# 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** I will walk through this with four essential steps.

**Step 1:** The Intake and Context Agent

- **Input:** The code changes (PR Diff) and the project's Readme.
- **Output:** A summary of what the code is trying to do (e.g. "This updates the page animations").
- **Success Criteria:** Correctly identifies which part of the system.
- **Failure Handling:** If the code is too messy to read, it notifies the users so that a professional can run a check manually.

**Step 2:** The Security & Quality Audit

- **Input:** The code changes + a list of "Security Rules" (OWASP).
- **Output:** A checklist of passes/fails (e.g., "Pass: No hardcoded passwords. Fail: Missing error handling on line 42").
- **Success Criteria:** Catches critical bugs before a human even sees it.
- **Failure Handling:** If the AI is unsure, it marks the line as "Needs Senior Review."

**Step 3:** Automated Deployment Orchestrator

- **Input:** The "Passed" audit report from the security and quality audit.
- **Output:** Commands to push the code to "Staging" (the testing site) or "Production" (the live site).
- **Success Criteria:** Code reaches the server without crashing.
- **Failure Handling:** If the server doesn't respond, it stops the deployment immediately.

**Step 4:** The Health Monitor (Rollback Agent)

- **Input:** Server metrics (like "Is the site slow?" or "Are there errors?").
- **Output:** A "Keep Live" or "Rollback" command.

- **Success Criteria:** If the site crashes, it automatically puts the old version back in under 60 seconds.
- **Failure Handling:** Alerts the engineering team on Slack immediately if a rollback fails.

**Question 1.2**

- **Blocking:** Step 3 must wait for step 2. The security of the code has to be confirmed before it gets deployed.
- **Parallel:** Security and coding style can be checked at the same time. -**Critical Decision Point:** Between Step 2 and 3. This is the "Gate." If the AI gives a low confidence score, a human must click "Approve."

**Question 1.3** The data that needs to be passed across each stage are as follows:

- **From 1 to 2:** The context (i.e What the code is for.).
- **From 2 to 3:** The security token, proof that the code is safe.
- **from 3 to 4:** The "Version ID" (So the monitor knows which version to roll back to if things break).

---

# 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: Designing the Prompts

**Prompt 1: The Intake and Context Agent System Persona:** You are a Software Engineer and Code Reviewer, and your goal is to identify code changes in the project. You are precise and are to carry out your task with scrutiny.

**Input:**

- `Project_context` (i.e python/Flask app)
- `Code_Diff` (Actual code changes)

**Expected Output format:** Text- Straightforward and detailed description of changes made (e.g "This updates the login button.")

**Good vs. Bad:**

- Good- "The search algorithm in `algo.py` was changed from linear to binary search. The button style in `styles.css` was updated, the button colour was changed from grey to black."
- Bad- "The files `algo.py` and `styles.css` where changed."

**Error Handling:**

**Prompt 2: The Security & Quality Audit Agent System Persona:** You are an expert Full-Stack Security Engineer and Code Reviewer. Your goal is to identify vulnerabilities and logic flaws in code diffs. You are pedantic, precise, and prioritize security over features.

**Input:**

- `Project_context` (i.e python/Flask app)
- `Code_Diff` (Actual code changes)
- `Security_Standards` (e.g. OWASP)

**Expected Output format:** ```json { "risk_level": "High/Med/Low", "issues": [{"line": 24, "type": "SQL Injection", "fix": "Use parameterized queries"}], "confidence_score": 0.95 }

**Good vs. Bad:**

- Good- "Specific line numbers and actionable fixes.
- Bad- "The code looks okay" or general advice without context.

**Error Handling:** If the code diff is incomplete or unparseable, return an error code `ERR_INVALID_DIFF` and request the full file.

## Question 2.2:Handling Challenging Scenarios

- **Obscure Libraries:** I would use a "RAG" (Retrieval-Augmented Generation) approach. The prompt would instruct the AI to first search the provided library documentation (attached as context) before reviewing the code.

- **Security Reviews:** I would use "Chain of Thought" prompting. I'll ask the AI to: "1. List all entry points for user input. 2. Trace how that input reaches the database. 3. Identify where sanitization is missing."

- **Performance Analysis:** I would feed the AI the `EXPLAIN ANALYZE` output from the database. This gives the AI the actual execution plan, not just the code, to see where the "bottleneck" is.

- **Legacy Code:** I would provide the AI with a "Legacy Context Map." The prompt would say: "Do not suggest modern Python 3.12 features; this environment is restricted to Python 2.7. Focus only on logic fixes."

## Question 2.3: Ensuring Consistency

To ensure the AI doesn't give different answers every time, I would:

1. **Lower the "Temperature":** Set the AI model temperature to 0.0 or 0.1 to make it predictable and "boring" rather than creative.
2. **Give Examples:** Include 3-5 examples of "Perfect Reviews" within the prompt so the AI has a template to copy.
3. **Validate Output:** Use a separate script (or Pydantic in Python) to verify that the AI's response is valid JSON before it moves to the next step.

---

# 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: Making the System Reusable

To make the system work across 50+ microservices and diverse teams, I would implement a "Configuration-as-Code" model:

- **Centralized Policy Engine:** Instead of hardcoding rules, use a global config file (YAML/JSON). Each team can "subscribe" to certain rules (e.g., "Financial Team" turns on strict HIPAA compliance checks, while "Internal Tools" uses standard checks).

- **Adapter Pattern for CI/CD:** Build "adapters" for different platforms. The core AI logic stays the same, but the "handshake" changes whether the team uses GitHub Actions, GitLab, or Jenkins.

- **Prompt Templating:** Use a library like LangChain to manage prompts. The system detects the programming language (Python, Java, Go) and automatically injects the correct language-specific security rules into the AI's "persona."

- **Secret Management:** Integrate with tools like HashiCorp Vault or AWS Secrets Manager to ensure the AI system never handles raw credentials, keeping it compliant with industry standards (SOC2/ISO27001).

### Question 3.2: Continuous Improvement

The system should get "smarter" every week by closing the feedback loop:

- **The "Human-in-the-Loop" Signal:** When a developer clicks "Approve" on a PR that the AI flagged as "High Risk," the system logs this as a False Positive. This data is fed back into the prompt engineering team to refine the AI's instructions.

- **Production Correlation:** If a deployment passes the AI review but causes a production crash, the system performs an Automated Post-Mortem. It analyzes the

"missed" bug and adds a new "Guardrail" to the AI's checklist to prevent that specific bug from ever slipping through again.

- **A/B Testing Prompts:** eriodically run two versions of the AI reviewer (e.g., Prompt A vs. Prompt B) and measure which one results in fewer developer "dismissals" and faster approval times.

- **Developer Sentiment:** Include a simple "Was this review helpful?" button. If a specific agent gets low ratings, it is automatically flagged for retraining.

---

## 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: Prioritization & 6-Month Roadmap

My implementation strategy follows a "crawl-walk-run" approach to build trust with the engineering team while delivering immediate value.

- **Months 1-2: MVP (Minimum Viable Product)**

- **Focus:** Automating the "Security & Quality Audit" for a single high-velocity Python team.

- **Goal:** Reduce manual review time for style and common security flaws.

- **Success Metrics:** Average PR cycle time reduced by 30%; 0 critical security leaks in the pilot service.

- **Months 3-4: Pilot Expansion & Automation**

- **Focus:** Integrate the AI Agent with GitHub Actions and implement the "Health Monitor" for automated rollbacks in the Staging environment.

- **Goal:** Establish the "automated gate" where AI can block a deployment if security or health checks fail.

- **Success Metrics:** 90% accuracy in automated rollback triggers; 50% reduction in total review time across 5 teams.

- **Months 5-6: Full Scale & Platformization**

- **Focus:** Scale to 50+ microservices and add support for Java and Go. Launch a "Self-Service Dashboard" for teams to customize their AI review rules.

- **Goal:** Standardize the deployment lifecycle across the entire organization.

- **Success Metrics:** < 4-hour average PR review time company-wide; > 80% developer satisfaction rate.

## Question 4.2: Risk Mitigation

- **AI Logic Errors:** I will implement a "Confidence Threshold." If the AI's confidence score is below 85%, it is forced to flag a human senior reviewer rather than making a solo decision.
- **System Downtime:** The architecture will use a **"Fail-Open"** mechanism. If the AI service is unreachable, the pipeline reverts to manual human approval so that critical hotfixes are never blocked by the tool itself.
- **Team Resistance:** To prevent friction, we will launch a "Transparency Dashboard" showing exactly why the AI makes certain decisions, and provide a "One-Click Dispute" button for developers to challenge AI feedback.

- **Compliance/Audit:** Every decision, prompt, and output will be logged in an immutable audit trail to satisfy SOC2 and industry-specific compliance requirements.

## Question 4.3: Tool Selection

To ensure reliability and speed, I will integrate the system with the following industry-standard tools:

| Category | Tool | Purpose |
| --- | --- | --- |
| Code Review | GitHub / GitLab | Primary interface for code hosting and PR comments via Webhooks. |
| CI/CD | GitHub Actions | Orchestrates the movement of code between Dev, Staging, and Prod. |
| Monitoring | Datadog / Prometheus | Provides the real-time telemetry (error rates, latency) for rollback decisio |
| Security | Snyk / SonarQube | Used as the primary "engine" for scanning, with AI used to interpret and s |
| Communication | Slack | Real-time alerts for the engineering team regarding deployment status o |