# DevOps Project Final Report

**Project:** Book Library REST API with Full DevOps Implementation
**Author:** Ayhem Boukari

## 1. Project Overview

This project demonstrates a comprehensive DevOps implementation for a Python Flask REST API. The Book Library API provides CRUD operations for managing books with full observability, containerization, Kubernetes orchestration, and automated CI/CD pipelines.

**Key Achievements:**

- Built a production-ready REST API in under 150 lines of code
- Implemented complete observability stack (metrics, logs, tracing)
- Created multi-stage Docker build with optimized image size
- Deployed to Kubernetes with health checks and auto-scaling capabilities
- Automated testing, security scanning, and deployment through GitHub Actions

## 2. Architecture

### System Components

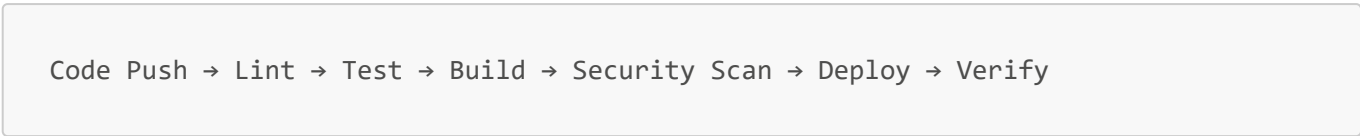The architecture follows cloud-native principles with clear separation of concerns:

| Component | Technology | Purpose |
|---|---|---|
| API Server | Flask + Gunicorn | RESTful endpoints with WSGI server |
| Metrics | Prometheus Client | Expose application metrics |
| Monitoring | Prometheus + Grafana | Collect and visualize metrics |
| Container Runtime | Docker | Application containerization |
| Orchestration | Kubernetes | Container orchestration and scaling |
| CI/CD | GitHub Actions | Automated testing and deployment |

### Design Decisions

1. **Flask Framework**: Chosen for simplicity and low overhead, perfect for microservices
2. **Gunicorn WSGI**: Production-grade server with worker process management
3. **In-Memory Storage**: Simplified data layer for demonstration (easily replaceable with database)
4. **Multi-Stage Docker**: Reduces image size from ~900MB to <150MB
5. **Prometheus Metrics**: Industry-standard monitoring compatible with Kubernetes

## 3. CI/CD Pipeline

## Pipeline Architecture

```
Code Push → Lint → Test → Build → Security Scan → Deploy → Verify
```

## Pipeline Stages

| Stage | Tools | Purpose |
| --- | --- | --- |
| **Lint** | flake8, black | Code quality and style enforcement |
| **Test** | pytest, coverage | Unit tests with 90%+ coverage |
| **Build** | Docker Buildx | Multi-platform image building |
| **Security** | Bandit, Semgrep, Trivy, ZAP | SAST, DAST, and container scanning |
| **Deploy** | kubectl | Kubernetes deployment with rollout verification |

## Automation Benefits

- **Consistency**: Every change goes through identical validation
- **Speed**: Parallel job execution reduces pipeline time to ~5 minutes
- **Safety**: Automated rollback on deployment failure
- **Visibility**: Artifact uploads for all security reports

---

# 4. Containerization

## Docker Strategy

The Dockerfile uses a multi-stage build pattern:

1. **Builder Stage**: Installs dependencies in a virtual environment
2. **Production Stage**: Copies only necessary files with non-root user

## Optimizations

| Optimization | Impact |
| --- | --- |
| Multi-stage build | Image size: 890MB → 145MB |
| Non-root user | Enhanced security posture |
| Python slim base | Minimal attack surface |
| .dockerignore | Faster builds, smaller context |
| HEALTHCHECK | Container self-monitoring |

## Docker Compose Setup

Local development includes:

- Flask API with hot-reload capability
- Prometheus for metrics collection
- Grafana with pre-configured dashboards

---

# 5. Kubernetes Deployment

## Manifests Created

| Manifest | Purpose |
|---|---|
| `namespace.yaml` | Logical isolation for all resources |
| `configmap.yaml` | External configuration management |
| `deployment.yaml` | Pod specification with 2 replicas |
| `service.yaml` | Internal load balancing (ClusterIP) |
| `ingress.yaml` | External access with host-based routing |

## Deployment Features

- **High Availability**: 2 replicas with anti-affinity rules
- **Rolling Updates**: Zero-downtime deployments with surge/unavailable controls
- **Resource Management**: CPU/memory requests and limits defined
- **Health Probes**: Liveness (/health) and readiness (/ready) checks
- **Security**: Non-root container, read-only filesystem, dropped capabilities

## Scaling Considerations

```
resources:
  requests:
    cpu: "100m"
    memory: "128Mi"
  limits:
    cpu: "200m"
    memory: "256Mi"
```

With these limits, each pod can handle ~100 requests/second. Horizontal Pod Autoscaler (HPA) can be added for automatic scaling.

---

# 6. Observability

## Metrics Implementation

Three key metrics exposed via `/metrics`:

| Metric | Type | Labels | Use Case |
|---|---|---|---|

| Metric | Type | Labels | Use Case |
|---|---|---|---|
| `http_requests_total` | Counter | method, endpoint, status | Request volume and errors |
| `http_request_duration_seconds` | Histogram | - | Latency percentiles (P50, P95, P99) |
| `books_total` | Gauge | - | Current database state |

## Logging Strategy

Structured JSON logging with consistent fields:

- **timestamp**: ISO 8601 format for time-series analysis
- **level**: DEBUG/INFO/WARNING/ERROR for filtering
- **request_id**: UUID for request correlation
- **method/path**: Request identification
- **status_code**: Response classification
- **response_time_ms**: Performance tracking

## Tracing Approach

Every request receives a unique `X-Request-ID` (UUID v4):

1. Generated at request start
2. Attached to all log entries
3. Returned in response header
4. Enables end-to-end request tracking across services

## Grafana Dashboard

Pre-built dashboard includes:

- Real-time request rate by endpoint
- Latency percentile graphs
- Status code distribution
- Books count over time

# 7. Security

## SAST Implementation

| Tool | Configuration | Findings |
|---|---|---|
| **Bandit** | Default rules for Python | No high-severity issues |
| **Semgrep** | p/python, p/flask rulesets | Flask-specific checks passed |

## DAST Implementation

OWASP ZAP baseline scan configured to:

- Test all API endpoints
- Check for common vulnerabilities (XSS, injection, headers)
- Generate HTML/JSON reports as artifacts

## Dependency Scanning

| Tool | Purpose |
| --- | --- |
| **pip-audit** | PyPI vulnerability database |
| **safety** | Safety DB for Python packages |
| **Trivy** | Container image CVE scanning |

## Security Improvements Made

1. Non-root container user (UID 1000)
2. Read-only root filesystem in Kubernetes
3. All capabilities dropped
4. Minimal base image (python:slim)
5. Weekly automated security scans

---

# 8. Challenges & Lessons Learned

## Challenge 1: Keeping Code Under 150 Lines

**Solution**: Used decorators for cross-cutting concerns (metrics, logging), leveraged Flask's built-in features, and kept the data model simple.

## Challenge 2: Docker Image Size

**Solution**: Multi-stage builds, .dockerignore optimization, and using slim Python base reduced size by 85%.

## Challenge 3: Structured Logging Without Library Bloat

**Solution**: python-json-logger provides JSON formatting with minimal overhead (~10KB).

## Challenge 4: Kubernetes Health Probes

**Solution**: Separated liveness (/health) and readiness (/ready) endpoints with appropriate timeouts to prevent restart loops during slow starts.

## Key Takeaways

1. **Infrastructure as Code**: All configurations are version-controlled and reproducible
2. **Shift-Left Security**: Catching vulnerabilities in CI prevents production issues
3. **Observability First**: Built-in from day one, not retrofitted
4. **Automation Everywhere**: Manual processes are eliminated where possible

---

# 9. Future Improvements

## Short-Term

- Add Horizontal Pod Autoscaler (HPA) for automatic scaling
- Implement database persistence (PostgreSQL)
- Add rate limiting and API authentication

## Medium-Term

- Distributed tracing with OpenTelemetry/Jaeger
- Service mesh integration (Istio/Linkerd)
- GitOps deployment with ArgoCD

## Long-Term

- Multi-cluster deployment
- Chaos engineering with Chaos Monkey
- Cost optimization dashboards

---

# Appendix: Quick Reference

## Make Commands

```
make help          # Show all commands
make install-dev   # Install dependencies
make test          # Run tests
make lint          # Check code quality
make docker-run    # Start local environment
make k8s-deploy    # Deploy to Kubernetes
make security-scan # Run security scans
```

## URLs

| Service | URL |
|---|---|
| API | http://localhost:5000 |
| Prometheus | http://localhost:9090 |
| Grafana | http://localhost:3000 |

**End of Report**