# Technical Challenge - Code Review and Deployment Pipeline Orchestration

**Format:** Structured interview with whiteboarding/documentation
**Assessment Focus:** Problem decomposition, AI prompting strategy, system design

**Please Fill in your Responses in the Response markdown boxes**

---

## Challenge Scenario

You are tasked with creating an AI-powered system that can handle the complete lifecycle of code review and deployment pipeline management for a mid-size software company. The system needs to:

**Current Pain Points:**

- Manual code reviews take 2-3 days per PR
- Inconsistent review quality across teams
- Deployment failures due to missed edge cases
- Security vulnerabilities slip through reviews
- No standardized deployment process across projects
- Rollback decisions are manual and slow

**Business Requirements:**

- Reduce review time to <4 hours for standard PRs
- Maintain or improve code quality
- Catch 90%+ of security vulnerabilities before deployment
- Standardize deployment across 50+ microservices
- Enable automatic rollback based on metrics
- Support multiple environments (dev, staging, prod)
- Handle both new features and hotfixes

---

## Part A: Problem Decomposition (25 points)

**Question 1.1:** Break this challenge down into discrete, manageable steps that could be handled by AI agents or automated systems. Each step should have:

- Clear input requirements
- Specific output format

- Success criteria
- Failure handling strategy

**Question 1.2:** Which steps can run in parallel? Which are blocking? Where are the critical decision points?

**Question 1.3:** Identify the key handoff points between steps. What data/context needs to be passed between each phase?

# Response Part A:

**Question 1.1: Discrete Steps**

1. **PR Ingestion & Context Extraction**

   - **Input:** Webhook event (PR created/updated), Repository URL, Diff content.
   - **Output:** Structured JSON containing file changes, commit messages, and author context.
   - **Success:** Successfully fetched all changed files and metadata.
   - **Failure:** Retry on API rate limits; alert if repository is inaccessible.

2. **Automated Static Analysis & Security Scan**

   - **Input:** Source code files from the PR.
   - **Output:** Report of linting errors, vulnerability findings (CVEs), and code coverage metrics.
   - **Success:** Scanners complete within defined timeout; report generated.
   - **Failure:** Fail the pipeline if critical severity issues are found; log scanner errors.

3. **AI Code Review Agent**

   - **Input:** Code diffs, static analysis report, coding standards document.
   - **Output:** List of review comments (suggestions, bugs, praise) mapped to specific lines of code.
   - **Success:** AI generates relevant, non-hallucinated comments; response time < 2 mins.
   - **Failure:** Fallback to rule-based checks if AI service is down; flag for human review if confidence is low.

4. **Human Review & Approval Gate**

   - **Input:** AI review comments, Test results, Static analysis report.
   - **Output:** Approval/Rejection decision, additional human comments.
   - **Success:** PR is either merged or returned for changes within 4 hours.
   - **Failure:** Escalate to team lead if PR sits unreviewed for > 4 hours.

5. **Artifact Build & Containerization**

   - **Input:** Merged code, Dockerfile/Build specs.

- **Output:** Immutable artifact (Docker image, binary) with version tag.
- **Success:** Build passes; artifact pushed to registry.
- **Failure:** Notify developer of build error; block deployment.

6. **Automated Deployment (Staging/Prod)**

- **Input:** Artifact version, Environment config (Kubernetes manifests, Terraform).
- **Output:** Running application in target environment.
- **Success:** Health checks pass (HTTP 200); pods are stable.
- **Failure:** Auto-rollback to previous version if health checks fail within 5 mins.

### Question 1.2: Parallelism & Critical Path

- **Parallel Steps:**
  - Step 2 (Static Analysis) and Step 3 (AI Review) can run simultaneously after Step 1.
  - Unit tests (part of Step 5 pre-requisites) can run in parallel with the review process.
- **Blocking Steps:**
  - Step 1 (Ingestion) blocks everything.
  - Step 4 (Human Approval) blocks Step 5 (Merge/Build).
  - Step 5 (Build) blocks Step 6 (Deployment).
- **Critical Decision Points:**
  - **Post-Analysis:** Should we even bother the AI? (If static analysis fails hard, reject immediately).
  - **Post-AI Review:** Is the confidence high enough to auto-approve trivial changes? (Risk policy).
  - **Post-Deployment:** Rollback or Stay? (Based on health metrics).

### Question 1.3: Key Handoffs

- **Step 1 -> 2 & 3:** Passes the *raw code diffs* and *metadata*.
- **Step 2 -> 3:** Passes *linting results* (so AI doesn't comment on things the linter already caught).
- **Step 3 -> 4:** Passes *suggested comments* and *summary* to the Human Reviewer.
- **Step 5 -> 6:** Passes the *Artifact ID (SHA/Tag)*. This is critical—we deploy *exactly* what we built.

---

# Part B: AI Prompting Strategy (30 points)

**Question 2.1:** For 2 consecutive major steps you identified, design specific AI prompts that would achieve the desired outcome. Include:

- System role/persona definition
- Structured input format
- Expected output format

- Examples of good vs bad responses
- Error handling instructions

**Question 2.2:** How would you handle the following challenging scenarios with your AI prompts:

- **Code that uses obscure libraries or frameworks**
- **Security reviews for code**
- **Performance analysis of database queries**
- **Legacy code modifications**

**Question 2.3:** How would you ensure your prompts are working effectively and getting consistent results?

# Response Part B:

**Question 2.1: AI Prompts for Two Major Steps**

**Step A: AI Code Reviewer (The "Critic")**

- **System Role:** "You are a Senior Software Engineer and Security Specialist. Your goal is to ensure code quality, security, and maintainability."
- **Input Format:**
```
{
  "file_path": "src/api/user_controller.ts",
  "diff": "...",
  "context": "User authentication logic using JWT.",
  "linter_errors": []
}
```
- **Prompt:**

  > "Review the following code diff. Focus ONLY on the changed lines but consider the surrounding context.
  >     1. Identify any SECURITY vulnerabilities (SQLi, XSS, IDOR).
  >     2. Identify PERFORMANCE bottlenecks (N+1 queries, expensive loops).
  >     3. Check for adherence to Clean Code principles (naming, modularity).
  > Output a JSON list of comments. Each comment must have a
  > 'line_number', 'severity' (INFO, WARN, CRITICAL), and 'message'."

- **Expected Output:** JSON list of review comments.
- **Error Handling:** If the diff is too large, request a summary or chunked processing.

**Step B: Test Case Generator (The "QA")**

- **System Role:** "You are a QA Automation Engineer expert in Jest and Pytest."
- **Input Format:** Code function/class definition.

- **Prompt:**

> "Analyze the following function. Generate 3 unit test cases:
>   1. A Happy Path (Positive) test.
>   2. A Negative test (Invalid input).
>   3. An Edge Case test (Boundary values, nulls).
> Provide the code in [Language]. Mock any external database or API calls."

- **Expected Output:** Executable test code snippet.

## Question 2.2: Handling Challenging Scenarios

- **Obscure Libraries:**
  - **Strategy:** Retrieval Augmented Generation (RAG). We will index the documentation of the obscure library and inject relevant snippets into the prompt context before asking for a review.
- **Security Reviews:**
  - **Strategy:** Use a specialized "Red Team" prompt persona. "Act as a malicious attacker trying to exploit this code." Also, chain-of-thought prompting: "First, list all user inputs. Second, trace where they go. Third, check for sanitization."
- **Performance Analysis:**
  - **Strategy:** Ask the AI to explain the Big-O complexity of the changed code. "Analyze the time and space complexity of this loop."
- **Legacy Code:**
  - **Strategy:** "Explain-then-Refactor". First ask the AI to summarize what the legacy code does in plain English. Once confirmed, ask it to suggest safe refactoring steps.

## Question 2.3: Ensuring Effectiveness

- **Golden Set Evaluation:** Maintain a dataset of 50 "bad" code snippets with known bugs. Run the AI against this set daily. If it misses a bug it previously caught, we have a regression.
- **User Feedback Loop:** Every AI comment in the PR has a "Helpful" / "Not Helpful" button. We track the "Acceptance Rate" of AI comments. Low acceptance rate triggers a prompt review.

---

# Part C: System Architecture & Reusability (25 points)

**Question 3.1:** How would you make this system reusable across different projects/teams? Consider:

- Configuration management
- Language/framework variations

- Different deployment targets (cloud providers, on-prem)
- Team-specific coding standards
- Industry-specific compliance requirements

**Question 3.2:** How would the system get better over time based on:

- False positive/negative rates in reviews
- Deployment success/failure patterns
- Developer feedback
- Production incident correlation

# Response Part C:

### Question 3.1: System Reusability

To make this system reusable across 50+ microservices and different teams:

1. **Configuration as Code (.github/ai-review.yml):**
   - Each repository will have a config file defining its language (Python/Node), strictness level (High/Low), and specific rules.
   - The system reads this config to dynamically adjust the AI prompts (e.g., "Use Python best practices").
2. **Plugin Architecture:**
   - The "Static Analysis" step should be a pluggable interface. We can swap `ESLint` for `Pylint` or `SonarQube` without changing the core workflow.
3. **Containerized Agents:**
   - The review logic runs in Docker containers. This ensures that whether the team uses Java 8 or Node 20, the review environment is consistent.
4. **Cloud-Agnostic Deployment:**
   - Use Terraform/OpenTofu for infrastructure. The deployment step accepts a "Target Provider" parameter (AWS, Azure, On-Prem) and applies the corresponding modules.

### Question 3.2: Continuous Improvement

- **False Positive/Negative Rates:**
  - We log every time a human explicitly dismisses an AI comment (False Positive). We analyze these logs to tweak the system instructions (e.g., "Stop complaining about missing docstrings in test files").
- **Deployment Patterns:**
  - If a deployment fails immediately after an AI approval, that PR is flagged as a "Missed Catch". We feed this example back into the prompt's few-shot examples: "Here is a bug you missed last time, don't miss it again."
- **Production Incidents:**

- Correlate PagerDuty alerts with recent deployments. If a specific microservice causes frequent incidents, the system automatically increases the "Review Strictness" for that repo (requiring 2 human reviewers instead of 1).

---

# Part D: Implementation Strategy (20 points)

**Question 4.1:** Prioritize your implementation. What would you build first? Create a 6-month roadmap with:

- MVP definition (what's the minimum viable system?)
- Pilot program strategy
- Rollout phases
- Success metrics for each phase

**Question 4.2:** Risk mitigation. What could go wrong and how would you handle:

- AI making incorrect review decisions
- System downtime during critical deployments
- Integration failures with existing tools
- Resistance from development teams
- Compliance/audit requirements

**Question 4.3:** Tool selection. What existing tools/platforms would you integrate with or build upon:

- Code review platforms (GitHub, GitLab, Bitbucket)
- CI/CD systems (Jenkins, GitHub Actions, GitLab CI)
- Monitoring tools (Datadog, New Relic, Prometheus)
- Security scanning tools (SonarQube, Snyk, Veracode)
- Communication tools (Slack, Teams, Jira)

# Response Part D:

**Question 4.1: Implementation Roadmap (6 Months)**

- **Month 1: MVP (The "Assistant")**
  - **Goal:** AI comments on PRs, but does not block merging.
  - **Scope:** GitHub Action that runs on PR open. Uses OpenAI API to suggest simple style fixes.
  - **Success Metric:** AI runs on 100% of new PRs; < 5s latency.
- **Month 2-3: Pilot Program (The "Gatekeeper")**
  - **Goal:** AI can block merges for Security issues.

- **Scope:** Integrate Static Analysis (SonarQube). Onboard 2 pilot teams. Add "Security" prompt persona.
- **Success Metric:** 30% reduction in human review time for pilot teams.
- **Month 4-5: Deployment Automation (The "Pipeline")**
  - **Goal:** Standardized CI/CD for all 50 services.
  - **Scope:** Build reusable GitHub Actions workflows for Build/Deploy. Implement "One-Click Rollback".
  - **Success Metric:** Deployment frequency doubles; Failure rate < 5%.
- **Month 6: Full Rollout & Optimization**
  - **Goal:** System-wide adoption.
  - **Scope:** Roll out to all teams. Enable auto-rollback based on Datadog metrics.
  - **Success Metric:** Review turnaround < 4 hours (Business Goal met).

## Question 4.2: Risk Mitigation

- **Risk: AI Making Incorrect Decisions (Hallucinations)**
  - **Mitigation: Human-in-the-loop is mandatory.** AI never auto-merges code in the first 6 months. It only *approves* or *requests changes*. A human must always press the final Merge button.
- **Risk: System Downtime During Deploy**
  - **Mitigation: Blue/Green Deployment.** We spin up the new version (Green) alongside the old (Blue). We only switch traffic if Green is healthy. If Green fails, traffic stays on Blue. Zero downtime.
- **Risk: Resistance from Developers**
  - **Mitigation: "Bot Noise" Control.** If the AI comments too much, devs will hate it. We will tune the prompt to be "Conservative"—only comment if 90% sure. We will also brand it as a "Helper" not a "Policeman".

## Question 4.3: Tool Selection

- **Code Review: GitHub** (Industry standard, great API).
- **CI/CD: GitHub Actions** (Native integration with code, easy to template).
- **AI Model: OpenAI GPT-4o** (Best reasoning capability for code) or **Claude 3.5 Sonnet** (Excellent coding performance).
- **Static Analysis: SonarQube** (Covers security, bugs, and code smells for many languages).
- **Monitoring: Datadog** (Unified view of logs, metrics, and traces).