

Product Requirements Document (PRD)

ACR-QA v2.0: Language-Agnostic Code Review Platform

Ahmed Mahmoud Abbas

Student ID: 222101213

King Salman International University (KSIU)

Supervisor: Dr. Samy AbdelNabi

November 23, 2025

Document Information

Field	Value
Project Name	ACR-QA v2.0 (Automated Code Review & Quality Assurance)
Owner	Ahmed Mahmoud Abbas (Student ID: 222101213)
Supervisor	Dr. Samy AbdelNabi
Institution	King Salman International University (KSIU)
Timeline	October 2025 - June 2026 (8 months)
Version	2.0 (Platform Version)
Last Updated	November 23, 2025

Contents

1 Executive Summary	4
1.1 Product Vision	4
1.2 Problem Statement	4
1.3 Core Innovation (v2.0 Differentiators)	4
2 Product Objectives & Success Metrics	4
2.1 Primary Goals	4
2.2 Academic Requirements	5
3 User Personas & Use Cases	5
3.1 Primary Personas	5
3.1.1 Persona 1: University Instructor (Dr. Sarah)	5
3.1.2 Persona 2: Small Team Tech Lead (Omar)	5
3.1.3 Persona 3: Open-Source Maintainer (Fatima)	6
3.2 User Workflows	6
3.2.1 Workflow 1: Student Submits PR (Primary)	6
3.2.2 Workflow 2: Team Lead Reviews Dashboard	6
3.2.3 Workflow 3: Maintainer Configures Rules	6
4 Functional Requirements	6
4.1 Core Features (MVP - Phase 1, Nov-Jan 2026)	6
4.1.1 F1: Python Code Analysis	6
4.1.2 F2: Canonical Findings Schema	7
4.1.3 F3: RAG-Enhanced AI Explanations	7
4.1.4 F4: GitHub PR Integration	8
4.1.5 F5: Provenance Database	9
4.1.6 F6: Dashboard & Reporting	9
4.2 Extended Features (Phase 2, Feb-Jun 2026)	10
4.2.1 F7: Multi-Language Support	10
4.2.2 F8: Evaluation Framework	10
4.2.3 F9: User Study Tools	10
4.2.4 F10: Feedback Loop & Learning	11
5 Non-Functional Requirements	11
5.1 Performance	11
5.2 Scalability	11
5.3 Security & Privacy	11
5.4 Reliability	11
5.5 Usability	12
6 Technical Architecture	12
6.1 System Components	12
6.2 Technology Stack	12
6.3 Data Models	13
7 Implementation Roadmap	13
7.1 Phase 1: Foundation (Oct-Jan 2026) - CURRENT	13
7.2 Phase 2: Expansion (Feb-Jun 2026)	14

8 Success Criteria & Acceptance Tests	14
8.1 MVP Acceptance (End of Phase 1)	14
8.1.1 Test 1: GitHub PR Integration	14
8.1.2 Test 2: Canonical Schema	14
8.1.3 Test 3: RAG Explanations	15
8.1.4 Test 4: Evaluation	15
8.1.5 Test 5: User Study	15
9 Risk Management	15
10 Open Questions & Decisions Needed	16
10.1 Immediate (Week 1)	16
10.2 Phase 2 (Jan-Mar)	16
11 Appendices	16
11.1 Glossary	16
11.2 References	16
11.3 Document Change Log	16

1 Executive Summary

1.1 Product Vision

ACR-QA v2.0 is a language-agnostic, on-premises code review platform that automatically detects bad practices, security vulnerabilities, design anti-patterns, and style violations in pull requests. It uses Retrieval-Augmented Generation (RAG) with Cerebras AI to provide evidence-grounded, natural language explanations that help developers understand and fix issues. The MVP focuses on analyzing pull request diffs only, not entire repositories, to keep scope feasible for an 8-month academic project.

1.2 Problem Statement

- Current Pain:** Code review quality varies by reviewer availability; manual reviews miss security issues; commercial tools cost \$10k-50k/year; cloud-based tools can't handle proprietary code
- Target Users:** University instructors (grading student PRs), small dev teams (5-20 engineers), open-source maintainers, technical recruiters
- Key Gap:** Existing tools lack explanations or use generic AI that hallucinates incorrect guidance

1.3 Core Innovation (v2.0 Differentiators)

- Canonical Findings Schema:** Universal JSON format normalizes outputs from disparate tools (Ruff, Semgrep, ESLint) across languages
- RAG-Enhanced Explanations:** Evidence-grounded prompts reduce hallucinations 42-68% vs. direct LLM calls
- Provenance-First Architecture:** Stores raw tool outputs, LLM prompts/responses, and user feedback for reproducible evaluation
- Adapter SDK:** Pluggable language support (Python first, JavaScript/Java next) via standardized interface

2 Product Objectives & Success Metrics

2.1 Primary Goals

Goal	Metric	Target	Timeline
Multi-Language Platform	Languages supported	supported Python (100%), +1 language	Nov 2025 / Mar 2026
Detection Quality	Precision (high-severity rules)	$\geq 70\%$	Feb 2026
	False positive rate	<30%	Feb 2026
AI Explanation Quality	User rating (1-5 scale)	≥ 3.0 median	Feb 2026
	LLM vs template preference	LLM rated ≥ 0.5 higher	Feb 2026

Goal	Metric	Target	Timeline
Performance	Analysis latency (PR <200 lines)	≤ 90 seconds	Jan 2026
Deployment	On-prem setup time	≤ 30 minutes	Apr 2026
Cost	Total recurring cost	\$0 (zero)	Ongoing
PR Review Experience	All findings posted as inline PR comments	100% of detected issues appear as GitHub review comments	Jan 2026

2.2 Academic Requirements

1. Working software system with Docker Compose packaging
2. Evaluation report: precision/recall on 80+ labeled issues
3. User study: 8-10 participants rating explanation usefulness
4. Adapter SDK documentation proving extensibility
5. Demonstration video showing end-to-end workflow

3 User Personas & Use Cases

3.1 Primary Personas

3.1.1 Persona 1: University Instructor (Dr. Sarah)

- **Context:** Teaches Software Engineering to 120 students; receives 300+ PRs/semester
- **Pain:** Can't manually review all PRs; students repeat same mistakes; no time for personalized feedback
- **Jobs-to-be-Done:** Automatically grade PRs for code quality; provide consistent feedback; track student progress
- **Success Criteria:** Reduces review time from 10min/PR to 2min/PR; students understand why code failed

3.1.2 Persona 2: Small Team Tech Lead (Omar)

- **Context:** Leads 8-person dev team at Egyptian startup; can't afford SonarQube Enterprise
- **Pain:** Junior devs push bad code; manual reviews miss security issues; cloud tools violate data policy
- **Jobs-to-be-Done:** Enforce quality gates on PRs; educate juniors via AI explanations; deploy on-prem
- **Success Criteria:** Catches SQL injection before production; costs \$0/month; runs on existing server

3.1.3 Persona 3: Open-Source Maintainer (Fatima)

- **Context:** Maintains popular Python library; receives PRs from 100+ external contributors
- **Pain:** Contributors ignore style guide; duplicate code gets merged; explaining issues wastes time
- **Jobs-to-be-Done:** Auto-comment on PRs with guidance; reduce back-and-forth; maintain code quality
- **Success Criteria:** 50% fewer “please fix style” comments; contributors self-correct before re-submission

3.2 User Workflows

3.2.1 Workflow 1: Student Submits PR (Primary)

1. Student opens PR with homework solution
2. GitHub webhook triggers ACR-QA analysis (30-90s)
3. System posts comment: “Found 3 issues: [AI explanations with line numbers]”
4. Student reads explanations, fixes code, pushes update
5. Re-analysis confirms fixes, instructor reviews only logic

3.2.2 Workflow 2: Team Lead Reviews Dashboard

1. Tech lead opens ACR-QA dashboard (React UI)
2. Views trend: “Security findings down 40% this month”
3. Clicks finding: sees code, AI explanation, false positive button
4. Marks 2 findings as FP (overly strict rules for their domain)
5. System adjusts thresholds for future PRs

3.2.3 Workflow 3: Maintainer Configures Rules

1. Maintainer edits `rules.yml` in repo
2. Adds custom rule: “No `requests.get()` without timeout”
3. Commits rule definition with rationale + example
4. Next PR triggers analysis using new rule
5. AI explanation cites the custom rule definition

4 Functional Requirements

4.1 Core Features (MVP - Phase 1, Nov-Jan 2026)

4.1.1 F1: Python Code Analysis

Description: Detect 5 categories of issues in Python code

Categories:

- **Bad Practices:** Mutable defaults, unused variables, dead code
- **Style Violations:** PEP 8 compliance, import ordering, line length
- **Design Smells:** Too many parameters (>5), large classes (>300 lines), high complexity (CC >10)
- **Security Issues:** SQL injection patterns, unsafe eval/exec, weak crypto
- **Code Duplication:** Token-based similarity (>80% match over 50+ tokens)

Tools Used: Ruff, Vulture, Radon, Semgrep, Bandit, jscpd

Input: Python files (.py) from PR diff

Output: Canonical findings with severity, confidence, line numbers. Severity rules: security issues and potential crashes = high, design smells affecting maintainability = medium, stylistic issues and minor formatting problems = low.

4.1.2 F2: Canonical Findings Schema

Description: Universal JSON format normalizing all tool outputs

Schema:

```
{
  "finding_id": "uuid-v4",
  "rule_id": "SOLID-001",
  "category": "design | security | style | duplication | unused",
  "severity": "high | medium | low",
  "confidence": 0.87,
  "file": "app/auth.py",
  "line": 42,
  "column": 10,
  "language": "python",
  "evidence": {
    "snippet": "def authenticate(user, pass, token, session, db):",
    "tool_output": {"radon": {"complexity": 12}},
    "context_before": ["line 39", "line 40", "line 41"],
    "context_after": ["line 43", "line 44", "line 45"]
  },
  "explanation": "AI-generated natural language...",
  "explanation_source": "llm | template",
  "timestamp": "2025-11-23T17:30:00Z"
}
```

Normalizer: Maps Ruff, Semgrep, Vulture, Radon, jscpd → canonical format

Rationale: Enables language-agnostic dashboard, cross-language comparisons, consistent API

4.1.3 F3: RAG-Enhanced AI Explanations

Description: Generate natural language explanations using Cerebras LLM with evidence-grounding

Process:

1. Load rule definition from `rules.yml` (description, rationale, remediation, examples)
2. Retrieve 3-6 lines of code context around issue
3. Construct prompt: “Explain using ONLY this rule definition and code context”

4. Call Cerebras API (llama3.1-8b, temperature=0.3, max tokens=150)
5. Validate: Check if response cites rule id; if not, use template fallback
6. Log prompt + response to provenance DB

Rules Catalog (rules.yml):

```
SOLID-001:
  name: "Too Many Parameters"
  category: "design"
  severity: "medium"
  description: "Functions with >5 parameters violate Single Responsibility"
  rationale: "Complex signatures indicate function does too much"
  remediation: "Extract parameters into dataclass or config object"
  example_good: |
    @dataclass
    class Config:
        user: str; token: str
        def auth(cfg: Config): ...
  example_bad: |
    def auth(user, pass, token, session, db): ...
```

Cost: ~\$0.0014 per PR analysis (50-200 explanations at \$0.60/1M tokens)

Fallback: Template-based explanation if LLM confidence <0.7 or API fails

4.1.4 F4: GitHub PR Integration

Description: Automatically analyze PRs and post findings as comments

Trigger Options:

- **GitHub Action (Phase 1):** Workflow file in .github/workflows/
- **Webhook Endpoint (Phase 2):** Flask/FastAPI server receives PR events
- **Manual trigger:** PR comment acr-qa review starts analysis (optional mode for demos)

Flow:

1. PR opened/updated → GitHub triggers action/webhook
2. Fetch PR diff via GitHub API (pygithub library)
3. Extract changed files + line ranges
4. Run analysis on changed code only (not entire repo)
5. Compute severity for each finding (high/medium/low) and sort comments by severity so the most critical issues appear first
6. Post findings as PR review comments with line annotations
7. Store PR metadata + findings in database

Comment Format:

```
**ACR-QA Detected Issue**
**Rule**: SOLID-001 (Too Many Parameters)
**Severity**: Medium
**File**: 'app/auth.py:42'
```

This function has 5 parameters, which violates the Single Responsibility Principle. Complex parameter lists indicate the function `is` doing too much `and` becomes hard to test `and` maintain.

****Suggested Fix**:**
Extract related parameters into a dataclass:

```
@dataclass
class AuthConfig:
    username: str
    password: str
    token: str

def authenticate(config: AuthConfig, session, db): ...

[Mark as False Positive](#) | [View Details](#)
```

Rate Limiting: Max 1 analysis/PR/minute to avoid spam

4.1.5 F5: Provenance Database

Description: PostgreSQL stores all analysis data for reproducibility and evaluation
Tables:

- **analyses:** PR metadata, timestamp, status, total findings count
- **findings:** Canonical finding objects (see F2 schema)
- **raw_outputs:** Original JSON from each tool (Ruff, Semgrep, etc.)
- **llm_interactions:** Prompts sent, responses received, model, temperature, cost
- **feedback:** User marks (false positive, helpful, unclear)

Retention: Unlimited (storage ~85MB for entire project)

Backup: Docker volume persists across container restarts

4.1.6 F6: Dashboard & Reporting

Description: Terminal UI + Markdown reports for viewing results
Dashboard Features (dashboard.py):

- Table view: All findings with severity, file, line, explanation preview
- Summary stats: Total findings, breakdown by category, cost spent
- Filter: By severity, category, language, date range
- Action buttons: Mark false positive, view provenance

Markdown Reports (generate_report.py):

- Executive summary (total issues, critical count)

- Findings list (grouped by file, severity)
 - AI explanations included
 - Trend comparison (if multiple analyses exist)
- Provenance Export (export_provenance.py):**

- JSON file with complete audit trail
- Includes: All tool outputs, LLM prompts/responses, timestamps, costs
- Used for academic evaluation and thesis defense

4.2 Extended Features (Phase 2, Feb-Jun 2026)

4.2.1 F7: Multi-Language Support

Priority Languages:

- JavaScript/TypeScript (March 2026)
- Java (Optional, April 2026)

JavaScript Adapter:

- Tools: ESLint, TypeScript compiler API, Semgrep, jscpd
- Rules: `rules_js.yml` with JS-specific patterns
- Same canonical schema + dashboard integration

Adapter SDK:

- Abstract base class: `LanguageAdapter`
- Required methods: `run_tools()`, `normalize_findings()`, `get_extensions()`
- Documentation: Step-by-step guide to add new languages

4.2.2 F8: Evaluation Framework

Seeded Dataset: 80-100 manually labeled issues

- 20 duplications, 20 style, 20 design, 20 security, 20 unused code
- Ground truth: True Positive (TP) or False Positive (FP)

Metrics Calculation (compute_metrics.py): Precision = $\frac{TP}{TP+FP}$ Recall = $\frac{TP}{TP+FN}$ F1 Score = $\frac{2 \times (Precision \times Recall)}{Precision + Recall}$

CI Integration: Nightly runs compute metrics on seeded dataset

Target: Precision $\geq 70\%$ for high-severity rules

4.2.3 F9: User Study Tools

Comparison Setup: 20 findings with dual explanations (LLM + template)

Google Form: Code snippet, two explanations (randomized order), rating scale

Questions:

1. “Rate Explanation A usefulness (1-5)”
2. “Rate Explanation B usefulness (1-5)”
3. “Which is clearer? A / B / Equal”
4. “Would you follow this guidance? Yes / No”

Target: 8-10 participants, $\geq 3.0/5.0$ median rating

4.2.4 F10: Feedback Loop & Learning

False Positive Marking: Developer clicks “Not a bug” in dashboard

Per-Repo Overrides: Store rule adjustments (e.g., “Disable SOLID-001 for auth.py”)

Threshold Tuning: If 80% of users mark rule FP → increase detection threshold

Version Tracking: Log rule version used for each finding (enables A/B testing)

5 Non-Functional Requirements

5.1 Performance

Metric	Target	Measurement
Analysis Latency	$\leq 90\text{ms}$ for PR < 200 lines	90th percentile
	<i>Scope: Larger PRs may take longer and are out-of-scope for formal evaluation.</i>	
LLM Response Time	$\leq 600\text{ms}$ per explanation	Median
Database Query Time	$\leq 100\text{ms}$ for dashboard load	95th percentile
Concurrent PRs	10 simultaneous analyses	Stress test

Table 2: Performance Requirements

5.2 Scalability

- **MVP Scope:** Single Docker Compose instance (1 worker)
- **Future:** Redis queue enables horizontal worker scaling (out of scope for graduation)
- **Storage Growth:** $\sim 70\text{KB}$ per PR $\times 1000$ PRs = 70MB (trivial)

5.3 Security & Privacy

- **On-Premises Deployment:** No code leaves customer infrastructure
- **API Keys:** Stored in .env file (not in Git); Docker secrets in production
- **LLM Data:** Code snippets sent to Cerebras API (documented in privacy policy)
- **Only minimal code context** (the offending snippet and a few surrounding lines) is sent to the LLM, never entire repositories, to reduce exposure risk
- **Local LLM Option:** Architecture supports offline mode with template explanations only

5.4 Reliability

- **Uptime:** Not applicable (on-prem, no SLA)
- **Error Handling:** Graceful degradation (LLM fails → use template)
- **Data Integrity:** PostgreSQL transactions ensure atomic saves

- **Backup:** Docker volume persists; users responsible for backup strategy

5.5 Usability

- **Setup Time:** 30 minutes from git clone to first analysis
- **Documentation:** README, architecture docs, API reference, video tutorial
- **Error Messages:** Plain English (no stack traces to end users)
- **Accessibility:** Terminal UI supports screen readers (basic)

6 Technical Architecture

6.1 System Components

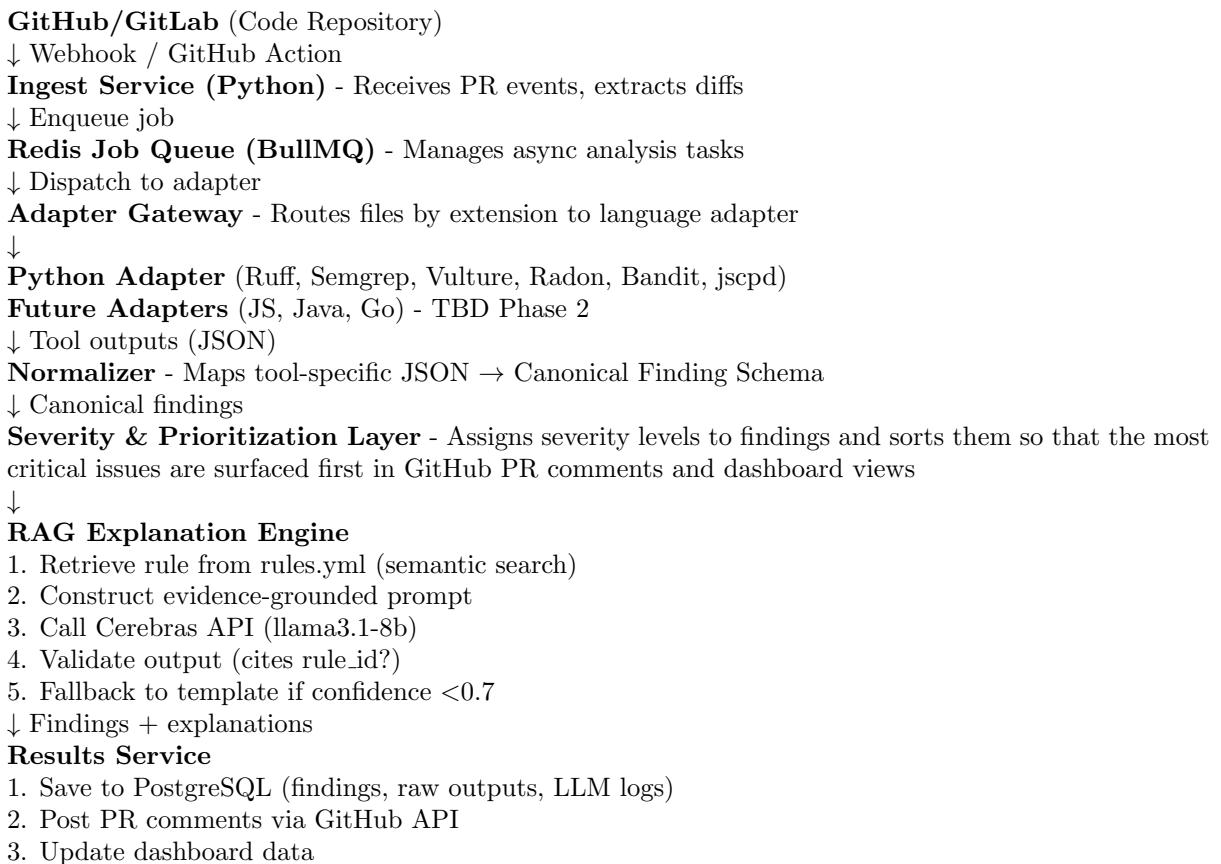


Figure 1: System Architecture

6.2 Technology Stack

Layer	Technology	Rationale
Language	Python 3.11+	Rich ecosystem, fast prototyping, AST support
Database	PostgreSQL 15+	JSONB for raw outputs, vector search for RAG

Layer	Technology	Rationale
Queue	Redis 7 + BullMQ	Async job processing, proven for CI/CD tools
LLM API	Cerebras (llama3.1-8b)	Free tier, 60-70× faster than OpenAI, \$0.60/1M tokens
Containerization	Docker Compose	On-prem deployment, zero-config startup
Static Analysis	Ruff, Semgrep, Vulture, Radon, Bandit, jscpd	Industry-standard tools, broad coverage
Dashboard	Rich (Python terminal UI)	Fast development, no frontend complexity
VCS Integration	GitHub API (pygithub)	Primary platform, GitLab in Phase 2
Testing	pytest	Standard Python testing framework

6.3 Data Models

Canonical Finding (Core Data Structure):

```
@dataclass
class CanonicalFinding:
    finding_id: str          # UUID
    rule_id: str              # e.g., "SOLID-001"
    category: FindingCategory # Enum: design, security, style,
                               duplication, unused
    severity: Severity        # Enum: high, medium, low
    confidence: float         # 0.0 - 1.0
    file: str                 # Relative path
    line: int
    column: int
    language: str             # e.g., "python"
    evidence: Evidence
    explanation: str
    explanation_source: str   # "llm" or "template"
    timestamp: datetime

@dataclass
class Evidence:
    snippet: str              # Code line causing issue
    tool_output: dict          # Raw JSON from tool
    context_before: List[str]  # 3 lines before
    context_after: List[str]   # 3 lines after
```

Severity is defined as: high = security or bug risk (e.g., injections, unsafe calls, crashes), medium = design and maintainability issues (e.g., long functions, too many parameters), and low = style and cosmetic issues (e.g., formatting, naming). This prioritization is used to order findings in PR comments and reports.

7 Implementation Roadmap

7.1 Phase 1: Foundation (Oct-Jan 2026) - CURRENT

Month	Deliverable	Status
Oct 2025	Python adapter, Docker setup, database schema	Complete
Nov 2025	Canonical schema, normalizer, GitHub Action, evidence-grounded prompts	In Progress
Dec 2025	RAG retrieval, dashboard polish, provenance export, severity scoring, PR comment templates	Planned
Jan 2026	Seeded dataset, evaluation metrics, user study prep, manual acr-qa review trigger	Planned

7.2 Phase 2: Expansion (Feb-Jun 2026)

Month	Deliverable
Feb 2026	User study execution, precision/recall computation, threshold tuning
Mar 2026	JavaScript/TypeScript adapter, multi-language demo
Apr 2026	Production hardening, on-prem deployment guide, monitoring
May 2026	Performance optimization, stress testing, Grafana dashboards
Jun 2026	Final report, demo video, documentation, submission

8 Success Criteria & Acceptance Tests

8.1 MVP Acceptance (End of Phase 1)

8.1.1 Test 1: GitHub PR Integration

- Open test PR with 10 code issues
- System analyzes within 90 seconds
- Posts 10 comments with AI explanations
- Each comment cites a rule id
- Findings are ordered so that high-severity issues appear at the top of the PR review
- Provenance DB logs all LLM interactions

8.1.2 Test 2: Canonical Schema

- Run Ruff, Semgrep, Vulture on sample code
- Normalizer produces findings with universal rule ids
- Database stores findings in canonical format
- Dashboard displays unified view across tools

8.1.3 Test 3: RAG Explanations

- Load 20 rules from `rules.yml`
- Generate explanations for 20 diverse findings
- 100% of explanations cite correct rule id
- <10% require template fallback

8.1.4 Test 4: Evaluation

- Label 80 findings as TP/FP
- Compute precision: $\geq 70\%$ for high-severity
- Compute recall: $\geq 60\%$ overall
- Document methodology in thesis

8.1.5 Test 5: User Study

- 8-10 participants complete rating form
- Median usefulness rating $\geq 3.0/5.0$
- LLM explanations rated ≥ 0.5 higher than templates
- Statistical significance (t-test, $p < 0.05$)

9 Risk Management

Risk	Probability	Impact	Mitigation
Cerebras API downtime	Medium	High	Template fallback; store all prompts for replay
GitHub rate limits	Low	Medium	Cache PR diffs; batch comment posts
Low precision ($<70\%$)	Medium	High	Tune thresholds conservatively; focus on high-confidence rules
User study recruitment fails	Medium	Medium	Expand to online (Reddit, GitHub); offer small incentive
Scope creep (too many languages)	High	High	Gate: Python must hit 70% precision before adding JS
Database migration issues	Low	Low	Version schema; test migrations in staging
Enterprise feature creep (e.g., full codebase context engine, advanced analytics)	Medium	High	Limit scope to PR diffs, 1-2 languages, and clearly documented non-goals; defer full-repo context and enterprise features beyond graduation

10 Open Questions & Decisions Needed

10.1 Immediate (Week 1)

- **GitHub Action vs Webhook?** → Recommendation: Action first (simpler)
- **Rules.yml structure?** → Recommendation: YAML (version controlled), DB later
- **Template fallback format?** → Recommendation: Jinja2 templates with same structure as LLM output

10.2 Phase 2 (Jan-Mar)

- **Second language: JS or Java?** → Depends on user study feedback
- **Self-hosted LLM option?** → Optional; Ollama + Llama 3.1 8B documented
- **GitLab support priority?** → Low unless user requests
- **Should future versions add full-repository context analysis, or remain PR-diff focused to keep complexity and costs low?**

11 Appendices

11.1 Glossary

Canonical Finding Normalized detection result in universal JSON schema

RAG (Retrieval-Augmented Generation) LLM technique that injects retrieved context into prompts

Provenance Complete audit trail of analysis (tool outputs, prompts, responses)

Adapter Language-specific module that runs tools and normalizes outputs

False Positive (FP) Detection flagged as issue but is actually correct code

True Positive (TP) Detection correctly identifies a real code issue

11.2 References

1. CustomGPT (2025) “RAG API vs Traditional LLM APIs” - 42-68% hallucination reduction
2. IEEE (2023) “Towards Multi-Language Static Code Analysis” - Adapter pattern validation
3. Semgrep Documentation - Pattern-based rule engine design
4. Johnson et al. (2013) “Why don’t developers use static analysis?” - User study methodology

11.3 Document Change Log

Date	Version	Changes	Author
Nov 23, 2025	1.0	Initial PRD creation	Ahmed Abbas

End of Product Requirements Document