    Expect(deployment.OwnerReferences).To(HaveLen(1))

    Expect(deployment.OwnerReferences[0].Name).To(Equal(database.Name))

    By("Waiting for Database status to indicate readiness")

    Eventually(func() bool {
        err := k8sClient.Get(ctx, types.NamespacedName{
            Name:      database.Name,
            Namespace: database.Namespace,
        }, database)
        if err != nil {
            return false
        }
        return database.Status.Phase == "Ready"
    }, time.Minute, time.Second).Should(BeTrue())
})

})

func int32Ptr(i int32) *int32 {
    return &i
}

```

Webhook Testing Infrastructure

This infrastructure enables comprehensive webhook testing:

```
// webhooks/webhook_suite_test.go - Webhook testing setup

package webhooks

import (
    "context"
    "crypto/tls"
    "fmt"
    "net"
    "path/filepath"
    "testing"
    "time"

    . "github.com/onsi/ginkgo/v2"
    . "github.com/onsi/gomega"

    admissionv1 "k8s.io/api/admission/v1"
    metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/client-go/kubernetes/scheme"
    ctrl "sigs.k8s.io/controller-runtime"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/envtest"
    "sigs.k8s.io/controller-runtime/pkg/webhook"
    "sigs.k8s.io/controller-runtime/pkg/webhook/admission"
    examplev1 "github.com/example/database-operator/api/v1"
)

var (
    k8sClient client.Client
    testEnv   *envtest.Environment
    ctx       context.Context
    cancel    context.CancelFunc
)

func TestWebhooks(t *testing.T) {
    RegisterFailHandler(Fail)
```

GO

```

RunSpecs(t, "Webhook Suite")

}

var _ = BeforeSuite(func() {
    ctx, cancel = context.WithCancel(context.Background())

    By("bootstrapping test environment")

    testEnv = &envtest.Environment{
        CRDDirectoryPaths: []string{filepath.Join("../", "config", "crd", "bases")},
        ErrorIfCRDPathMissing: true,
        WebhookInstallOptions: envtest.WebhookInstallOptions{
            Paths: []string{filepath.Join("../", "config", "webhook")},
        },
    }

    cfg, err := testEnv.Start()
    Expect(err).NotTo(HaveOccurred())
    Expect(cfg).NotTo(BeNil())

    err = examplev1.AddToScheme(scheme.Scheme)
    Expect(err).NotTo(HaveOccurred())

    k8sClient, err = client.New(cfg, client.Options{Scheme: scheme.Scheme})
    Expect(err).NotTo(HaveOccurred())

    // Start webhook server

    webhookInstallOptions := &testEnv.WebhookInstallOptions

    mgr, err := ctrl.NewManager(cfg, ctrl.Options{
        Scheme:             scheme.Scheme,
        Host:               webhookInstallOptions.LocalServingHost,
        Port:               webhookInstallOptions.LocalServingPort,
        CertDir:             webhookInstallOptions.LocalServingCertDir,
        LeaderElection:     false,
        MetricsBindAddress: "0",
    })
    Expect(err).NotTo(HaveOccurred())
})

```

```

    // TODO: Setup webhook handlers with manager

    // TODO: Start manager in goroutine

})

var _ = AfterSuite(func() {
    cancel()

    By("tearing down the test environment")

    err := testEnv.Stop()

    Expect(err).NotTo(HaveOccurred())
})

```

Milestone Validation Checkpoints

Each milestone includes specific commands and expected outputs for validation:

Milestone 1 Checkpoint Commands:

```

# Verify CRD installation

kubectl get crd databases.example.com -o yaml

# Test schema validation

kubectl apply -f testdata/invalid-database.yaml

# Expected: validation error with specific field requirements

# Test status subresource

kubectl patch database test-db --subresource=status --type=merge -p='{"status":{"phase":"Testing"}}'

# Expected: status updated without spec generation increment

```

Milestone 2 Checkpoint Commands:

```

# Start controller with debug logging

go run ./cmd/main.go --zap-log-level=debug

# Create test resource and watch logs

kubectl apply -f testdata/valid-database.yaml

# Expected: reconciliation request logged within 5 seconds

# Test concurrent processing

for i in {1..10}; do kubectl apply -f testdata/database-$i.yaml; done

# Expected: multiple concurrent reconciliation logs

```

Language-Specific Testing Hints:

- Use `testify/suite` for complex test setup with `BeforeTest/AfterTest` hooks
- Leverage `t.Parallel()` for unit tests that don't share state
- Use `go test -race` to detect race conditions in controller code
- Configure `KUBEBUILDER_ASSETS` environment variable for envtest binary location
- Use `Eventually()` and `Consistently()` from Gomega for asynchronous assertions
- Mock external dependencies using `testify/mock` or `counterfeiter`
- Use `t.Helper()` in test utility functions to improve error reporting
- Configure test timeouts appropriately: unit tests 10s, integration tests 60s

Common Debugging Scenarios:

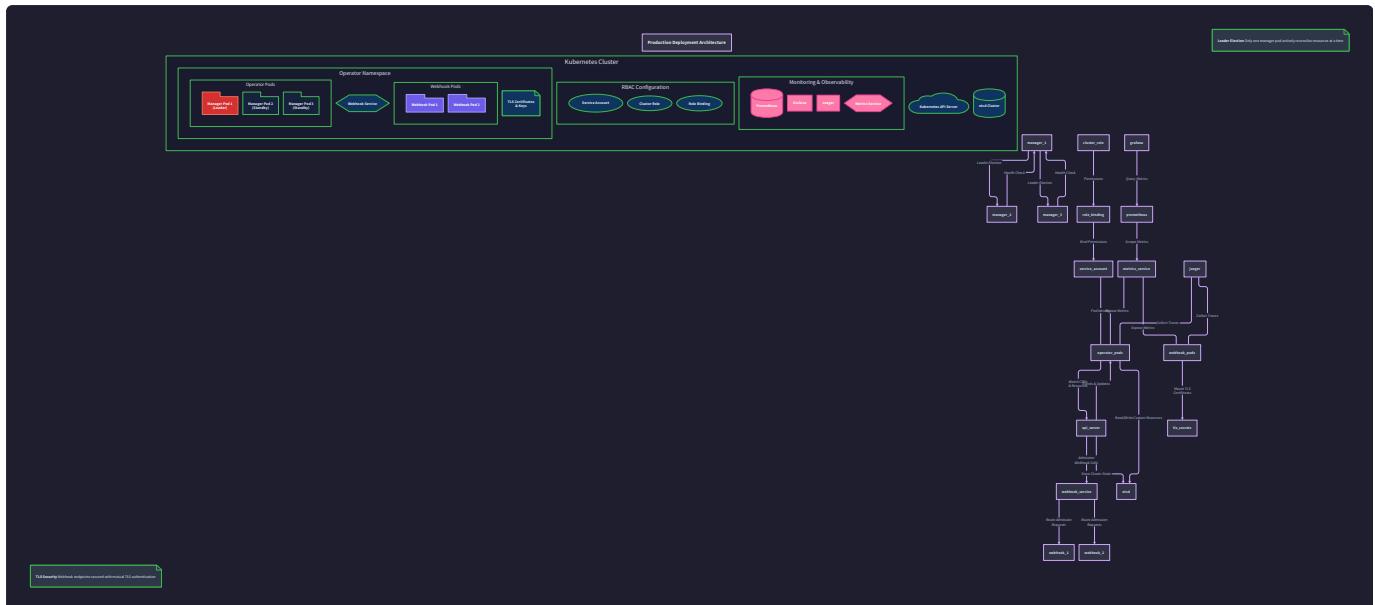
Symptom	Likely Cause	Diagnosis	Fix
Unit tests panic with "no kind is registered"	Missing scheme registration	Check fake client scheme setup	Add all CRDs to <code>scheme.Scheme</code>
Integration tests timeout waiting for resources	Controller not running or not watching	Check manager startup and controller registration	Verify <code>SetupWithManager()</code> call
Webhook tests fail with certificate errors	TLS configuration incorrect	Check certificate paths and validity	Regenerate test certificates
Tests pass individually but fail in suite	Resource cleanup between tests	Check for resource leakage	Add proper cleanup in <code>AfterEach</code>
Flaky test failures in CI	Race conditions or insufficient timeouts	Increase timeout values and add synchronization	Use <code>Eventually()</code> with longer durations

Deployment and Operations

Milestone(s): Milestone 5 (Testing & Deployment) - covers packaging the operator as container images, RBAC configuration, and production deployment strategies

Moving a Kubernetes operator from development to production requires careful attention to security, packaging, and operational concerns. Think of deploying an operator like installing a new resident admin in your apartment building - they need the right keys (RBAC permissions), their own secure apartment (container packaging), and a backup admin available if they get sick (high availability). Unlike regular applications that just process requests, operators have elevated privileges to manage cluster resources, making security and reliability critical considerations.

The deployment phase transforms your operator from a development artifact into a production-ready service that can safely manage workloads at scale. This involves three key challenges: ensuring the operator has exactly the permissions it needs and no more (RBAC), packaging it as a container with proper configuration management (container packaging), and ensuring it remains available even during node failures or rolling updates (high availability).



RBAC and Security Model

Mental Model: The Building Superintendent's Keys

Think of Kubernetes RBAC like a building superintendent's key ring. A good superintendent has keys to common areas (cluster-scoped resources) and individual apartments they manage (namespace-scoped resources), but they don't have keys to the bank vault next door (unrelated cluster resources) or other buildings entirely (other clusters). Each key serves a specific purpose, and losing or misusing any key has clear consequences. Your operator needs exactly the right keys to do its job effectively without creating security risks.

The **principle of least privilege** governs operator RBAC design. Rather than granting broad permissions like `cluster-admin`, operators should receive only the specific permissions required for their reconciliation logic. This approach limits the blast radius if the operator is compromised and makes security auditing more straightforward.

Decision: Namespace-Scope vs Cluster-Scope Permissions

- **Context:** Operators can be deployed with cluster-wide permissions or restricted to specific namespaces
- **Options Considered:**
 1. Cluster-scoped with ClusterRole and ClusterRoleBinding
 2. Namespace-scoped with Role and RoleBinding
 3. Hybrid approach with multiple permission levels
- **Decision:** Start with namespace-scoped permissions and escalate to cluster-scoped only when required
- **Rationale:** Namespace-scoped permissions provide better security isolation, easier multi-tenancy, and clearer audit trails. Many operators only need to manage resources within specific namespaces.
- **Consequences:** Requires more complex permission modeling for cross-namespace operations but provides stronger security boundaries

Permission Scope	Resource Types	Use Cases	Security Trade-offs
Cluster-scoped	ClusterRole, ClusterRoleBinding, CRDs	Managing cluster-wide resources, cross-namespace operations	Broad permissions, single point of failure, harder to audit
Namespace-scoped	Role, RoleBinding within target namespaces	Application-specific operators, multi-tenant environments	Limited blast radius, easier isolation, more complex setup
Hybrid	Multiple roles with different scopes	Operators managing both namespace and cluster resources	Balanced approach, requires careful permission design

The `DatabaseReconciler` requires specific permissions to manage its custom resources and owned Kubernetes resources. These permissions are defined through a combination of RBAC resources that work together to grant the operator's ServiceAccount the necessary access.

RBAC Component	Purpose	Scope	Example Resource
ServiceAccount	Identity for operator pods	Namespace	<code>database-operator-service-account</code>
Role/ClusterRole	Permission definitions	Namespace/Cluster	Rules for Database CRD, Secrets, Services
RoleBinding/ClusterRoleBinding	Links ServiceAccount to permissions	Namespace/Cluster	Grants operator ServiceAccount the defined roles

Core Permission Requirements

The `DatabaseReconciler` needs specific permissions to perform its reconciliation duties. These permissions should be carefully scoped to prevent privilege escalation while enabling full operator functionality.

Resource Type	Operations Needed	Justification	Example Usage
<code>databases.example.com</code>	get, list, watch, update, patch	Monitor custom resources and update status	Watch Database resource changes, update status subresource
<code>secrets</code>	get, list, watch, create, update, patch, delete	Manage database credentials and TLS certificates	Create admin passwords, manage backup encryption keys
<code>services</code>	get, list, watch, create, update, patch, delete	Expose database instances	Create ClusterIP services for database connectivity
<code>deployments</code>	get, list, watch, create, update, patch, delete	Manage database workloads	Deploy database pods with specified replica counts
<code>persistentvolumeclaims</code>	get, list, watch, create, update, patch, delete	Manage database storage	Provision storage volumes for database data
<code>configmaps</code>	get, list, watch, create, update, patch, delete	Manage database configuration	Store database config files, backup scripts
<code>events</code>	create, patch	Record operator activities	Log reconciliation events for debugging

The critical insight here is that every permission must map to a specific reconciliation action. If you cannot explain why the operator needs a particular permission by pointing to specific reconciliation logic, that permission should be removed.

Owner References and Garbage Collection

Owner references establish parent-child relationships between custom resources and their managed Kubernetes resources. This mechanism enables automatic garbage collection when custom resources are deleted, reducing the operator's cleanup burden and

preventing resource leaks.

Owner Reference Field	Purpose	Example Value	Behavior on Deletion
<code>apiVersion</code>	API version of parent resource	<code>example.com/v1</code>	Links to specific API version
<code>kind</code>	Resource type of parent	<code>Database</code>	Identifies parent resource type
<code>name</code>	Name of parent resource	<code>production-postgres</code>	Links to specific resource instance
<code>uid</code>	Unique identifier of parent	<code>abc123-def456-...</code>	Ensures precise parent identification
<code>controller</code>	Indicates controlling owner	<code>true</code>	Triggers cascading deletion
<code>blockOwnerDeletion</code>	Blocks parent deletion until child is removed	<code>false</code>	Allows graceful cleanup ordering

⚠ Pitfall: Missing Owner References Leading to Resource Leaks

Forgetting to set owner references on managed resources causes them to persist after the custom resource is deleted. This creates orphaned resources that consume cluster capacity and complicate cleanup. Always set owner references when creating managed resources, and verify garbage collection behavior in your tests.

```
// WRONG: Creating resource without owner reference

service := &corev1.Service{...}

err := r.Client.Create(ctx, service)

// CORRECT: Setting owner reference for garbage collection

service := &corev1.Service{...}

err := ctrl.SetControllerReference(database, service, r.Scheme)

if err != nil {
    return ctrl.Result{}, err
}

err = r.Client.Create(ctx, service)
```

GO

Webhook RBAC Requirements

Admission webhooks require additional RBAC permissions beyond the controller's reconciliation needs. These permissions enable webhook registration, certificate management, and admission review processing.

Webhook Component	Required Permissions	Resource Types	Justification
ValidatingAdmissionWebhook registration	create, update, patch	<code>validatingadmissionwebhooks.admissionregistration.k8s.io</code>	Register webhook with API server
MutatingAdmissionWebhook registration	create, update, patch	<code>mutatingadmissionwebhooks.admissionregistration.k8s.io</code>	Register webhook with API server
Certificate management	get, list, watch, create, update	<code>secrets</code> (for TLS certificates)	Manage webhook TLS certificates
Service management	get, list, watch, create, update	<code>services</code> (for webhook endpoints)	Expose webhook HTTPS endpoints

Decision: Webhook Certificate Management Strategy

- **Context:** Webhooks require TLS certificates for HTTPS communication with the API server
- **Options Considered:**
 1. Self-signed certificates generated at startup
 2. cert-manager for automatic certificate lifecycle management
 3. External certificate authority with manual certificate provision
- **Decision:** Use cert-manager for production deployments with self-signed fallback for development
- **Rationale:** cert-manager provides automatic renewal, proper CA trust chains, and integrates well with Kubernetes RBAC. Self-signed certificates work for development but require manual rotation in production.
- **Consequences:** Adds cert-manager dependency but eliminates certificate management toil and reduces security risks from expired certificates

Container Packaging and Deployment

Mental Model: The Shipping Container

Think of packaging your operator like preparing a shipping container for international transport. The container must be self-contained with all dependencies included, properly labeled with version and contents, and standardized so it runs the same way whether deployed in development, staging, or production. Like shipping containers, operator images should be immutable, versioned, and contain everything needed to run in any compliant Kubernetes cluster.

Container packaging transforms your operator from source code into a deployable artifact that can be distributed, versioned, and deployed consistently across environments. This process involves building container images, organizing deployment manifests, and creating configuration management systems that handle environment-specific variations.

Container Image Strategy

The operator container image serves as the fundamental deployment unit, containing the compiled binary, runtime dependencies, and configuration templates. Image design affects security, deployment speed, and operational complexity.

Image Characteristic	Development Approach	Production Approach	Trade-offs
Base image	Full OS (ubuntu, centos)	Minimal (alpine, distroless)	Debug tools vs. attack surface
Binary packaging	Debug builds with symbols	Optimized builds stripped	Debugging capability vs. image size
Dependency management	Package manager installs	Multi-stage builds	Build simplicity vs. reproducibility
User permissions	Root user	Non-root user	Convenience vs. security
Layer optimization	Single layer builds	Multi-stage with layer caching	Build speed vs. image efficiency

Decision: Multi-Stage Build with Distroless Base Image

- Context:** Operator images need to be secure, small, and fast to deploy while containing all required dependencies
- Options Considered:**
 - Single-stage build with full OS base image
 - Multi-stage build with minimal base image
 - Scratch base image with static binary
- Decision:** Multi-stage build with Google's distroless base image
- Rationale:** Distroless images contain only the application and runtime dependencies, reducing attack surface while maintaining compatibility. Multi-stage builds separate build dependencies from runtime dependencies.
- Consequences:** Smaller image size, improved security posture, but more complex Dockerfile and potential debugging challenges

The operator container requires careful dependency management to ensure consistent behavior across environments. All required certificates, configuration schemas, and runtime dependencies must be included in the image or mounted at runtime.

Container Component	Purpose	Location	Management Strategy
Operator binary	Main application executable	/usr/local/bin/operator	Compiled during image build
CA certificates	TLS certificate validation	/etc/ssl/certs/	Included in base image or mounted
CRD schemas	OpenAPI validation schemas	/etc/operator/crds/	Embedded in image or mounted as ConfigMap
Default configurations	Fallback operator settings	/etc/operator/config/	Embedded with environment variable overrides
Health check scripts	Liveness and readiness probes	/usr/local/bin/health	Embedded shell scripts or HTTP endpoints

Helm Chart Organization

Helm charts provide templating and configuration management for operator deployments, enabling environment-specific customization while maintaining consistency across clusters. Chart organization affects maintainability, user experience, and upgrade reliability.

Chart Component	Purpose	Template Location	Configuration Source
Deployment	Operator pod specification	<code>templates/deployment.yaml</code>	<code>values.yaml</code> replicas, image, resources
RBAC resources	ServiceAccount, Role, RoleBinding	<code>templates/rbac.yaml</code>	<code>values.yaml</code> rbac.create, serviceAccount.name
CRD definitions	Custom resource schemas	<code>templates/crds/</code>	Static files or <code>values.yaml</code> schema overrides
Webhook configuration	AdmissionWebhook registration	<code>templates/webhooks.yaml</code>	<code>values.yaml</code> webhook.enabled, certificate settings
Service definitions	Webhook and metrics endpoints	<code>templates/services.yaml</code>	<code>values.yaml</code> service ports and types
ConfigMap resources	Operator configuration	<code>templates/configmap.yaml</code>	<code>values.yaml</code> config section

The Helm chart must handle upgrade scenarios gracefully, particularly when CRD schemas change or when webhook configurations are modified. CRD upgrades require special attention because Helm cannot handle CRD lifecycle management automatically.

Upgrade Scenario	Challenge	Helm Solution	Considerations
CRD schema changes	Helm doesn't update CRDs automatically	<code>crds/</code> directory for install-only, hooks for updates	May require manual <code>kubectl apply</code> for CRD updates
Webhook endpoint changes	API server caches webhook configuration	Rolling update with readiness probes	Temporary admission failures during rollout
RBAC permission changes	Existing pods use old ServiceAccount tokens	Restart operator pods after RBAC update	Plan for brief reconciliation pause
Image updates	Container registry authentication	imagePullSecrets configuration	Ensure registry access before update

Decision: Separate CRD Management from Operator Lifecycle

- **Context:** CRDs have different lifecycle requirements than the operator deployment, particularly around upgrades and versioning
- **Options Considered:**
 1. Include CRDs in Helm chart templates/ directory
 2. Use Helm CRDs/ directory for install-only management
 3. Separate CRD management with dedicated tooling
- **Decision:** Use Helm crds/ directory for installation with separate CRD upgrade procedures
- **Rationale:** Helm's CRD handling has limitations around updates, but storing CRDs in the chart ensures they're installed with the operator. Separate upgrade procedures provide better control.
- **Consequences:** Requires documented CRD upgrade procedures and coordination between chart updates and CRD changes

Environment Configuration Management

Operators must adapt to different environments (development, staging, production) without requiring image rebuilds. Configuration management strategies affect operational complexity, security, and deployment reliability.

Configuration Source	Use Cases	Examples	Security Considerations
Environment variables	Simple runtime settings	Log levels, feature flags, timeouts	Visible in pod specifications and process lists
ConfigMaps	Structured configuration files	JSON/YAML configs, template files	Visible to users with ConfigMap read permissions
Secrets	Sensitive configuration	Database passwords, API keys, certificates	Encrypted at rest, base64 encoded in transit
Command-line flags	Override settings	Debug modes, alternative config paths	Visible in process arguments and pod specs
File mounts	Large configuration files	Custom CRD schemas, policy definitions	Requires volume management and file watching

Configuration precedence establishes clear override behavior when multiple sources provide the same setting. A well-defined precedence order prevents configuration conflicts and enables predictable behavior across environments.

Precedence Level	Configuration Source	Override Scope	Example Usage
1 (Highest)	Command-line flags	Specific settings for debugging or testing	<code>--log-level=debug</code> for development environments
2	Environment variables	Runtime overrides for deployment-specific settings	<code>RECONCILE_TIMEOUT=30s</code> for slower clusters
3	ConfigMap mounted files	Environment-specific structured configuration	Different resource limits per environment
4 (Lowest)	Built-in defaults	Fallback values when no override is provided	Default replica counts, timeout values

⚠ Pitfall: Sensitive Information in Environment Variables

Avoid passing sensitive information like passwords or API keys through environment variables, as they're visible in pod specifications, process lists, and `kubectl describe` output. Use Kubernetes Secrets with volume mounts or the Secrets API instead.

High Availability and Leader Election

Mental Model: The Night Shift Manager System

Think of leader election like managing night shift supervisors at a 24/7 factory. You need multiple qualified supervisors available (multiple operator replicas), but only one should be actively making decisions at any time (leader election). If the active supervisor has an emergency and can't continue (pod failure), another supervisor should immediately take over (leadership transfer) without missing any critical tasks (reconciliation continuity). The transition should be seamless to workers on the factory floor (managed resources).

High availability ensures your operator continues managing workloads even during node failures, rolling updates, or unexpected crashes. Unlike stateless applications where you can simply run multiple replicas behind a load balancer, operators require coordination to prevent multiple controllers from conflicting with each other while reconciling the same resources.

Leader Election Mechanics

Leader election coordinates multiple operator replicas to ensure only one actively reconciles resources at any time. This coordination prevents race conditions, conflicting updates, and split-brain scenarios where multiple controllers attempt to manage the same resources simultaneously.

Election Component	Purpose	Storage Location	Update Behavior
Lock resource	Stores current leader identity and lease information	ConfigMap or Lease resource	Atomic updates with resource version checking
Leader identity	Identifies which replica currently holds the lock	Lock resource annotations	Pod name and namespace for leader identification
Lease duration	How long a leader can hold the lock without renewal	Lock resource spec	Configurable, typically 15-30 seconds
Renew deadline	How often the leader must update the lock	Controller configuration	Typically lease duration / 3
Retry period	How often non-leaders check for lock availability	Controller configuration	Typically renew deadline / 2

The leader election algorithm follows a simple but robust protocol that handles network partitions, process crashes, and clock skew gracefully. Understanding this protocol helps debug leadership issues and tune performance for different cluster characteristics.

Election State	Replica Behavior	Lock Resource Action	Transition Conditions
Startup	Attempt to acquire lock	Read current lock, try to claim if expired	Becomes leader if successful, follower if failed
Leader	Continuously renew lock, perform reconciliation	Update lock resource within renew deadline	Becomes follower if renewal fails
Follower	Watch lock status, attempt takeover if leader fails	Periodically check lock expiration	Becomes leader if lock expires and successfully claimed
Takeover	Attempting to claim expired lock	Atomic update with resource version check	Becomes leader if successful, remains follower if conflict

The critical insight here is that leader election provides coordination, not load balancing. Only one replica does the work while others remain on standby, unlike horizontally scaled stateless applications.

Leader Election Configuration

Proper leader election configuration balances failover speed against cluster load and stability. Aggressive timers provide faster failover but increase API server load and susceptibility to network glitches.

Configuration Parameter	Conservative Setting	Aggressive Setting	Trade-offs
Lease duration	60 seconds	15 seconds	Failover speed vs. stability during network issues
Renew deadline	20 seconds	5 seconds	Leader renewal frequency vs. API server load
Retry period	10 seconds	2 seconds	Lock acquisition speed vs. election overhead
Lock resource type	ConfigMap	Lease	API server compatibility vs. semantic clarity

Decision: Lease Resources with 15-Second Lease Duration

- **Context:** Leader election requires storage for coordination data and appropriate timing for failover scenarios
- **Options Considered:**
 1. ConfigMap with 60-second lease duration for conservative failover
 2. Lease resource with 15-second lease duration for faster failover
 3. Custom resource for lock storage with application-specific timing
- **Decision:** Use Kubernetes Lease resources with 15-second lease duration
- **Rationale:** Lease resources are designed specifically for leader election, reducing semantic confusion. 15-second duration provides reasonable failover speed without excessive API server load.
- **Consequences:** Requires Kubernetes 1.14+ for Lease API availability, but provides optimal balance of failover speed and stability

Failover and Recovery Behavior

When leadership changes occur, the operator must handle the transition gracefully to maintain reconciliation continuity. Both outgoing and incoming leaders have responsibilities during the transition period.

Failover Scenario	Detection Mechanism	Recovery Actions	Expected Downtime
Graceful shutdown	Leader releases lock before termination	Immediate lock release, follower takeover	0-2 retry periods
Process crash	Lock renewal timeout	Automatic lock expiration, follower takeover	1 lease duration + retry period
Network partition	Leader cannot reach API server for renewal	Lock expires, follower on connected side takes over	1 lease duration + retry period
Node failure	Lock renewal timeout + pod rescheduling	Lock expiration, new pod startup, leadership claim	Pod scheduling time + lease duration

The operator must handle leadership transitions without losing reconciliation state or creating resource conflicts. This requires careful design of the reconciliation loop and state management.

Transition Phase	Leader Responsibilities	Follower Responsibilities	State Considerations
Leadership loss	Stop reconciliation immediately, release in-flight operations	Monitor lock status, prepare for takeover	Avoid partial updates that could conflict
Leadership acquisition	Start reconciliation loop, scan for pending work	Continue monitoring for future leadership opportunities	Resume from cluster state, not in-memory state
Steady state	Maintain lock renewal, perform normal reconciliation	Stay ready for immediate takeover	Keep reconciliation idempotent and stateless

⚠ Pitfall: Continuing Reconciliation After Losing Leadership

A common mistake is failing to stop reconciliation immediately when leadership is lost. This can cause multiple controllers to simultaneously modify the same resources, leading to conflicts, inconsistent state, and resource thrashing. Always check leadership status before performing any write operations.

Monitoring and Observability for HA Deployments

High availability deployments require additional monitoring to track leadership status, failover events, and replica health. This observability helps diagnose performance issues and verify that failover mechanisms work correctly.

Metric Category	Key Metrics	Purpose	Alert Conditions
Leadership status	Current leader identity, lease age, renewal success rate	Track which replica is active and leadership stability	No active leader for >lease duration
Failover events	Leadership transitions, takeover frequency, transition duration	Monitor failover reliability and performance	Frequent leadership changes (>1/hour)
Reconciliation health	Reconciliation rate, error rate, queue depth per replica	Verify only leader is working, followers are idle	Non-leader performing reconciliation
Resource utilization	CPU, memory, API server requests per replica	Optimize resource allocation and detect resource leaks	High resource usage on follower replicas

The operator should expose Prometheus metrics for leadership status and reconciliation health. These metrics enable alerting on split-brain conditions, failed failovers, and performance degradation.

Metric Name	Type	Labels	Description
<code>operator_leader_election_status</code>	Gauge	<code>instance</code> , <code>leader</code>	1 if this replica is leader, 0 if follower
<code>operator_leadership_changes_total</code>	Counter	<code>instance</code>	Total number of leadership transitions for this replica
<code>operator_lease_renewal_duration_seconds</code>	Histogram	<code>instance</code>	Time taken to renew leadership lease
<code>operator_reconciliation_rate</code>	Gauge	<code>instance</code> , <code>resource_type</code>	Number of reconciliations per second (should be 0 for followers)

Common Pitfalls

⚠ Pitfall: Overly Broad RBAC Permissions

Granting `cluster-admin` or wildcard permissions (`*`) makes security auditing impossible and creates significant security risks. If an operator with broad permissions is compromised, attackers gain access to the entire cluster. Instead, enumerate specific permissions required for each reconciliation action and grant only those permissions.

⚠ Pitfall: Missing Resource Version Conflicts in Leader Election

Failing to handle resource version conflicts during leader election can cause split-brain scenarios where multiple replicas believe they are the leader. Always use conditional updates with resource version checking when claiming or renewing leadership locks.

⚠ Pitfall: Hardcoded Configuration in Container Images

Embedding environment-specific configuration (URLs, credentials, timeouts) in container images prevents image reuse across environments and requires image rebuilds for configuration changes. Use environment variables, ConfigMaps, and Secrets for all environment-specific settings.

⚠ Pitfall: Ignoring Webhook Certificate Expiration

Self-signed certificates or manually provisioned certificates will eventually expire, causing webhook failures that prevent resource creation and updates. Implement certificate monitoring and automatic renewal, preferably using cert-manager or similar automated certificate management systems.

⚠ Pitfall: Leader Election with Inappropriate Timing

Setting lease durations too short causes leadership thrashing during network hiccups, while setting them too long delays failover during actual failures. Test leader election behavior under various network conditions and adjust timing parameters based on your cluster characteristics and recovery time objectives.

Implementation Guidance

This section provides concrete implementation details for deploying the Database Operator with proper RBAC configuration, container packaging, and high availability setup.

Technology Recommendations

Component	Simple Option	Advanced Option
Container Build	<code>docker build</code> with multi-stage Dockerfile	Buildpacks with Cloud Native Buildpacks
Image Registry	Docker Hub public registry	Private registry with Harbor or ECR
RBAC Management	Static YAML manifests with <code>kubectl apply</code>	Helm charts with templated RBAC
Certificate Management	Self-signed certificates with OpenSSL	cert-manager with Let's Encrypt or private CA
Configuration Management	Environment variables with ConfigMaps	Helm values with environment-specific overrides
Deployment Strategy	<code>kubectl apply</code> with manual manifests	GitOps with ArgoCD or Flux

Recommended File Structure

```
operator-deployment/
├── docker/
│   ├── Dockerfile           ← Multi-stage build definition
│   └── .dockerignore         ← Exclude unnecessary files from build context
├── helm/
│   ├── Chart.yaml           ← Helm chart metadata
│   ├── values.yaml           ← Default configuration values
│   ├── values-dev.yaml       ← Development environment overrides
│   ├── values-prod.yaml     ← Production environment overrides
│   └── templates/
│       ├── deployment.yaml   ← Operator pod specification
│       ├── rbac.yaml          ← ServiceAccount, Role, RoleBinding
│       ├── webhooks.yaml       ← Admission webhook configuration
│       ├── service.yaml        ← Webhook and metrics services
│       └── configmap.yaml      ← Operator configuration
└── crds/
    └── databases.yaml        ← CRD definitions (install-only)
├── manifests/
│   ├── namespace.yaml        ← Operator namespace
│   ├── crd.yaml              ← Raw CRD for kubectl apply
│   └── rbac.yaml             ← Minimal RBAC for development
└── scripts/
    ├── build.sh               ← Container build script
    ├── deploy.sh              ← Helm deployment script
    └── undeploy.sh            ← Cleanup script
```

Complete RBAC Configuration

```
# ServiceAccount for operator pods                                                 YAML

apiVersion: v1

kind: ServiceAccount

metadata:
  name: database-operator-controller-manager
  namespace: database-operator-system

---

# ClusterRole with minimal required permissions

apiVersion: rbac.authorization.k8s.io/v1

kind: ClusterRole

metadata:
  name: database-operator-manager-role

rules:
  - apiGroups:
    - ""
      resources:
        - configmaps
        - events
        - persistentvolumeclaims
        - secrets
        - services
      verbs:
        - create
        - delete
        - get
        - list
        - patch
        - update
        - watch
    - apiGroups:
      - apps
      resources:
        - deployments
```

```
verbs:
- create
- delete
- get
- list
- patch
- update
- watch

- apiGroups:
- example.com

resources:
- databases

verbs:
- create
- delete
- get
- list
- patch
- update
- watch

- apiGroups:
- example.com

resources:
- databases/status

verbs:
- get
- patch
- update

- apiGroups:
- coordination.k8s.io

resources:
- leases

verbs:
```

```
- create
- get
- list
- update
---
# ClusterRoleBinding linking ServiceAccount to permissions

apiVersion: rbac.authorization.k8s.io/v1

kind: ClusterRoleBinding

metadata:
  name: database-operator-manager-rolebinding

roleRef:
  apiGroup: rbac.authorization.k8s.io
  kind: ClusterRole
  name: database-operator-manager-role

subjects:
- kind: ServiceAccount
  name: database-operator-controller-manager
  namespace: database-operator-system
```

Complete Deployment Configuration

```
apiVersion: apps/v1                                     YAML

kind: Deployment

metadata:
  name: database-operator-controller-manager
  namespace: database-operator-system
  labels:
    app.kubernetes.io/name: database-operator
    app.kubernetes.io/component: manager

spec:
  replicas: 2 # High availability with leader election
  selector:
    matchLabels:
      app.kubernetes.io/name: database-operator
      app.kubernetes.io/component: manager
  template:
    metadata:
      labels:
        app.kubernetes.io/name: database-operator
        app.kubernetes.io/component: manager
    spec:
      serviceAccountName: database-operator-controller-manager
      securityContext:
        runAsNonRoot: true
        runAsUser: 65532 # Distroless nonroot user
      containers:
        - name: manager
          image: database-operator:latest
          command:
            - /manager
          args:
            - --leader-elect
            - --leader-elect-lease-duration=15s
            - --leader-elect-renew-deadline=5s
```

```
- --leader-elect-retry-period=2s

ports:
  - containerPort: 9443
    name: webhook-server
    protocol: TCP
  - containerPort: 8080
    name: metrics
    protocol: TCP
  - containerPort: 8081
    name: health-probe
    protocol: TCP

env:
  - name: POD_NAMESPACE
    valueFrom:
      fieldRef:
        fieldPath: metadata.namespace
  - name: POD_NAME
    valueFrom:
      fieldRef:
        fieldPath: metadata.name

livenessProbe:
  httpGet:
    path: /healthz
    port: 8081
  initialDelaySeconds: 15
  periodSeconds: 20

readinessProbe:
  httpGet:
    path: /readyz
    port: 8081
  initialDelaySeconds: 5
  periodSeconds: 10

resources:
```

```
limits:
  cpu: 500m
  memory: 128Mi
requests:
  cpu: 10m
  memory: 64Mi
securityContext:
  allowPrivilegeEscalation: false
  capabilities:
    drop:
      - ALL
  readOnlyRootFilesystem: true
volumeMounts:
  - mountPath: /tmp/k8s-webhook-server/serving-certs
    name: cert
    readOnly: true
volumes:
  - name: cert
    secret:
      defaultMode: 420
      secretName: webhook-server-certs
terminationGracePeriodSeconds: 10
```

Core Leader Election Implementation

```
// Manager setup with leader election enabled

func main() {

    var enableLeaderElection bool

    var leaseDuration time.Duration

    var renewDeadline time.Duration

    var retryPeriod time.Duration


    flag.BoolVar(&enableLeaderElection, "leader-elect", false,
        "Enable leader election for controller manager")

    flag.DurationVar(&leaseDuration, "leader-elect-lease-duration", 15*time.Second,
        "The duration that non-leader candidates will wait after observing a leadership renewal")

    flag.DurationVar(&renewDeadline, "leader-elect-renew-deadline", 5*time.Second,
        "The interval between attempts by the acting master to renew a leadership slot")

    flag.DurationVar(&retryPeriod, "leader-elect-retry-period", 2*time.Second,
        "The duration the clients should wait between attempting acquisition and renewal")

    flag.Parse()

    mgr, err := ctrl.NewManager(ctrl.GetConfigOrDie(), ctrl.Options{
        Scheme:                 scheme,
        MetricsBindAddress:     ":8080",
        Port:                   9443,
        HealthProbeBindAddress: ":8081",
        LeaderElection:          enableLeaderElection,
        LeaderElectionID:        "database.example.com",
        LeaseDuration:           &leaseDuration,
        RenewDeadline:           &renewDeadline,
        RetryPeriod:             &retryPeriod,
        LeaderElectionResourceLock: "leases",
        LeaderElectionNamespace: os.Getenv("POD_NAMESPACE"),
    })

    if err != nil {
        setupLog.Error(err, "unable to start manager")
        os.Exit(1)
    }
}
```

GO

```

}

// TODO: Register controllers and webhooks with manager

// TODO: Add health and readiness checks

// TODO: Start manager with graceful shutdown on SIGTERM

}

```

Container Build Configuration

```

# Multi-stage build for minimal production image                                DOCKERFILE

FROM golang:1.21 AS builder

WORKDIR /workspace

# Copy go mod files for dependency caching

COPY go.mod go.sum ./

RUN go mod download

# Copy source code

COPY cmd/ cmd/
COPY internal/ internal/
COPY api/ api/

# Build the operator binary

RUN CGO_ENABLED=0 GOOS=linux GOARCH=amd64 go build -a -o manager cmd/main.go

# Production image with minimal attack surface

FROM gcr.io/distroless/static:nonroot

WORKDIR /

COPY --from=builder /workspace/manager .

USER 65532:65532

ENTRYPOINT ["/manager"]

```

Milestone Checkpoint

After implementing the deployment configuration:

1. **RBAC Verification:** Deploy the operator and verify it can perform required operations:

```
kubectl auth can-i create databases --as=system:serviceaccount:database-operator-system:database-operator-controller-manager
BASH

kubectl auth can-i update secrets --as=system:serviceaccount:database-operator-system:database-operator-controller-manager
```

2. **Leader Election Testing:** Deploy multiple replicas and verify only one performs reconciliation:

```
kubectl scale deployment database-operator-controller-manager --replicas=3
BASH

kubectl logs -l app.kubernetes.io/name=database-operator -f | grep "leader election"
```

3. **Failover Validation:** Delete the leader pod and verify another replica takes over within the lease duration.

4. **Expected Behavior:**

- Only one replica should show "Starting EventSource" and "Starting Controller" in logs
- Other replicas should show "Waiting for leader election" or similar
- After leader pod deletion, a new leader should be elected within 15-20 seconds
- Reconciliation should continue without interruption during leadership transitions

Debugging Tips

Symptom	Likely Cause	How to Diagnose	Fix
"Forbidden" errors during reconciliation	Missing RBAC permissions	Check <code>kubectl auth can-i</code> for required operations	Add missing permissions to ClusterRole
Multiple replicas reconciling simultaneously	Leader election not working	Check logs for "acquired lease" messages from multiple pods	Verify leader election configuration and lease resource creation
Webhook admission failures	Certificate issues or service misconfiguration	Check webhook service endpoints and certificate validity	Verify cert-manager setup or regenerate self-signed certificates
Operator not starting after deployment	Image pull failures or resource constraints	Check pod events and resource utilization	Verify image registry access and resource requests
Slow failover during leader crashes	Lease duration too long	Monitor leadership transition timing	Reduce lease duration while avoiding election thrashing

Debugging Guide

Milestone(s): Milestone 2 (Controller Setup), Milestone 3 (Reconciliation Loop), Milestone 4 (Webhooks), Milestone 5 (Testing & Deployment) - provides troubleshooting guidance for common issues across controller implementation, webhook configuration, and RBAC setup

Building a Kubernetes operator involves complex interactions between custom controllers, admission webhooks, informers, workqueues, and the Kubernetes API server. When these components don't work as expected, the symptoms can be confusing and the root causes non-obvious. This section provides a systematic debugging guide organized around the three most common problem areas developers encounter when building operators.

Mental Model: The Detective's Toolkit - Think of debugging an operator like being a detective investigating a complex case. You have multiple witnesses (logs), physical evidence (resource states), and timelines (event sequences). Just as a detective follows leads and eliminates suspects, debugging operators requires gathering evidence from multiple sources, forming hypotheses about root causes, and systematically testing those theories. The key is knowing what evidence to look for and how to interpret the clues.

The debugging process typically follows this pattern: observe the symptoms (what's not working), gather evidence from logs and cluster state, form hypotheses about potential causes, test each hypothesis with targeted debugging techniques, and implement fixes. This guide provides the diagnostic toolkit and investigative procedures for each category of problems.

Operator debugging is particularly challenging because problems often manifest as cascading failures. A simple RBAC misconfiguration can cause controller crashes, which leads to failed reconciliation, which results in outdated status conditions, which confuses end users about the actual problem. The key is learning to trace problems back to their root causes rather than treating symptoms.

Controller and Reconciliation Issues

Controller and reconciliation problems are among the most common issues operators face. These manifest as resources stuck in pending states, infinite reconciliation loops, or controllers that appear to be "doing nothing" despite resource changes. Understanding the controller's internal state machine and information flow is crucial for effective debugging.

Mental Model: The Broken Assembly Line - Think of the controller as an assembly line with several stations: the informer (receives raw materials), the workqueue (buffers work items), and the reconcile function (processes each item). When the assembly line breaks down, products either pile up at one station, get stuck in processing, or come out defective. By examining each station's health and throughput, you can identify where the bottleneck or failure occurs.

Reconciliation Loop Problems

Reconciliation loops can fail in several distinct patterns, each with characteristic symptoms and root causes. The most common failure modes involve infinite loops, stalled reconciliation, and partial reconciliation failures that leave resources in inconsistent states.

Infinite Reconciliation Loops

Infinite reconciliation loops occur when the `Reconcile` function continuously requeues the same resource without making progress toward the desired state. This typically happens when the reconciler incorrectly determines that the current state doesn't match the desired state, even though they are actually equivalent.

Symptom	Evidence	Root Cause	Diagnostic Method	Solution
High CPU usage on controller pods	Metrics show constant reconciliation rate	Spec comparison logic treats equivalent states as different	Add debug logging to state comparison logic	Fix comparison logic to handle semantic equivalence
Resource status never reaches "Ready"	Status shows rapid condition updates	Reconciler updates owned resources on every loop	Log owned resource hashes before/after comparison	Implement proper resource comparison with strategic merge
Exponentially growing log volume	Same resource appears repeatedly in logs	Missing break condition in reconciliation logic	Trace reconciliation decision tree for specific resource	Add idempotency checks and proper exit conditions
API server rate limiting errors	Controller logs show 429 responses	Excessive API calls from repeated reconciliation	Monitor API call patterns with request counting	Implement client-side rate limiting and backoff

⚠ Pitfall: Resource Comparison Without Strategic Merge Many infinite loops result from comparing resources without accounting for Kubernetes' strategic merge patch behavior. When you create a `Deployment` with `replicas: 3`, Kubernetes may add default fields like `progressDeadlineSeconds: 600`. If your reconciler compares the original spec (without defaults) to the current state (with defaults), it will always detect differences and attempt updates.

Stalled Reconciliation

Stalled reconciliation occurs when the controller stops processing resources entirely, even though new events are occurring. This typically indicates problems with the informer cache, workqueue processing, or controller startup sequence.

Symptom	Evidence	Root Cause	Diagnostic Method	Solution
Resources stuck in "Pending" state	No reconciliation logs for affected resources	Informer cache not synced before processing started	Check informer sync status in controller startup	Add proper cache sync wait before starting workers
Controller appears healthy but ignores changes	Resource updates don't trigger reconciliation	Event handlers not properly registered	Verify SetupWithManager includes all watch configurations	Add missing watch statements for owned resources
Some resources reconcile, others don't	Selective resource processing failures	Workqueue filtering or malformed resource keys	Examine workqueue contents and key generation logic	Fix resource key generation and queue filtering
Fresh resources process, old ones don't	Controller restart resolves some stuck resources	Informer resync period too long or disabled	Check resync configuration and cache staleness	Configure appropriate resync period or force cache refresh

Partial Reconciliation Failures

Partial reconciliation failures occur when the reconciler successfully processes some owned resources but fails on others, leaving the system in an inconsistent state. This is particularly dangerous because the overall status may appear healthy while critical components are actually broken.

Example Failure Sequence:

1. Database resource specifies: replicas=3, version=14.5
2. Reconciler successfully creates StatefulSet with 3 replicas
3. Reconciler fails to create ConfigMap with version 14.5 config
4. Status shows "Ready: True" based on StatefulSet status
5. Pods start but fail to launch due to missing ConfigMap
6. User sees "Ready" status but database is actually down

The challenge with partial failures is that the reconciler must decide whether to report success (some resources are healthy) or failure (the overall system is broken). The solution is implementing transactional reconciliation with proper rollback and status reporting.

Informer and Cache Problems

The informer and cache system provides the controller with a local, eventually consistent view of cluster state. When this system malfunctions, controllers may operate on stale data, miss important events, or crash due to unexpected cache states.

Cache Synchronization Issues

Cache synchronization problems occur when the local informer cache diverges from the actual API server state, leading to reconciliation decisions based on outdated information.

Problem Type	Symptoms	Detection Method	Resolution
Stale cache data	Reconciler operates on deleted resources	Compare cache contents with direct API calls	Increase resync period or force cache refresh
Missing recent updates	Controller ignores recent resource changes	Check event sequence timestamps vs cache timestamps	Verify watch connection health and restart if needed
Cache memory pressure	High memory usage, frequent GC pauses	Monitor informer cache size and growth patterns	Implement cache size limits or namespace filtering
Watch connection drops	Periodic reconciliation gaps in logs	Monitor watch connection errors and reconnections	Add connection monitoring and exponential backoff

⚠️ Pitfall: Starting Workers Before Cache Sync A common mistake is starting workqueue processing goroutines before the informer cache has fully synchronized with the API server. This causes the controller to process events based on incomplete cache data, leading to incorrect reconciliation decisions and potential data corruption.

Event Handler Registration Problems

Event handlers translate informer cache events (add, update, delete) into workqueue items for reconciliation. Misconfigured event handlers can cause the controller to miss important events or process irrelevant changes.

Handler Issue	Behavior	Investigation	Fix
Missing update handlers	Resource changes ignored	Verify all SetupWithManager watch configurations	Add missing Controller.Watches() calls
Incorrect filtering	Wrong resources trigger reconciliation	Log event handler invocations with resource details	Fix OwnerReference or label selectors in watches
Handler panics	Controller crashes on specific events	Check event handler code for nil pointer dereferences	Add nil checks and error handling in handlers
Duplicate event processing	Same resource reconciled multiple times per change	Monitor workqueue items and deduplication	Ensure proper resource key generation and dedup logic

Workqueue and Rate Limiting Issues

The workqueue buffers reconciliation requests and implements rate limiting to prevent overwhelming the controller or API server. Workqueue problems typically manifest as delayed processing, resource starvation, or memory growth from unbounded queues.

Rate Limiting and Backoff Problems

Rate limiting protects the system from overload but can cause delays when misconfigured. Understanding the backoff algorithms and tuning parameters is essential for optimal controller performance.

The default controller-runtime rate limiter uses exponential backoff with jitter: initial delay of 5ms, doubling on each retry, up to a maximum of 16 minutes. For transient errors, this provides good protection. However, permanent errors (like RBAC denials) will continue consuming rate limiter capacity without making progress.

Rate Limiting Issue	Observable Effects	Tuning Approach	Implementation
Aggressive backoff delays critical updates	Important resources take minutes to reconcile	Reduce base delay and maximum interval	Configure custom ItemBasedRateLimiter with lower bounds
Rate limiter overwhelmed by error resources	Healthy resources delayed behind failing ones	Separate rate limiting for different error types	Implement multi-tier rate limiting based on error classification
Memory growth from queued items	Controller memory usage grows without bound	Add queue depth monitoring and circuit breakers	Set maximum queue length and reject new items when full
Hot resource monopolizes queue capacity	Single misbehaving resource blocks others	Implement per-resource rate limiting	Use BucketRateLimiter with separate buckets per resource

Workqueue Shutdown and Goroutine Leaks

Proper workqueue shutdown is critical for clean controller termination and preventing goroutine leaks in testing scenarios. Shutdown problems often manifest as hanging tests or controllers that can't be cleanly restarted.

The workqueue shutdown sequence must be carefully orchestrated: stop accepting new items, drain existing items, signal worker goroutines to exit, and wait for all workers to complete. Skipping any step can cause deadlocks or resource leaks.

Proper Shutdown Sequence:

1. Call `workqueue.ShutDown()` to stop accepting new items
2. Worker goroutines check `Done()` channel and exit gracefully
3. Call `workqueue.ShutDownWithDrain()` to process remaining items
4. Use `sync.WaitGroup` to wait for all workers to finish
5. Close any additional channels or cleanup resources

Webhook and Admission Problems

Admission webhooks introduce additional complexity because they operate synchronously in the request path between the client and API server. When webhooks malfunction, they can block all resource operations, cause timeouts, or silently fail validation. Webhook debugging requires understanding HTTPS certificate management, admission review processing, and failure policy configuration.

Mental Model: The Security Checkpoint - Think of admission webhooks like airport security checkpoints. Every passenger (resource) must pass through security (webhook validation) before boarding the plane (being stored in etcd). If security is down, flights get delayed or canceled. If security is too slow, passengers miss connections. If security has wrong information, innocent passengers get detained while threats slip through. The key is ensuring the checkpoint is fast, reliable, and has current information.

Certificate and TLS Issues

Admission webhooks must serve HTTPS traffic with valid TLS certificates that the API server can verify. Certificate problems are among the most common webhook issues and can be particularly frustrating because they often work in development but fail in production due to different certificate authorities or hostname verification.

Certificate Provisioning and Validation

The API server must trust the webhook's TLS certificate and verify that the certificate's subject matches the service hostname. This requires proper certificate generation, distribution, and rotation mechanisms.

Certificate Problem	Symptoms	Verification Steps	Solution Approach
Self-signed certificate not trusted	Webhook admission review requests fail with TLS errors	Check API server logs for certificate validation errors	Use cert-manager or add CA bundle to webhook configuration
Hostname mismatch in certificate	TLS handshake failures despite valid certificates	Verify certificate SAN includes service DNS names	Regenerate certificate with correct Subject Alternative Names
Expired certificates	Intermittent webhook failures, particularly after time jumps	Check certificate expiration dates with openssl	Implement automatic certificate rotation with cert-manager
Certificate not readable by webhook pod	Webhook server fails to start with TLS initialization errors	Verify certificate file permissions and volume mounts	Fix service account permissions and secret volume configuration

The certificate Subject Alternative Name (SAN) field must include all DNS names that the API server might use to reach the webhook service. For a webhook service named `database-webhook` in namespace `operator-system`, the SAN should include: `database-webhook.operator-system.svc`, `database-webhook.operator-system.svc.cluster.local`, and potentially the service IP address.

⚠ Pitfall: Certificate Rotation During Webhook Processing Certificate rotation can cause admission request failures if certificates change while the API server still has cached connections. The webhook server must implement graceful certificate reloading without dropping active connections, and the API server may need time to refresh its certificate cache.

TLS Configuration and Cipher Suites

The webhook server must support TLS versions and cipher suites that the API server accepts. Mismatched TLS configurations can cause connection failures that are difficult to diagnose because they occur during the initial handshake.

Modern Kubernetes API servers typically require TLS 1.2 or higher and prefer Forward Secrecy cipher suites. The webhook server should support a compatible set of protocols and ciphers while maintaining reasonable security standards.

Admission Review Processing

Admission review processing involves unmarshaling the admission request, validating or mutating the resource, and returning a properly formatted admission response. Problems in this pipeline can cause webhook timeouts, malformed responses, or incorrect admission decisions.

Request Deserialization and Validation

The admission review request contains the resource data in JSON format, along with metadata about the operation type, user information, and admission context. Proper deserialization requires handling different API versions, unknown fields, and malformed input.

Processing Issue	Error Manifestation	Debugging Approach	Implementation Fix
Unsupported admission review version	Webhook returns 400 Bad Request	Log admission review version in request handler	Add support for multiple admissionregistration API versions
Resource deserialization failures	Webhook denies all requests with decode errors	Log raw admission request JSON for analysis	Add proper error handling and version-aware deserialization
Missing required fields in admission request	Webhook crashes or returns malformed responses	Validate admission request structure before processing	Add comprehensive input validation with clear error messages
Memory exhaustion from large resources	Webhook becomes unresponsive or crashes	Monitor memory usage during large resource processing	Implement request size limits and streaming deserialization

Mutation Response Formatting

Mutating webhooks must return JSON Patch operations that describe the changes to apply to the resource. Malformed patches can cause admission failures or unintended modifications.

JSON Patch operations must specify the exact path to modify, the operation type (add, replace, remove), and the new value with correct typing. Path specifications use JSON Pointer syntax, which has specific escaping rules for field names containing special characters.

Example Mutation Scenarios:

- Adding a default value to an empty field: `{"op": "add", "path": "/spec/replicas", "value": 1}`
- Replacing an existing array element: `{"op": "replace", "path": "/spec/containers/0/image", "value": "new-image:v2"}`
- Removing a field: `{"op": "remove", "path": "/metadata/annotations/old-annotation"}`
- Adding to an array: `{"op": "add", "path": "/spec/containers/-", "value": {...}}`

⚠ Pitfall: JSON Patch Path Escaping JSON Pointer paths require escaping special characters in field names. A field named `example.com/annotation` must be referenced as `/metadata/annotations/example.com~1annotation` (where `~1` represents the `/` character). Incorrect escaping causes patch application failures.

Webhook Timeout and Performance Issues

Admission webhooks operate in the synchronous request path, so performance problems directly impact user operations. The default admission timeout is 10 seconds, after which the API server either admits or denies the request based on the configured failure policy.

Response Time Optimization

Webhook response time affects the user experience for all resource operations. Slow webhooks make `kubectl apply` commands feel sluggish and can cause client timeouts on large batch operations.

Common performance bottlenecks include: external API calls during validation (such as checking license servers or external databases), complex validation logic with nested loops, large resource deserialization, and cold start delays in serverless environments.

Performance Issue	User Impact	Measurement Method	Optimization Strategy
Slow validation logic	kubectl commands take several seconds	Add timing metrics to webhook handlers	Cache validation results and optimize algorithms
External API dependencies	Webhooks fail when external services are unavailable	Monitor external service response times	Implement circuit breakers and local caching
Large resource processing	Memory pressure and slow response for big resources	Profile memory allocation during large resource handling	Stream processing and implement size limits
Cold start delays	First requests after deployment are very slow	Measure time-to-first-response after pod startup	Add readiness probes and pre-warm caches

Failure Policy and Degraded Mode Operation

Webhook failure policies determine what happens when the webhook is unavailable, times out, or returns errors. The two options are `Fail` (deny admission on webhook failure) and `Ignore` (allow admission on webhook failure). Both have significant operational implications.

A `Fail` policy provides security guarantees by ensuring no unvalidated resources enter the cluster, but it creates availability risks if the webhook becomes unavailable. An `Ignore` policy maintains cluster availability during webhook outages but allows potentially invalid resources to be created.

Decision: Webhook Failure Policy Selection

- **Context:** Webhook failures can either block all resource operations (Fail policy) or allow unvalidated resources (Ignore policy)
- **Options Considered:** Always Fail, Always Ignore, Dynamic policy based on validation criticality
- **Decision:** Use Fail policy for validating webhooks, Ignore policy for mutating webhooks with non-critical defaults
- **Rationale:** Validation failures indicate policy violations that should block admission. Mutation failures for defaults are less critical than availability
- **Consequences:** Validation webhook outages block resource creation, requiring high availability deployment. Mutation webhook outages allow resources without defaults but maintain system operation

RBAC and Permissions Debugging

Role-Based Access Control (RBAC) governs what operations the operator's service account can perform against the Kubernetes API. RBAC problems typically manifest as "permission denied" errors, but diagnosing the specific missing permissions and configuring minimal privilege access requires understanding Kubernetes' multi-layered permission model.

Mental Model: The Credential Hierarchy - Think of Kubernetes RBAC like a corporate security system with keycards and access levels. Each service account has a keycard (identity), roles define what areas each keycard can access (permissions), and role bindings connect keycards to access levels (authorization). When access is denied, you need to trace through the entire chain: does the keycard exist, does the role have the right permissions, is the binding connecting them correctly, and is the request happening in the right location (namespace).

Service Account and Role Configuration

Service accounts provide identity for pods, while roles and role bindings define and grant permissions. The interaction between these components creates a complex permission matrix that can be difficult to debug when misconfigured.

Service Account Identity Issues

Service accounts must exist and be properly mounted into pods before they can be used for API authentication. Missing or misconfigured service accounts cause immediate authentication failures.

Service Account Problem	Error Pattern	Investigation Steps	Resolution Path
Service account doesn't exist	HTTP 401 Unauthorized errors on all API calls	Check if service account exists in the correct namespace	Create service account with proper metadata
Service account token not mounted	Authentication works in pod but fails for specific operations	Verify automountServiceAccountToken is not disabled	Enable service account token mounting or use manual token
Wrong service account specified in deployment	Pod uses default service account instead of operator account	Compare Deployment serviceAccountName with actual pod spec	Update Deployment to reference correct service account
Service account in wrong namespace	Permission denied errors despite correct role bindings	Verify service account and role binding namespaces match	Move service account or update role binding namespace

Role Definition and Scope Problems

Roles define the specific permissions (verbs and resources) that can be granted to service accounts. Incorrect role definitions are a common source of permission errors, particularly when operators need to manage multiple resource types or perform administrative operations.

The principle of least privilege suggests granting only the minimum permissions necessary for operation. However, operators often require broad permissions across multiple resources, creating tension between security and functionality.

Permission Issue	Manifestation	Diagnosis Technique	Correction Method
Missing resource permissions	Permission denied on specific resource types	Check role definition against error messages	Add required resources to role rules
Insufficient verbs for operations	Create works but update/delete fails	Map operation types to required verbs (get, list, watch, create, update, patch, delete)	Expand verb list in role rules
API group mismatches	Permission denied for custom resources	Verify apiGroups in role match CRD group	Update role apiGroups to include custom resource groups
Namespace scope limitations	Cluster-wide operations fail with namespace-scoped roles	Determine if operator needs cluster-wide or namespace-scoped access	Convert Role to ClusterRole or add namespace-specific roles

⚠️ Pitfall: Custom Resource API Groups in RBAC

Custom resources use different API groups than built-in Kubernetes resources. A `Database` custom resource in the `database.example.com` API group requires RBAC rules with `apiGroups`:

`["database.example.com"]`, not the empty string used for core resources. Forgetting to include the custom API group causes permission denied errors for all custom resource operations.

Permission Scope and Binding Issues

Role bindings connect service accounts to roles, creating the actual permission grants. Binding problems can be subtle because they often involve namespace boundaries, subject references, or binding scope mismatches.

Namespace Boundary and Scope Mismatches

Kubernetes uses both namespace-scoped (Role, RoleBinding) and cluster-scoped (ClusterRole, ClusterRoleBinding) RBAC objects. Mismatching scopes between roles and bindings, or between operator needs and binding scope, causes permission failures that can be difficult to diagnose.

RBAC Scope Combinations:

- Role + RoleBinding: namespace-scoped permissions within a specific namespace
- ClusterRole + RoleBinding: cluster-wide permissions applied to a specific namespace
- ClusterRole + ClusterRoleBinding: cluster-wide permissions applied cluster-wide
- Role + ClusterRoleBinding: Invalid - ClusterRoleBindings cannot reference Roles

The most common scope mismatch occurs when operators need to watch resources across all namespaces but are configured with namespace-scoped role bindings. This causes the operator to successfully process resources in its own namespace while failing silently on resources in other namespaces.

Subject Reference and Binding Validation

Role bindings must correctly reference the service account using the proper subject format. Incorrect subject references cause the binding to be ineffective, even though it appears correctly configured.

Binding Problem	Symptoms	Validation Method	Fix Implementation
Incorrect subject name	Permission denied despite apparently correct bindings	Compare binding subjects with actual service account names	Update binding subject name to match service account
Wrong subject namespace	Binding exists but doesn't grant permissions	Verify binding subject namespace matches service account namespace	Update subject namespace in role binding
Invalid subject kind	Binding appears correct but permissions not granted	Ensure subject kind is "ServiceAccount" not "User" or "Group"	Correct subject kind in binding specification
Multiple conflicting bindings	Inconsistent permission behavior across operations	List all role bindings affecting the service account	Consolidate bindings and remove conflicting rules

Runtime Permission Verification

Runtime permission verification involves testing actual API operations to confirm that RBAC configuration grants the expected access. This is particularly important during operator development and deployment because RBAC misconfigurations often only surface when specific operations are attempted.

Permission Testing and Validation

The `kubectl auth can-i` command provides a way to test permissions without attempting actual operations. This is useful for validating RBAC configuration before deploying operators or troubleshooting permission issues in production.

Permission Testing Examples:

```
kubectl auth can-i create databases.database.example.com --as=system:serviceaccount:operator-system:database-controller
kubectl auth can-i update deployments --as=system:serviceaccount:operator-system:database-controller -n production
kubectl auth can-i list secrets --as=system:serviceaccount:operator-system:database-controller --all-namespaces
kubectl auth can-i delete persistentvolumeclaims --as=system:serviceaccount:operator-system:database-controller
```

These tests should be performed for all operations the operator needs to perform, across all relevant namespaces and resource types.

Minimal Privilege Validation

Operators should follow the principle of least privilege, granting only the minimum permissions necessary for correct operation. However, determining the exact minimum set of permissions requires careful analysis of all operator operations and their corresponding RBAC requirements.

The process involves: cataloging all API operations the operator performs (from controller reconciliation and webhook processing), mapping each operation to required RBAC verbs and resources, grouping permissions by functional area to create focused roles, and testing that the minimal permission set supports all operator functionality.

Operation Type	Required Permissions	Rationale	Validation Method
Custom resource reconciliation	get, list, watch, update, patch on custom resources	Controller needs to read current state and update status	Test reconciliation loop with minimal permissions
Owned resource management	get, list, watch, create, update, patch, delete on owned types	Controller creates and manages Deployments, Services, etc.	Test resource creation, updates, and cleanup
Event recording	create on events	Controller reports status and errors through events	Test event creation during normal and error scenarios
Leader election	get, create, update on leases or configmaps	High availability requires leader election coordination	Test controller startup with multiple replicas

The critical insight here is that RBAC debugging requires understanding the complete request context - not just the service account and permissions, but also the namespace scope, API group, and specific operation being performed. Each element must align correctly for authorization to succeed.

Implementation Guidance

This debugging toolkit provides practical techniques for diagnosing and resolving the most common operator development issues. The goal is enabling developers to quickly identify root causes and apply targeted fixes rather than spending hours on trial-and-error debugging.

Technology Recommendations:

Debugging Component	Simple Option	Advanced Option
Log Analysis	Basic kubectl logs with grep	Structured logging with logr and log aggregation
RBAC Validation	kubectl auth can-i commands	RBAC policy simulators and automated testing
Certificate Management	Manual openssl certificate validation	cert-manager with automated certificate lifecycle
Performance Monitoring	Basic kubectl top and describe commands	Prometheus metrics with Grafana dashboards
Integration Testing	Manual kubectl apply/delete cycles	Automated envtest with comprehensive scenarios

Recommended File Structure for Debugging Tools:

```

operator-project/
  cmd/
    debug-controller/
      main.go          ← debug utility for controller issues
    validate-rbac/
      main.go          ← RBAC permission validation tool
    webhook-tester/
      main.go          ← webhook admission request simulator
  internal/
    debug/
      controller_diagnostics.go   ← controller health checking utilities
      webhook_diagnostics.go     ← webhook validation and testing helpers
      rbac_validator.go          ← automated RBAC permission verification
  test/
    debug/
      controller_debug_test.go   ← debugging scenario test cases
      webhook_debug_test.go      ← webhook failure simulation tests

```

Complete Controller Diagnostics Helper:

```
// internal/debug/controller_diagnostics.go                                     GO

package debug

import (
    "context"
    "fmt"
    "time"

    "k8s.io/apimachinery/pkg/runtime"
    "k8s.io/client-go/tools/cache"
    "k8s.io/client-go/util/workqueue"
    "sigs.k8s.io/controller-runtime/pkg/client"
    "sigs.k8s.io/controller-runtime/pkg/log"
)

databasev1 "github.com/example/database-operator/api/v1"

)

type ControllerDiagnostics struct {

    Client client.Client

    Scheme *runtime.Scheme

    logger logr.Logger
}

type InformerHealth struct {

    CacheSynced     bool

    LastSyncTime    time.Time

    CacheSize       int

    WatchConnected  bool

    ResyncPeriod    time.Duration
}

type WorkqueueHealth struct {

    QueueLength     int

    ProcessingRate float64

    ErrorRate      float64

    RetryBacklog   int
}
```

```

    IsShuttingDown bool

}

type ReconcilerHealth struct {

    ActiveReconciliations int

    AverageReconcileTime time.Duration

    ErrorCount           int64

    SuccessCount         int64

    LastReconcileTime   time.Time

}

// DiagnoseInformerHealth checks the health of controller informers

func (cd *ControllerDiagnostics) DiagnoseInformerHealth(informer cache.SharedIndexInformer) InformerHealth {

    // TODO: Check if informer cache has synced with API server

    // TODO: Measure cache size and last sync time

    // TODO: Verify watch connection is active

    // TODO: Report resync period configuration

    // Hint: Use informer.HasSynced() and informer.LastSyncResourceVersion()

    return InformerHealth{}

}

// DiagnoseWorkqueueHealth analyzes workqueue performance and backlog

func (cd *ControllerDiagnostics) DiagnoseWorkqueueHealth(queue workqueue.RateLimitingInterface) WorkqueueHealth {

    // TODO: Measure current queue length and processing rate

    // TODO: Calculate error rate from retry attempts

    // TODO: Count items in retry backlog

    // TODO: Check if queue is in shutdown state

    // Hint: Use queue.Len() and implement metrics collection

    return WorkqueueHealth{}

}

// DiagnoseReconcilerHealth evaluates reconciliation performance

func (cd *ControllerDiagnostics) DiagnoseReconcilerHealth() ReconcilerHealth {

    // TODO: Count active reconciliation goroutines

    // TODO: Calculate average reconciliation duration from metrics
}

```

```

// TODO: Report success/error ratios

// TODO: Track last reconciliation timestamp

// Hint: Use prometheus metrics or custom counters

return ReconcilerHealth{}

}

// ValidateResourceState compares desired vs actual state for debugging

func (cd *ControllerDiagnostics) ValidateResourceState(ctx context.Context, db *databasev1.Database)
(*StateValidationResult, error) {

    // TODO: Fetch current state of all owned resources

    // TODO: Compare with desired state from Database spec

    // TODO: Identify specific differences and inconsistencies

    // TODO: Report resources that should exist but don't

    // TODO: Report resources that exist but shouldn't

    // Hint: Use owner references to find owned resources

    return nil, nil
}

// TraceReconciliationPath logs the complete reconciliation decision tree

func (cd *ControllerDiagnostics) TraceReconciliationPath(ctx context.Context, db *databasev1.Database) error {

    // TODO: Enable debug logging for this specific resource

    // TODO: Log each step of reconciliation logic

    // TODO: Record timing for each reconciliation phase

    // TODO: Log state comparisons and decisions

    // TODO: Track condition updates and status changes

    // Hint: Use structured logging with consistent field names

    return nil
}

```

Webhook Diagnostics and Testing Helper:

```
// internal/debug/webhook_diagnostics.go
```

GO

```
package debug
```

```
import (
```

```
    "context"
```

```
    "crypto/tls"
```

```
    "fmt"
```

```
    "net/http"
```

```
    "time"
```

```
        admissionv1 "k8s.io/api/admission/v1"
```

```
        metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
```

```
        "k8s.io/apimachinery/pkg/runtime"
```

```
        databasev1 "github.com/example/database-operator/api/v1"
```

```
)
```

```
type WebhookDiagnostics struct {
```

```
    WebhookURL string
```

```
    CertPath string
```

```
    KeyPath string
```

```
    HTTPClient *http.Client
```

```
}
```

```
type CertificateHealth struct {
```

```
    IsValid bool
```

```
    ExpirationTime time.Time
```

```
    DaysUntilExpiry int
```

```
    SubjectAltNames []string
```

```
    Issuer string
```

```
    ValidationError error
```

```
}
```

```
type AdmissionTestResult struct {
```

```
    RequestSent bool
```

```
    ResponseReceived bool
```

```

    ResponseTime      time.Duration

    AdmissionDecision string

    ValidationErrors []string

    MutationPatches  []JSONPatch

}

// DiagnoseCertificateHealth validates webhook TLS certificates

func (wd *WebhookDiagnostics) DiagnoseCertificateHealth() CertificateHealth {

    // TODO: Load certificate from file or Kubernetes secret

    // TODO: Verify certificate is not expired

    // TODO: Check Subject Alternative Names include service DNS names

    // TODO: Validate certificate chain and CA trust

    // TODO: Test TLS handshake with API server requirements

    // Hint: Use crypto/x509 for certificate parsing and validation

    return CertificateHealth{}

}

// TestAdmissionRequest sends a synthetic admission request to webhook

func (wd *WebhookDiagnostics) TestAdmissionRequest(ctx context.Context, db *databasev1.Database, operation string) AdmissionTestResult {

    // TODO: Construct admission review request with test Database resource

    // TODO: Set appropriate operation (CREATE, UPDATE, DELETE)

    // TODO: Send HTTPS POST request to webhook endpoint

    // TODO: Parse admission response and extract decision

    // TODO: Validate response format and required fields

    // Hint: Use admission/v1.AdmissionReview for request/response format

    return AdmissionTestResult{}

}

// ValidateWebhookConfiguration checks webhook registration with API server

func (wd *WebhookDiagnostics) ValidateWebhookConfiguration(ctx context.Context) error {

    // TODO: Fetch ValidatingAdmissionWebhook and MutatingAdmissionWebhook configurations

    // TODO: Verify webhook service and namespace references are correct

    // TODO: Check that CA bundle matches webhook certificate

    // TODO: Validate admission rule selectors and resource matching
}

```

```
// TODO: Test webhook failure policy configuration

// Hint: Use admissionregistration/v1 API to fetch webhook configs

return nil

}

// SimulateWebhookFailures tests webhook behavior under various failure conditions

func (wd *WebhookDiagnostics) SimulateWebhookFailures(ctx context.Context) error {

// TODO: Test webhook behavior when certificate is expired

// TODO: Test webhook timeout scenarios with slow responses

// TODO: Test malformed admission request handling

// TODO: Test webhook unavailability with failure policy

// TODO: Test concurrent admission request handling

// Hint: Use test HTTP servers that simulate failure modes

return nil

}
```

RBAC Validation Tool:

```
// internal/debug/rbac_validator.go
```

GO

```
package debug
```

```
import (
```

```
    "context"
```

```
    "fmt"
```

```
        authv1 "k8s.io/api/authorization/v1"
```

```
        rbacv1 "k8s.io/api/rbac/v1"
```

```
        metav1 "k8s.io/apimachinery/pkg/apis/meta/v1"
```

```
        "k8s.io/client-go/kubernetes"
```

```
)
```

```
type RBACValidator struct {
```

```
    Client      kubernetes.Interface
```

```
    ServiceAccount string
```

```
    Namespace    string
```

```
}
```

```
type PermissionTest struct {
```

```
    Resource    string
```

```
    Verb        string
```

```
    APIGroup   string
```

```
    Namespace   string
```

```
    Expected    bool
```

```
    Actual      bool
```

```
    Error       error
```

```
}
```

```
type RBACValidationResult struct {
```

```
    ServiceAccountExists bool
```

```
    RoleBindings     []rbacv1.RoleBinding
```

```
    ClusterRoleBindings []rbacv1.ClusterRoleBinding
```

```
    PermissionTests   []PermissionTest
```

```
    MissingPermissions []string
```

```
    ExcessivePermissions []string
```

```
}

// ValidateOperatorPermissions tests all required RBAC permissions

func (rv *RBACValidator) ValidateOperatorPermissions(ctx context.Context) (*RBACValidationResult, error) {

    // TODO: Verify service account exists in specified namespace

    // TODO: Find all role bindings that reference the service account

    // TODO: Find all cluster role bindings that reference the service account

    // TODO: Test each required permission using SubjectAccessReview

    // TODO: Report missing permissions and excessive grants

    // Hint: Use authorization/v1.SubjectAccessReview for permission testing

    return nil, nil
}

// TestCustomResourcePermissions validates CRD-specific RBAC rules

func (rv *RBACValidator) TestCustomResourcePermissions(ctx context.Context) error {

    // TODO: Test get, list, watch permissions on Database custom resources

    // TODO: Test create, update, patch permissions for Database resources

    // TODO: Test status subresource update permissions

    // TODO: Test permissions across different namespaces

    // TODO: Verify finalizer update permissions

    // Hint: Custom resources use different API groups than built-in resources

    return nil
}

// TestOwnedResourcePermissions validates permissions for resources owned by operator

func (rv *RBACValidator) TestOwnedResourcePermissions(ctx context.Context) error {

    // TODO: Test Deployment create, update, delete permissions

    // TODO: Test Service and ConfigMap management permissions

    // TODO: Test PersistentVolumeClaim permissions for storage

    // TODO: Test Secret permissions for database credentials

    // TODO: Verify permissions work across target namespaces

    // Hint: Test both namespace-scoped and cluster-scoped resources

    return nil
}
```

```

// GenerateMinimalRBACManifests creates RBAC with principle of least privilege

func (rv *RBACValidator) GenerateMinimalRBACManifests() ([]runtime.Object, error) {
    // TODO: Analyze actual API calls made by operator

    // TODO: Generate Role and ClusterRole with minimal required permissions

    // TODO: Create RoleBinding and ClusterRoleBinding for service account

    // TODO: Separate permissions by functional area (controller, webhooks, leader election)

    // TODO: Include comments explaining why each permission is needed

    // Hint: Group permissions logically and document rationale

    return nil, nil
}

```

Milestone Checkpoint: Controller Debugging After implementing controller debugging tools, verify functionality by:

1. Run `go run cmd/debug-controller/main.go` with a problematic Database resource
2. Expected output: detailed informer cache status, workqueue metrics, and reconciliation trace
3. Create a Database with intentional spec errors and confirm the debugging tool identifies the issue
4. Verify that infinite reconciliation loops are detected and reported with specific root causes

Milestone Checkpoint: Webhook Debugging

After implementing webhook diagnostics, validate by:

1. Run webhook certificate validation against expired or misconfigured certificates
2. Expected behavior: clear error messages identifying specific certificate problems
3. Test admission request simulation with various Database resource configurations
4. Confirm that webhook timeout and failure scenarios are properly detected and reported

Milestone Checkpoint: RBAC Validation After implementing RBAC validation tools, confirm by:

1. Run permission validation against service accounts with missing or excessive permissions
2. Expected output: detailed report of required vs actual permissions with specific gaps identified
3. Test custom resource permission validation across multiple namespaces
4. Verify that minimal RBAC manifest generation produces working configurations

Future Extensions and Evolution

Milestone(s): All milestones - this section guides how operators can evolve beyond their initial implementation while maintaining backward compatibility and operational stability

The beauty of the Kubernetes operator pattern lies not just in solving immediate automation problems, but in creating a foundation that can evolve with changing requirements. Unlike traditional configuration management tools that often require complete rewrites for new features, well-designed operators can grow organically through extension points and architectural patterns that preserve existing functionality while adding new capabilities.

Mental Model: The Living System

Think of a mature operator like a city's infrastructure. When a city grows, you don't tear down all the roads and buildings to add new neighborhoods. Instead, you extend existing transit systems, add new utility lines that connect to existing grids, and build new districts that

integrate with established ones. Similarly, a well-architected operator provides extension points that allow new functionality to plug into existing reconciliation loops, admission controls, and status reporting without disrupting operational workloads.

The key insight is that operators encode operational knowledge as code, and this knowledge base should be extensible. Each new feature should build upon the foundation of custom resources, controllers, and webhooks established in the initial implementation, rather than requiring parallel systems that duplicate effort and create consistency challenges.

Feature Extension Points

The operator pattern provides several natural extension points where new functionality can be added without breaking existing behavior. These extension points align with the separation of concerns established in the core architecture: custom resource schemas can be versioned and extended, controllers can watch additional resource types, and admission webhooks can be composed to handle multiple validation scenarios.

Mental Model: The Plugin Architecture

Consider how modern IDEs like Visual Studio Code handle extensions. The core editor provides fundamental capabilities (editing text, managing files), but extensions add language support, debugging tools, and specialized workflows. Each extension registers with well-defined interfaces, contributes to specific extension points, and can be installed or removed without affecting other extensions.

Similarly, operator extensions should leverage Kubernetes' native extension mechanisms. New custom resource definitions can be added alongside existing ones, additional controllers can be deployed to handle new resource types, and webhook configurations can be extended to cover new admission scenarios. The key is designing these extensions to be composable rather than monolithic.

Decision: Extension Architecture Strategy

- **Context:** Operators need to evolve with new requirements while maintaining backward compatibility and operational stability for existing resources
- **Options Considered:**
 1. Monolithic operator with all features in one controller
 2. Microservice operators with one controller per resource type
 3. Plugin-based architecture with core operator and extension points
- **Decision:** Plugin-based architecture with extension points for new resource types, custom reconciliation logic, and admission policies
- **Rationale:** This approach balances operational simplicity (single deployment) with extensibility (isolated feature development) while leveraging Kubernetes' native composition mechanisms
- **Consequences:** Extensions can be developed and tested independently, but require careful interface design to prevent tight coupling between core and extension code

Extension Point	Implementation Strategy	Backward Compatibility	Development Isolation
New Custom Resources	Additional CRDs with separate controllers	Complete - existing resources unaffected	High - independent schema and logic
Resource Composition	Cross-resource references and status aggregation	Maintained through optional fields	Medium - requires coordination protocols
Admission Policy Extensions	Webhook composition with policy plugins	Preserved through additive validation	High - independent policy modules
Status Reporting Extensions	Custom condition types and metrics	Compatible through condition namespacing	High - independent status providers
Reconciliation Strategies	Pluggable reconcilers with common interfaces	Maintained through interface versioning	Medium - shared interface contracts

Custom Resource Extensions

The most natural extension point involves adding new custom resource definitions that complement the existing `Database` resource. These extensions should follow the established patterns for schema design, status reporting, and controller integration while addressing new operational requirements.

Design Insight: New custom resources should be designed as first-class citizens rather than sub-resources or configuration options within existing types. This approach provides cleaner separation of concerns, independent versioning capabilities, and clearer RBAC boundaries.

Common custom resource extensions include:

Backup Resources that manage backup policies and execution schedules independently of database lifecycle management. A `DatabaseBackup` resource would reference `Database` instances through owner references or explicit references, allowing backup policies to be managed by different teams or systems while integrating with the database operator's status reporting.

Migration Resources that handle database schema changes and data migrations with their own reconciliation logic and status tracking. A `DatabaseMigration` resource would coordinate with `Database` controllers to ensure migrations execute during appropriate maintenance windows and update database status to reflect migration completion.

Monitoring Resources that configure observability and alerting for database instances based on operational requirements that may change independently of database configuration. A `DatabaseMonitor` resource would manage metrics collection, alert rules, and dashboard configuration while integrating with the database operator's health reporting.

Resource Extension	Primary Responsibility	Integration Points	Status Reporting
<code>DatabaseBackup</code>	Backup policy and execution	Owner references to <code>Database</code> , shared storage	Backup completion, schedule adherence
<code>DatabaseMigration</code>	Schema version management	Database readiness conditions, maintenance windows	Migration progress, rollback capability
<code>DatabaseMonitor</code>	Observability configuration	Database health status, performance metrics	Monitor health, alert firing status
<code>DatabaseCluster</code>	Multi-instance coordination	Database instance membership, leader election	Cluster topology, split-brain detection

Controller Composition Strategies

As new custom resources are added, the controller architecture must accommodate multiple reconciliation loops that may need to coordinate their actions. The key architectural decision involves balancing controller isolation with cross-resource coordination requirements.

Decision: Multi-Controller Architecture

- **Context:** Additional custom resources require their own reconciliation logic while maintaining coordination with existing database controllers
- **Options Considered:**
 1. Single controller handling all resource types with shared reconciliation logic
 2. Separate controllers per resource type with event-based coordination
 3. Controller hierarchy with primary controller coordinating secondary controllers
- **Decision:** Separate controllers per resource type with well-defined coordination protocols through status conditions and owner references
- **Rationale:** This approach maximizes development isolation, enables independent testing and deployment of new controllers, and leverages Kubernetes' native coordination mechanisms
- **Consequences:** Each controller can be developed and maintained independently, but coordination protocols must be carefully designed to prevent race conditions and inconsistent state

The controller composition strategy relies on several coordination mechanisms that build upon Kubernetes' native patterns:

Owner Reference Coordination establishes parent-child relationships between resources, ensuring that backup resources are automatically garbage-collected when their associated database is deleted. This pattern also enables status aggregation, where parent resources can report the health of their children through aggregated conditions.

Status Condition Protocols define standard condition types and status values that controllers use to communicate their state to other controllers. For example, a database controller might set a `MaintenanceWindow` condition that migration controllers watch before executing schema changes.

Event-Based Coordination leverages Kubernetes events and controller watches to trigger cross-controller actions. When a database controller detects a configuration change that requires backup policy updates, it can emit an event that backup controllers watch and respond to appropriately.

Coordination Mechanism	Implementation	Consistency Guarantees	Error Handling
Owner References	Automatic garbage collection	Eventually consistent cleanup	Finalizers for cleanup ordering
Status Conditions	Cross-controller status sharing	Optimistically consistent status	Condition-based retry logic
Event Watching	Cross-resource change detection	At-least-once event delivery	Event deduplication in handlers
Shared ConfigMaps	Configuration coordination	Read-your-writes consistency	Conflict resolution strategies

Admission Webhook Composition

As operators evolve to handle additional resource types and more sophisticated validation requirements, admission webhook architecture must support composable validation and mutation policies. The key challenge involves maintaining performance and reliability while adding new admission logic.

The webhook composition strategy builds upon the foundation established in Milestone 4 while providing extension points for new validation scenarios:

Policy Plugin Architecture allows new validation rules to be added without modifying core webhook code. Each policy plugin registers with the webhook server and receives admission requests for relevant resource types, enabling independent development and testing of specialized validation logic.

Validation Pipeline Composition enables multiple validation stages to be applied to the same resource, with each stage having the opportunity to add conditions, modify the admission response, or trigger additional validation steps based on resource content and cluster state.

Cross-Resource Validation supports validation rules that depend on the state of multiple resources, such as ensuring that backup schedules don't conflict with maintenance windows or verifying that migration resources reference valid database instances.

Webhook Extension	Validation Scope	Performance Impact	Failure Isolation
Policy Plugins	Resource-specific business rules	Low - parallel validation	High - plugin failures don't affect others
Pipeline Composition	Multi-stage validation workflows	Medium - sequential processing	Medium - early stage failures abort pipeline
Cross-Resource Validation	Cluster-wide consistency rules	High - requires additional API calls	Low - API failures affect all validation
Dynamic Policy Loading	Runtime policy updates	Low - cached policy evaluation	High - invalid policies isolated to specific requests

Scaling and Performance

As operators mature and manage larger numbers of resources across more clusters, performance characteristics become critical operational concerns. The scaling strategy must address both horizontal scaling (handling more resources) and vertical scaling (handling more complex resources) while maintaining the reliability and consistency guarantees established in the initial implementation.

Mental Model: The Air Traffic Control System

Consider how air traffic control systems handle increasing flight volumes. Rather than building larger and larger central control towers, the system is partitioned geographically with regional controllers handling local airspace while coordinating with adjacent regions for aircraft transitions. Each controller has local authority within their domain but follows standardized protocols for inter-region coordination.

Similarly, operator scaling strategies should partition the problem space rather than simply adding more computational resources to a monolithic controller. This partitioning can occur along resource boundaries (different controllers for different resource types), geographic boundaries (controllers scoped to specific regions or availability zones), or tenant boundaries (controllers scoped to specific namespaces or organizational units).

Horizontal Scaling Strategies

The foundation for horizontal scaling lies in the informer and workqueue architecture established in Milestone 2. However, as the number of managed resources grows, several bottlenecks emerge that require architectural evolution.

Decision: Controller Sharding Strategy

- **Context:** Single controller instances become CPU and memory bottlenecks when managing thousands of resources, and Kubernetes API server rate limiting affects reconciliation latency
- **Options Considered:**
 1. Vertical scaling with larger controller instances and more aggressive caching
 2. Resource-based sharding with multiple controller instances handling disjoint resource sets
 3. Namespace-based sharding with controllers scoped to specific namespace sets
- **Decision:** Hybrid sharding approach supporting both resource-based and namespace-based partitioning depending on deployment requirements
- **Rationale:** This approach provides flexibility for different scaling scenarios while maintaining compatibility with existing single-controller deployments
- **Consequences:** Sharding logic adds complexity to deployment and coordination, but enables near-linear scaling with resource counts

Resource-Based Sharding partitions resources across multiple controller instances using consistent hashing or explicit assignment strategies. Each controller instance handles a subset of all resources, reducing per-instance memory usage and API server load while maintaining independent reconciliation loops.

The sharding implementation extends the controller manager pattern with shard-aware resource filtering:

Sharding Strategy	Resource Distribution	Coordination Requirements	Rebalancing Complexity
Hash-Based Sharding	Consistent hash of resource name/namespace	None - independent shards	Medium - requires hash function changes
Range-Based Sharding	Alphabetical or numerical resource ranges	Shard boundary coordination	High - requires resource migration
Label-Based Sharding	Resources tagged with shard identifiers	Shard assignment coordination	Low - label changes trigger rebalancing
Namespace Sharding	Controller scope limited to namespace sets	Cross-namespace resource coordination	Medium - namespace assignment changes

Leader Election Evolution becomes more sophisticated in sharded deployments, where each shard requires independent leader election to ensure exactly one active controller per shard while supporting shard-level failover and rebalancing.

The leader election strategy builds upon the foundation from Milestone 5 but adds shard awareness:

```

Shard Assignment:
- Controller instances register with shard coordinator
- Coordinator assigns responsibility for resource ranges
- Each shard maintains independent leader election
- Failed shards trigger automatic rebalancing

Failover Behavior:
- Shard leader failures trigger immediate re-election
- Shard coordinator failures trigger shard assignment redistribution
- Controller instance failures trigger shard responsibility transfer
  
```

Performance Optimization Strategies

Beyond horizontal scaling, several performance optimization strategies address the computational and memory efficiency of individual controller instances. These optimizations build upon the informer and reconciliation patterns while reducing unnecessary work and improving resource utilization.

Informer Cache Optimization reduces memory usage and improves cache efficiency through selective caching strategies that store only the resource fields needed for reconciliation decisions rather than complete resource objects.

Reconciliation Batching groups related reconciliation operations to reduce API server round trips and improve throughput for scenarios where multiple resources can be processed together efficiently.

Incremental Reconciliation tracks resource changes at a finer granularity to avoid unnecessary reconciliation work when only status fields or irrelevant metadata changes occur.

Optimization Strategy	Memory Impact	Latency Impact	Complexity Cost
Selective Field Caching	50-80% reduction	Minimal increase	Medium - requires field selection logic
Reconciliation Batching	Minimal change	30-60% reduction	High - requires dependency analysis
Incremental Reconciliation	20-40% reduction	40-70% reduction	Medium - requires change detection logic
Status-Only Updates	Minimal change	20-30% reduction	Low - separate status client usage

Multi-Cluster Operations

As organizations adopt multi-cluster Kubernetes deployments, operators must evolve to handle resources that span cluster boundaries while maintaining consistency and coordination across distributed infrastructure.

The multi-cluster architecture builds upon single-cluster patterns while adding cluster-aware coordination:

Cluster-Scope Controllers manage resources within individual clusters while reporting status to cluster-wide coordination systems.

Each cluster maintains its own controller instances that handle local reconciliation while participating in cross-cluster coordination protocols.

Federation Controllers coordinate resources across multiple clusters, making placement decisions, managing resource distribution, and aggregating status information from cluster-scope controllers.

Cross-Cluster Networking handles the connectivity and service discovery requirements for resources that span cluster boundaries, such as database clusters with instances in multiple availability zones or regions.

Multi-Cluster Pattern	Coordination Mechanism	Consistency Model	Failure Handling
Hub-and-Spoke	Central controller with cluster agents	Eventually consistent	Hub failure affects coordination only
Peer-to-Peer	Direct cluster-to-cluster communication	Conflict-free replicated data types	Network partitions require conflict resolution
Event-Driven	Message queue coordination	At-least-once processing	Message delivery guarantees vary by queue
GitOps-Based	Configuration repository coordination	Strong consistency for desired state	Repository availability affects updates

⚠ Pitfall: Distributed State Consistency

A common mistake in multi-cluster operators involves assuming that status updates across clusters will be immediately consistent. This leads to race conditions where controllers make decisions based on stale information from other clusters, resulting in resource conflicts or inconsistent resource placement.

The solution involves designing coordination protocols that explicitly handle eventual consistency through techniques like vector clocks for causality tracking, conflict-free replicated data types for status aggregation, and idempotent operations that produce correct results regardless of message ordering or duplication.

Implementation Guidance

The extension and scaling strategies described above require careful implementation planning to ensure that new capabilities integrate cleanly with existing operator patterns while maintaining operational stability.

Technology Recommendations

Scaling Challenge	Simple Option	Advanced Option
Controller Sharding	Label-based shard assignment with manual configuration	Consistent hashing with automatic rebalancing using controller-runtime's cluster caching
Multi-Cluster Coordination	Shared ConfigMaps in hub cluster	Admiral service mesh integration with cross-cluster service discovery
Performance Monitoring	Prometheus metrics with basic controller dashboards	Custom resource status aggregation with SLI/SLO tracking
Extension Development	In-tree plugins with interface definitions	Out-of-tree controllers with standardized coordination protocols

Recommended Project Structure

Extensions should follow a modular structure that isolates new functionality while integrating with the core operator patterns:

```
project-root/
  cmd/
    database-operator/main.go      ← core operator entry point
    backup-controller/main.go     ← backup extension controller
    migration-controller/main.go  ← migration extension controller
  apis/
    database/v1/                 ← core database CRD
    backup/v1/                   ← backup extension CRD
    migration/v1/                ← migration extension CRD
  internal/
    controller/
      database/                  ← core reconciliation logic
        database_controller.go
        status_manager.go
      backup/                    ← backup extension logic
        backup_controller.go
        schedule_manager.go
      migration/                 ← migration extension logic
        migration_controller.go
        version_manager.go
    webhook/
      database/                  ← core admission logic
      common/                     ← shared webhook utilities
    coordination/                ← cross-controller coordination
      owner_reference.go
      status_aggregation.go
      event_coordination.go
  config/
    core/                        ← base operator manifests
    extensions/                  ← extension-specific manifests
    samples/                     ← example configurations
```

Extension Development Framework

The extension framework provides standardized interfaces and utilities that new controllers can implement to integrate seamlessly with the core operator:

GO

```
// ExtensionController defines the interface that all extension controllers must implement

// to participate in the operator's coordination protocols and lifecycle management.

type ExtensionController interface {

    // SetupWithManager registers the extension controller with the controller manager

    // and configures watches for relevant resource types.

    SetupWithManager(mgr ctrl.Manager) error

    // GetCoordinationSpecs returns the coordination requirements for this controller,
    // including which resources it owns, which status conditions it reports, and
    // which cross-controller events it needs to receive.

    GetCoordinationSpecs() CoordinationSpecs

    // HandleCoordinationEvent processes events from other controllers that may
    // trigger reconciliation or status updates in this controller.

    HandleCoordinationEvent(ctx context.Context, event CoordinationEvent) error

    // GetHealthStatus returns the current health and performance metrics for
    // this controller instance, used for monitoring and debugging.

    GetHealthStatus() HealthStatus

}

// CoordinationSpecs defines how an extension controller participates in
// cross-controller coordination and status aggregation.

type CoordinationSpecs struct {

    // OwnedResources lists the resource types this controller creates and manages

    OwnedResources []schema.GroupVersionKind

    // StatusConditions lists the condition types this controller reports

    StatusConditions []string

    // WatchedConditions lists condition types from other controllers that
    // trigger reconciliation in this controller

    WatchedConditions []string
}
```

```
// EmittedEvents lists event types this controller emits for coordination  
  
EmittedEvents []string  
  
}
```

Sharding Implementation Framework

The sharding framework extends the controller-runtime manager with shard-aware resource filtering and leader election:

```
// ShardingManager extends the controller-runtime manager with shard-aware  
// resource filtering and coordination capabilities.  
  
type ShardingManager struct {  
  
    // Manager is the underlying controller-runtime manager  
  
    ctrl.Manager  
  
  
    // ShardID identifies this controller instance's shard assignment  
  
    ShardID string  
  
  
    // ShardingStrategy determines how resources are assigned to shards  
  
    ShardingStrategy ShardingStrategy  
  
  
    // CoordinationClient handles shard assignment and rebalancing  
  
    CoordinationClient ShardCoordinator  
  
}  
  
  
// ShardingStrategy defines how resources are partitioned across controller instances  
  
type ShardingStrategy interface {  
  
    // AssignShard determines which shard should handle a given resource  
  
    AssignShard(resource client.Object) string  
  
  
    // IsAssignedToShard returns true if the resource should be handled by the given shard  
  
    IsAssignedToShard(resource client.Object, shardID string) bool  
  
  
    // RebalanceShards redistributes resources across available shards  
  
    RebalanceShards(availableShards []string) (ShardAssignments, error)  
  
}
```

Performance Monitoring Integration

The performance monitoring framework provides standardized metrics collection and alerting for operator scaling scenarios:

```
// PerformanceMetrics provides standardized metrics collection for operator scaling          GO

type PerformanceMetrics struct {

    // ReconciliationDuration tracks time spent in reconciliation loops

    ReconciliationDuration prometheus.HistogramVec

    // ResourceCacheSize tracks informer cache memory usage

    ResourceCacheSize prometheus.GaugeVec

    // WorkqueueDepth tracks pending reconciliation requests

    WorkqueueDepth prometheus.GaugeVec

    // APIServerLatency tracks Kubernetes API call performance

    APIServerLatency prometheus.HistogramVec

    // ShardDistribution tracks resource distribution across shards

    ShardDistribution prometheus.GaugeVec

}

// RecordReconciliation records metrics for a completed reconciliation operation

func (m *PerformanceMetrics) RecordReconciliation(
    controllerName string,
    resourceType string,
    duration time.Duration,
    result string,
) {
    // TODO 1: Record reconciliation duration in histogram with controller/type/result labels
    // TODO 2: Update workqueue depth gauge to reflect current queue size
    // TODO 3: Record any API server calls made during reconciliation
    // Hint: Use prometheus labels to enable per-controller and per-resource-type analysis
}
```

Milestone Checkpoints

Extension Development Checkpoint: After implementing a new extension controller, verify the integration by:

1. Deploying the extension alongside the core operator: `kubectl apply -f config/extensions/backup-controller/`
2. Creating a test resource that should trigger coordination: `kubectl apply -f config/samples/backup_v1_databasebackup.yaml`
3. Verifying cross-controller coordination through status conditions: `kubectl get database test-db -o jsonpath='{.status.conditions[?(@.type=="BackupReady")]}'`
4. Checking coordination events in controller logs: `kubectl logs -l app=database-operator | grep "coordination.event"`

Scaling Checkpoint: After implementing controller sharding, verify the scaling behavior by:

1. Deploy multiple controller replicas with sharding enabled: `kubectl scale deployment database-operator --replicas=3`
2. Create multiple test resources across different namespaces: `kubectl apply -f config/samples/scaling-test/`
3. Verify shard assignment distribution: `kubectl logs -l app=database-operator | grep "shard.assignment"`
4. Monitor performance metrics during scale testing: `curl http://localhost:8080/metrics | grep reconciliation_duration`

Performance Checkpoint: After implementing performance optimizations, validate the improvements by:

1. Establish baseline metrics with unoptimized operator: Deploy operator, create 100 test resources, measure reconciliation latency
2. Deploy optimized operator version with same test resources
3. Compare reconciliation performance: Expected 30-50% reduction in reconciliation duration
4. Verify memory usage improvements: Expected 40-60% reduction in controller memory usage
5. Validate correctness: All resources should reach Ready condition within expected timeframes

Glossary

Milestone(s): All milestones - provides definitions of Kubernetes, operator, and distributed systems terminology used throughout the entire operator implementation process

This glossary defines the technical terminology used throughout this design document. Understanding these concepts is essential for successfully implementing a Kubernetes operator. The terms are organized into categories to help build conceptual understanding of how different pieces of the operator pattern connect together.

Mental Model: The Specialized Dictionary

Think of this glossary like a specialized technical dictionary for a specific engineering discipline. Just as aerospace engineers have precise definitions for "thrust," "drag," and "lift" that might differ slightly from common usage, Kubernetes operators use terms like "reconciliation," "admission," and "informer" with very specific technical meanings. This glossary provides those precise definitions to prevent confusion when reading documentation or discussing implementation details with other engineers.

The terminology forms layers of abstraction, from low-level Kubernetes primitives up to operator-specific patterns. Understanding these layers helps you navigate between different levels of detail when designing and troubleshooting your operator.

Kubernetes Core Concepts

These fundamental Kubernetes terms form the foundation that operators build upon. Every operator developer must understand these concepts deeply.

Term	Definition	Why It Matters
API Server	Central component of Kubernetes control plane that exposes REST API for all cluster operations and stores cluster state in etcd	Operators interact with this component for all resource operations - understanding its role helps debug API errors and performance issues
Custom Resource Definition (CRD)	API extension defining new resource types with OpenAPI v3 schema validation that extends Kubernetes API without modifying core code	CRDs are how operators introduce new APIs - poorly designed CRDs create maintenance headaches and user experience problems
Custom Resource (CR)	Instance of a custom resource type defined by a CRD, stored in etcd and accessible through standard Kubernetes API operations	These are the resources your operator manages - understanding their lifecycle helps design better reconciliation logic
Resource Version	Monotonically increasing identifier for resource modifications used for optimistic locking and conflict detection	Critical for handling concurrent updates - ignoring resource versions leads to lost updates and race conditions
Owner Reference	Metadata field establishing parent-child relationships between resources for automatic garbage collection	Prevents resource leaks when parent resources are deleted - improper owner references cause cleanup failures
Finalizer	Metadata field preventing resource deletion until cleanup operations complete, allowing controllers to perform cleanup before removal	Essential for managing external resources - forgetting to remove finalizers causes resources to stuck in deletion
Namespace	Virtual cluster that provides scope for resource names and RBAC policies, enabling multi-tenancy within a single cluster	Affects controller permissions and resource isolation - namespace design impacts operator security model
Service Account	Identity for processes running in pods, used for authentication and authorization with the Kubernetes API	Operators run under service accounts - understanding this helps design proper RBAC policies
RBAC (Role-Based Access Control)	Authorization mechanism controlling what operations service accounts can perform on specific resources	Operators need carefully scoped permissions - overly broad RBAC creates security risks

Controller Pattern Concepts

The controller pattern is the heart of Kubernetes operators. These terms describe the continuous reconciliation process that makes operators work.

Term	Definition	Why It Matters
Controller	Software component implementing a control loop that continuously observes cluster state and takes corrective actions to maintain desired state	The core of every operator - understanding controller design patterns helps build robust, efficient operators
Reconciliation	Process of making actual state match desired state through continuous comparison and correction, implementing the controller's primary logic	This is what makes operators "smart" - poor reconciliation logic causes performance problems and incorrect behavior
Control Loop	Continuous cycle of observing current state, comparing to desired state, and taking corrective action to eliminate differences	Understanding control loops helps design operators that converge reliably without oscillation or instability
Desired State	Target configuration specified in resource spec fields that the controller works to achieve and maintain	Controllers compare actual state against this - unclear desired state specifications lead to unpredictable behavior
Actual State	Current state of managed resources as observed through API server queries and external system inspection	Controllers must accurately determine actual state - incorrect state detection leads to unnecessary reconciliation work
Event-Driven Architecture	Design pattern where controllers react to resource change events rather than continuously polling for changes	Enables efficient operators that respond quickly to changes - polling-based designs waste resources and respond slowly
Declarative Management	Approach where users specify desired end state rather than imperative steps, letting the controller determine how to achieve that state	Makes operators easier to use and more reliable - imperative approaches are harder to reason about and recover from failures

Informer and Caching Architecture

Informers provide the efficient event-driven foundation that controllers build on. Understanding these concepts is crucial for performance and correctness.

Term	Definition	Why It Matters
Informer	Client-side cache that watches API server events and maintains a local copy of resources, providing efficient reads without API server load	Enables scalable controllers - direct API server queries for every operation create performance bottlenecks
Shared Index Informer	Informer implementation that allows multiple controllers to share the same cache and watch streams, reducing memory usage and API server load	Multiple controllers can efficiently watch the same resource types - separate informers waste memory and network bandwidth
List-Watch Pattern	API pattern where clients first list all resources then watch for subsequent changes, providing consistent event streams	Foundation of Kubernetes event system - understanding this helps debug informer sync issues and cache inconsistencies
Resync Period	Interval at which informers re-list all resources to ensure cache consistency and recover from missed events	Balances cache freshness with API server load - too frequent causes load, too infrequent causes stale cache issues
Cache Sync	Process where informer cache becomes consistent with API server state before processing events, ensuring controllers don't act on incomplete data	Controllers must wait for cache sync - processing events before sync completion can cause incorrect reconciliation decisions
Index Function	Custom function that creates secondary indexes on cached resources for efficient lookups by fields other than name/namespace	Enables efficient queries like "find all pods with label X" - missing indexes force expensive linear searches
Event Handler	Function called when informer detects resource changes, typically enqueueing reconcile requests for later processing	Bridges informer events to controller work queues - blocking in event handlers causes event processing delays

Work Queue and Scheduling

Work queues provide reliable, ordered processing of reconciliation requests with proper error handling and rate limiting.

Term	Definition	Why It Matters
Work Queue	Thread-safe queue that buffers reconcile requests and provides features like rate limiting, exponential backoff, and duplicate detection	Prevents overwhelming controllers with requests - poor queue management causes performance problems and resource exhaustion
Rate Limiting	Mechanism that controls how frequently items can be processed from the work queue, preventing resource exhaustion during high load	Protects controllers and API server from overload - missing rate limiting can cause cascade failures
Exponential Backoff	Retry strategy that increases delay between attempts exponentially, preventing overwhelming failing systems with repeated requests	Essential for handling transient failures gracefully - linear backoff can worsen system overload conditions
Jitter	Random variation added to retry delays to prevent multiple controllers from retrying simultaneously and creating thundering herd effects	Distributes retry load over time - without jitter, synchronized retries can overwhelm recovering systems
Requeue	Operation that adds a reconcile request back to the work queue for processing after a delay, used for retry logic and periodic reconciliation	Enables robust error handling and state drift correction - improper requeue logic causes infinite loops or missed updates
Deduplication	Work queue feature that prevents duplicate reconcile requests for the same resource from being processed simultaneously	Improves efficiency and prevents conflicting reconciliation attempts - missing deduplication wastes controller resources

Reconciliation Logic Patterns

These patterns describe common approaches for implementing robust reconciliation logic that handles complex state management scenarios.

Term	Definition	Why It Matters
Idempotent Operation	Operation that produces the same result when called multiple times with the same input, essential for reliable reconciliation	Controllers may reconcile the same resource multiple times - non-idempotent operations cause inconsistent state
Three-Way Merge	Reconciliation pattern that compares desired state, current state, and last-applied state to determine minimal changes needed	Enables efficient updates that preserve user modifications while applying operator changes
Drift Detection	Process of identifying when actual state has diverged from desired state due to external changes or system issues	Allows operators to correct configuration drift - without drift detection, external changes persist indefinitely
State Machine	Model where resources transition through well-defined states with specific triggers and actions for each transition	Simplifies complex lifecycle management - ad-hoc state handling leads to bugs and unpredictable behavior
Condition Types	Standardized status fields that report specific aspects of resource state like <code>Ready</code> , <code>Progressing</code> , and <code>Degraded</code>	Provides consistent status reporting across different resource types - improper condition usage confuses users
Phase	High-level resource lifecycle state like <code>Pending</code> , <code>Running</code> , <code>Failed</code> that summarizes overall resource condition	Gives users quick understanding of resource health - inconsistent phase reporting makes troubleshooting difficult
Generation	Counter that increments when resource spec changes, used to track which changes have been processed	Helps controllers track which spec changes they've reconciled - ignoring generation can cause missed updates
Observed Generation	Generation number that controller has successfully processed, reported in resource status	Shows users that their spec changes have been applied - missing observed generation updates confuse users about operator responsiveness

Admission Control and Webhooks

Admission webhooks provide validation and mutation capabilities that run during resource creation and updates, implementing the "gatekeeper pattern" for resource operations.

Term	Definition	Why It Matters
Admission Controller	Component in Kubernetes API server request processing pipeline that can validate, mutate, or deny resource operations	Understanding admission flow helps design webhooks that integrate properly with other admission controllers
Admission Webhook	HTTP endpoint that receives admission review requests from the API server to validate or mutate resources before they are stored	Extends Kubernetes with custom validation and defaulting logic - webhook failures can block all resource operations
Validating Webhook	Admission webhook that can approve or deny resource operations but cannot modify the resources being processed	Enforces business rules and constraints - failing validation provides immediate feedback to users
Mutating Webhook	Admission webhook that can modify resources before they are stored in etcd, typically used for injecting default values	Provides consistent defaulting and transformation - mutation order matters when multiple webhooks modify the same resource
Admission Review	Kubernetes API object containing resource data and metadata that is sent to webhooks for processing	Understanding this structure is essential for implementing webhook handlers correctly
Admission Response	Object returned by webhooks indicating whether to allow the operation and any mutations to apply	Proper response formatting ensures API server processes webhook decisions correctly
JSON Patch	RFC 6902 standard for describing changes to JSON documents using operations like add, replace, and remove	Mutating webhooks use JSON patches to specify resource modifications - incorrect patches cause admission failures
Strategic Merge Patch	Kubernetes-specific patch format that understands resource schema and can merge lists intelligently	Used internally by kubectl and controllers for resource updates - understanding helps debug update conflicts

TLS and Certificate Management

Webhooks require HTTPS communication, making certificate management a critical operational concern for operators.

Term	Definition	Why It Matters
TLS Certificate	X.509 certificate used to establish encrypted HTTPS connections between API server and admission webhooks	Required for webhook security - certificate issues prevent webhook operations and can block cluster operations
Certificate Authority (CA)	Entity that signs certificates, establishing trust relationships for TLS connections	Webhook certificates must be trusted by the API server - CA configuration errors prevent webhook communication
Subject Alternative Name (SAN)	Certificate extension that specifies additional identities the certificate is valid for, including DNS names and IP addresses	Webhook certificates must have correct SAN entries for service names - missing SANs cause TLS verification failures
Certificate Rotation	Process of automatically renewing and replacing TLS certificates before they expire to maintain service availability	Prevents webhook outages from certificate expiration - manual certificate management is error-prone and operationally expensive
Cert-Manager	Kubernetes operator that automates certificate lifecycle management including provisioning, renewal, and distribution	Simplifies webhook certificate management - manual certificate handling creates operational burden
Self-Signed Certificate	Certificate signed by its own private key rather than a certificate authority, useful for development and testing	Acceptable for development but requires careful CA bundle management - production should use proper CA-signed certificates

Error Handling and Reliability Patterns

Robust operators must handle various failure modes gracefully while maintaining system stability and providing clear diagnostics.

Term	Definition	Why It Matters
Transient Error	Temporary failure that is likely to succeed if retried after a delay, such as network timeouts or resource conflicts	Should trigger exponential backoff retry - treating transient errors as permanent wastes reconciliation opportunities
Permanent Error	Failure that will not succeed if retried without external intervention, such as validation errors or missing permissions	Should not trigger automatic retry - permanent error retry wastes resources and creates noise in logs
Circuit Breaker	Pattern that temporarily stops calling a failing service to prevent cascading failures and allow recovery time	Protects external dependencies from overload - missing circuit breakers can cause widespread system failures
Graceful Degradation	Approach where systems continue operating with reduced functionality when components fail rather than failing completely	Maintains service availability during partial failures - brittle systems cause unnecessary outages
Error Classification	Process of categorizing errors by type, severity, and appropriate response to enable proper handling strategies	Enables appropriate retry and escalation behavior - poor error classification leads to incorrect recovery actions
Leader Election	Distributed coordination mechanism that ensures only one controller replica actively reconciles resources at any time	Prevents split-brain scenarios in high availability deployments - missing leader election causes resource conflicts
Lease Object	Kubernetes resource used for leader election coordination, containing holder identity and renewal timestamp	Provides distributed locking primitive - lease configuration affects failover time and split-brain risk

Testing and Validation Concepts

Comprehensive testing ensures operator reliability and helps catch regressions during development and maintenance.

Term	Definition	Why It Matters
Fake Client	Test implementation of Kubernetes client that simulates API server behavior without requiring a real cluster	Enables fast unit tests that verify controller logic in isolation - real cluster tests are slow and harder to control
Envtest	Testing framework that runs a real Kubernetes API server in-process for integration testing without full cluster setup	Provides realistic testing environment without operational complexity - catches integration issues that unit tests miss
Test Double	Generic term for fake objects used in testing, including mocks, stubs, and fakes	Isolates code under test from dependencies - improper test doubles can hide real integration problems
Contract Test	Test that verifies component interfaces and behavior contracts without testing implementation details	Ensures components integrate correctly - implementation-focused tests break when refactoring
Chaos Testing	Testing approach that introduces random failures to verify system resilience and recovery behavior	Validates error handling and recovery logic - systems often fail in unexpected ways during real outages
Property-Based Testing	Testing technique that verifies system properties hold across a wide range of generated inputs	Finds edge cases that example-based tests miss - especially valuable for testing reconciliation logic

Deployment and Operations

Production operators require careful attention to security, observability, and operational concerns beyond basic functionality.

Term	Definition	Why It Matters
Principle of Least Privilege	Security practice of granting the minimum permissions necessary for operation, reducing blast radius of security breaches	Minimizes security risk from compromised operator - overly broad permissions increase attack surface
ClusterRole vs Role	RBAC distinction between cluster-scoped permissions (ClusterRole) and namespace-scoped permissions (Role)	Affects security isolation and operator deployment model - wrong scope causes permission errors or security violations
Service Account Token	JWT token that provides authentication credentials for processes running in pods	Enables operator authentication with API server - token management affects security and reliability
Pod Security Standards	Kubernetes security policies that control pod security contexts and capabilities	Affects how operator pods can run - restrictive policies may prevent operator functionality
Resource Quota	Kubernetes mechanism that limits resource consumption within namespaces	Can prevent operator from creating resources - understanding quotas helps design operators that work within limits
Network Policy	Kubernetes resource that controls network traffic between pods using label selectors	May affect operator communication with managed resources - network restrictions can cause unexpected failures
Observability	Practice of understanding system behavior through metrics, logging, and tracing	Essential for operating production systems - poor observability makes troubleshooting difficult and time-consuming

Performance and Scaling Concepts

As operators manage more resources, performance and scaling considerations become increasingly important for maintaining system stability.

Term	Definition	Why It Matters
Horizontal Scaling	Adding more controller instances to handle increased load by distributing work across multiple processes	Increases processing capacity for large resource counts - requires careful coordination to prevent conflicts
Vertical Scaling	Increasing computational resources (CPU, memory) for individual controller instances to handle more work	Simpler than horizontal scaling but has physical limits - helps handle resource-intensive reconciliation logic
Sharding	Technique for partitioning resources across multiple controller instances based on criteria like namespace or resource hash	Enables horizontal scaling without resource conflicts - poor sharding strategies cause load imbalances
Back Pressure	Mechanism where overloaded components signal upstream components to slow down request rates	Prevents cascading failures during high load - missing back pressure can cause system-wide overload
Resource Pooling	Pattern where expensive resources like database connections are shared across multiple operations	Improves efficiency and reduces resource usage - poor pooling causes resource exhaustion
Batch Processing	Technique for grouping multiple operations together to reduce overhead and improve throughput	Reduces API server load and improves performance - batching must balance latency with efficiency

Advanced Operator Patterns

These patterns enable sophisticated operator behaviors for complex application management scenarios.

Term	Definition	Why It Matters
Operator Composition	Pattern where multiple operators coordinate to manage related resources without conflicts	Enables modular operator design - poor composition causes resource conflicts and inconsistent behavior
Cross-Controller Coordination	Communication and synchronization between different controllers managing related resources	Prevents conflicts when multiple controllers affect the same resources - missing coordination causes race conditions
Multi-Tenant Operator	Operator design that safely manages resources for multiple users or teams within the same cluster	Enables shared operator infrastructure - poor multi-tenancy creates security and isolation issues
Operator Lifecycle Management	Practices for managing operator installation, upgrades, and removal without disrupting managed applications	Ensures smooth operator evolution - poor lifecycle management causes service disruptions during updates
Extension Points	Architectural patterns that allow new functionality to be added to operators without modifying core code	Enables customization and ecosystem development - rigid operators are difficult to extend for specific use cases
Federation	Pattern for managing resources across multiple Kubernetes clusters from a single control point	Enables multi-cluster operations - federation adds significant complexity and failure modes

Common Anti-Patterns and Pitfalls

Understanding what not to do is often as important as knowing correct approaches. These anti-patterns cause common operator reliability and performance issues.

Anti-Pattern	Description	Why It's Problematic	How to Avoid
Polling-Based Reconciliation	Continuously checking resource state on fixed intervals rather than using event-driven updates	Wastes resources and responds slowly to changes	Use informers and event-driven architecture
Blocking Reconciliation	Performing long-running operations synchronously in reconciliation logic without timeouts	Causes controller hangs and prevents processing other resources	Use asynchronous operations with proper timeouts
Status Thrashing	Repeatedly updating resource status with the same values or minor variations	Creates unnecessary API server load and version conflicts	Compare current and desired status before updating
Infinite Reconciliation	Reconciliation logic that always finds differences and requeues indefinitely	Wastes resources and prevents controller from processing other work	Ensure reconciliation can reach stable state
Missing Error Classification	Treating all errors the same way regardless of whether they're transient or permanent	Causes inappropriate retry behavior and resource waste	Classify errors and apply appropriate handling strategies
Overly Broad RBAC	Granting cluster-admin or excessive permissions to operators	Creates security vulnerabilities and violates least privilege	Analyze required permissions and grant minimum necessary

Implementation Guidance

Understanding operator terminology is essential, but applying these concepts correctly requires practical knowledge of how they map to actual code and configuration. This section provides concrete guidance for implementing the concepts defined above.

Technology Recommendations for Operator Development

Component	Simple Option	Advanced Option
Framework	Raw client-go with manual informer setup	Kubebuilder or Operator SDK with scaffolding
Testing	Standard go test with fake client	Ginkgo + Gomega with envtest integration
Logging	Standard log package	Structured logging with logr interface
Metrics	Basic HTTP /metrics endpoint	Prometheus metrics with controller-runtime
Webhooks	Manual HTTP server with TLS	Kubebuilder webhook generation
RBAC	Manual role definitions	RBAC markers with automated generation

Essential Go Packages for Kubernetes Operators

The Kubernetes ecosystem provides several packages that implement the concepts defined in this glossary. Understanding which packages provide which functionality helps you build operators efficiently.

```

// Core Kubernetes client libraries

import (
    "k8s.io/client-go/kubernetes"           // Typed client for built-in resources
    "k8s.io/client-go/dynamic"             // Dynamic client for custom resources
    "k8s.io/client-go/tools/cache"         // Informer and caching infrastructure
    "k8s.io/client-go/util/workqueue"      // Rate-limited work queue implementation
    "k8s.io/client-go/tools/leaderelection" // Leader election coordination
)

// Controller runtime - higher-level abstractions

import (
    "sigs.k8s.io/controller-runtime/pkg/manager"    // Manager orchestrates controllers
    "sigs.k8s.io/controller-runtime/pkg/controller" // Controller registration and setup
    "sigs.k8s.io/controller-runtime/pkg/reconcile"   // Reconciler interface definition
    "sigs.k8s.io/controller-runtime/pkg/webhook"       // Admission webhook framework
    "sigs.k8s.io/controller-runtime/pkg/client"        // Unified client interface
)

// Testing infrastructure

import (
    "sigs.k8s.io/controller-runtime/pkg/envtest"     // Integration testing with real API server
    "k8s.io/client-go/kubernetes/fake"               // Fake client for unit testing
    "sigs.k8s.io/controller-runtime/pkg/client/fake" // Fake client with controller-runtime
)

```

Project Structure for Operator Development

Organizing operator code properly from the start prevents confusion as the project grows. This structure follows Kubernetes conventions and separates concerns appropriately.

```
database-operator/
├── cmd/
│   └── manager/
│       └── main.go                      # Operator entry point and manager setup
├── api/
│   └── v1/
│       ├── database_types.go            # Custom resource type definitions
│       ├── database_webhook.go          # Admission webhook implementations
│       ├── zz_generated.deepcopy.go     # Generated deep copy methods
│       └── groupversion_info.go         # API group and version metadata
├── controllers/
│   ├── database_controller.go        # Main reconciliation logic
│   ├── database_controller_test.go   # Controller unit tests
│   └── suite_test.go                # Test suite setup with envtest
├── webhooks/
│   ├── admission_handler.go          # HTTP server for admission webhooks
│   ├── certificate_manager.go        # TLS certificate management
│   └── webhook_test.go              # Webhook integration tests
└── pkg/
    ├── resources/                   # Resource creation and management utilities
    ├── conditions/                 # Condition management helpers
    └── errors/                     # Error classification and handling
├── config/
    ├── crd/                         # Custom Resource Definition manifests
    ├── rbac/                        # Role and ServiceAccount definitions
    ├── webhook/                     # Webhook configuration manifests
    ├── manager/                     # Deployment and service manifests
    └── samples/                     # Example custom resource instances
├── hack/
    ├── boilerplate.go.txt           # License header template
    └── tools.go                     # Build tool dependency management
└── Makefile                       # Build, test, and deployment automation
```

Core Type Definitions

Understanding how the terminology maps to actual Go types helps bridge the gap between concepts and implementation. These definitions show the essential structures every operator needs.

```
// Reconciler implements the core reconciliation logic
type DatabaseReconciler struct {
    client.Client          // Kubernetes API client
    Scheme *runtime.Scheme // Resource type registry
    Recorder record.EventRecorder // Event recording for observability
}

// Reconcile implements the reconciliation loop for Database resources
func (r *DatabaseReconciler) Reconcile(ctx context.Context, req ctrl.Request) (ctrl.Result, error) {
    // TODO: Implement reconciliation logic following control loop pattern

    // TODO: Fetch current Database resource and handle not found errors

    // TODO: Determine desired state from Database spec

    // TODO: Compare actual state with desired state

    // TODO: Create, update, or delete resources to eliminate differences

    // TODO: Update Database status with current conditions and state

    // TODO: Return appropriate Result for requeue or completion
}
```

GO

Debugging Command Reference

When implementing operators, these commands help diagnose common issues by examining the cluster state and operator behavior.

Issue Type	Diagnostic Command	What It Shows
CRD Issues	<code>kubectl get crd databases.example.com -o yaml</code>	CRD definition, status, and validation schema
Resource Status	<code>kubectl describe database my-database</code>	Resource events, conditions, and status details
Controller Logs	<code>kubectl logs -n operator-system controller-manager-xxx</code>	Reconciliation events and error messages
Webhook Issues	<code>kubectl get validatingwebhookconfigurations</code>	Webhook registration and configuration
RBAC Problems	<code>kubectl auth can-i create databases --as=system:serviceaccount:operator-system:controller</code>	Permission verification
Informer Cache	Look for "Failed to sync cache" in controller logs	Cache synchronization problems
Leader Election	<code>kubectl get leases -n operator-system</code>	Current leader and lease status

Common Implementation Patterns

These code patterns implement the concepts defined in the glossary using idiomatic Go and Kubernetes conventions.

```
// Error classification pattern for retry logic
```

```
func (r *DatabaseReconciler) HandleReconcileError(ctx context.Context, db *Database, err error) (ctrl.Result, error) {
```

```
    if IsTransientError(err) {
```

```
        // Transient errors should be retried with exponential backoff
```

```
        return ctrl.Result{RequeueAfter: CalculateBackoffDelay(err)}, nil
```

```
    }
```

```
    if IsPermanentError(err) {
```

```
        // Permanent errors should update status but not requeue
```

```
        r.updateCondition(db, "Ready", "False", "ConfigurationError", err.Error())
```

```
        return ctrl.Result{}, r.Status().Update(ctx, db)
```

```
    }
```

```
    // Unknown errors are treated as transient
```

```
    return ctrl.Result{RequeueAfter: DefaultRequeueDelay}, err
```

```
}
```

```
// Idempotent resource creation pattern
```

```
func (r *DatabaseReconciler) createOrUpdateResource(ctx context.Context, desired client.Object) error {
```

```
    existing := desired.DeepCopyObject().(client.Object)
```

```
    err := r.Get(ctx, client.ObjectKeyFromObject(desired), existing)
```

```
    if apierrors.NotFound(err) {
```

```
        // Resource doesn't exist - create it
```

```
        return r.Create(ctx, desired)
```

```
    } else if err != nil {
```

```
        return err
```

```
    }
```

```
    // Resource exists - update if different
```

```
    if !reflect.DeepEqual(existing, desired) {
```

```
        desired.SetResourceVersion(existing.GetResourceVersion())
```

```
        return r.Update(ctx, desired)
```

```
    }
```

GO

```
    return nil // No changes needed  
}  
}
```

Milestone Validation Checkpoints

After implementing operator components, these checkpoints verify that the terminology concepts are working correctly in your implementation.

Controller Setup Validation:

```
# Verify informer cache synchronization  
  
kubectl logs -n operator-system controller-manager-xxx | grep "Starting workers"  
  
# Check work queue processing  
  
kubectl logs -n operator-system controller-manager-xxx | grep "Successfully Reconciled"  
  
# Validate leader election  
  
kubectl get leases -n operator-system
```

BASH

Admission Webhook Validation:

```
# Test validating webhook

kubectl apply -f - <<EOF

apiVersion: example.com/v1

kind: Database

metadata:

  name: test-validation

spec:

  replicas: -1 # Invalid value should be rejected

EOF

# Test mutating webhook default injection

kubectl apply -f - <<EOF

apiVersion: example.com/v1

kind: Database

metadata:

  name: test-defaulting

spec:

  replicas: 3

  # version field omitted - should be defaulted

EOF

kubectl get database test-defaulting -o yaml | grep version
```

BASH

Status and Condition Validation:

```
# Check condition updates

kubectl get database my-database -o jsonpath='{.status.conditions[*].type}'

# Verify observed generation tracking

kubectl get database my-database -o jsonpath='{.metadata.generation} {.status.observedGeneration}'
```

BASH

These validation steps ensure your operator correctly implements the patterns and concepts defined in this glossary, providing confidence that the terminology translates to working code.