



TUNISIAN REPUBLIC
Ministry of Higher Education and Scientific Research
University of Carthage
**National Institute of Applied Sciences and
Technology**



Graduation Project

*In order to obtain the
National Engineering Diploma*

Specialty : Software Engineering

Enforcing Best Practices with LLM-IDE Integration

Presented by

Arij KOUKI

Hosted by

Google

INSAT Supervisor : Ms. YOUSSEF Rabaa
Company Supervisor : Ms. LOPEZ Irene

Presented on : -/-/2025

JURY

M. President FLEN (President)
Ms. Reviewer FLENA (Reviewer)

Academic Year : 2024/2025



TUNISIAN REPUBLIC
Ministry of Higher Education and Scientific Research
University of Carthage
**National Institute of Applied Sciences and
Technology**



Graduation Project
In order to obtain the
National Engineering Diploma
Specialty : Software Engineering

Enforcing Best Practices with LLM-IDE Integration

Presented by
Arij KOUKI

Hosted by

Company Supervisor	University Supervisor

Academic Year : 2024/2025

Acknowledgements

I wish to extend my deepest gratitude and appreciation to everyone who has contributed significantly to the successful completion of this project.

My sincere thanks go to Irene Lopez and Andrew Xue, my host and co-host at Google Zurich, for their invaluable guidance, support, and trust throughout my internship. Their mentorship and kindness made this experience not only enriching but truly transformative.

Thank you to the YouTube Developer Infrastructure team for their warm welcome, collaborative spirit, and continuous encouragement, and to my mentor Veronica Radu, whose insights and support greatly contributed to my personal and professional growth.

I also wish to thank Mrs. Rabaa Youssef, my academic supervisor, for accompanying me in this final academic milestone.

To the distinguished members of the jury, I am grateful for your time and consideration in reviewing my work. I hope this report lives up to the standards expected of a graduation project.

To the professors of the National Institute of Applied Sciences and Technology, thank you for playing a vital role in shaping my academic and professional foundation.

And to my family and friends, your unwavering belief in me has carried me through the most demanding moments. Thank you for your love, patience, and presence.

Finally, a quiet note of gratitude to myself for the perseverance and dedication that made this journey possible.

Résumé

Ce projet a été réalisé chez Google Zurich dans le cadre d'un Diplôme National d'Ingénieur en Génie Logiciel. Il explore l'intégration de l'intelligence artificielle générative dans le processus de développement logiciel en incorporant un agent basé sur un Large Language Model (LLM) au sein d'un Environnement de Développement Intégré (IDE) interne.

Cet agent effectue une analyse approfondie du code afin de détecter des violations complexes ou subjectives que les outils d'analyse statique traditionnels peuvent négliger. En générant des explications claires et compréhensibles ainsi que des suggestions concrètes, il aide les développeurs, notamment ceux contribuant à YouTube, à maintenir une haute qualité de code et à respecter les bonnes pratiques. Intégré de manière transparente au sein du flux de travail via une extension de l'IDE, l'agent améliore la productivité et contribue à la réduction de la dette technique sans perturber l'expérience de développement.

Ce travail met en évidence le potentiel des outils assistés par l'IA pour transformer l'expérience des développeurs et ouvre des perspectives pour l'avenir des environnements de développement intelligents.

Mots-clés : Génie Logiciel, Intelligence Artificielle Générative, Intégration IDE, Qualité du Code, Productivité des Développeurs

Abstract

This project was carried out at Google Zurich as part of a National Engineering Diploma in Software Engineering. It investigates the integration of generative artificial intelligence into the software development process by embedding a Large Language Model (LLM)-powered agent within an internal Integrated Development Environment (IDE).

The agent performs in-depth code analysis to detect complex or subjective violations that traditional static analysis tools may overlook. By generating clear, human-readable explanations and actionable suggestions, it supports developers, particularly those contributing to YouTube, in maintaining high code quality and adhering to best practices. Seamlessly integrated into the development workflow through an IDE extension, the agent enhances productivity and helps reduce technical debt without disrupting the coding experience.

This work demonstrates the potential of AI-assisted tooling to transform the developer experience and raises broader implications for the future of intelligent software engineering environments.

Keywords: Software Engineering, Generative AI, IDE Integration, Code Quality, Developer Productivity

Contents

Résumé	2
Abstract	3
List of Figures	7
List of Tables	8
List of Acronyms	9
General Introduction	1
I Project Overview	2
1 Host Company: Google	2
1.1.1 Presentation	2
1.1.2 Products and services	2
1.1.3 Focus Area	3
1.1.3.1 YouTube	3
1.1.3.2 YouTube Developer Infrastructure Team	3
2 Project Overview	4
1.2.1 Project Context	4
1.2.2 Problem Statement	4
1.2.3 Proposed Solution	5
3 Work Methodology	5
1.3.1 Agile Development Approach	5
1.3.2 Kanban Workflow	6
1.3.3 Development Process	7
1.3.4 Project Timeline	8
II Business understanding and project requirements	11
1 Foundations and Motivation	11
2.1.1 Foundational AI Concepts	11
2.1.1.1 Large Language Models (LLMs)	12
2.1.1.2 AI Agents	14

2.1.2	The AI Revolution in Software Engineering	14
2.1.2.1	AI-Enhanced SDLC	15
2.1.2.2	AI Applications Across the SDLC	15
2	State of the Art and Existing Solutions	16
2.2.1	State of the Art	16
2.2.1.1	AI-Powered Code Analysis Tools	16
2.2.1.2	Capabilities of Market Solutions	16
2.2.2	Existing Environment-Specific Approaches	17
2.2.2.1	Current Feedback Mechanisms	17
2.2.2.2	Gap Analysis: Market vs. Environment Solutions	17
2.2.2.3	Impact of Delayed Feedback	19
2.2.3	Opportunity for Framework-Specific AI Solutions	19
3	Project Requirements	19
2.3.1	Motivation for Use Case Analysis	20
2.3.1.1	System Use Cases	20
2.3.2	Functional Requirements	21
2.3.3	Non-Functional Requirements	21
III	System Design and Architecture	23
1	System Architecture Overview	23
3.1.1	High-Level Architecture	23
3.1.2	System Workflow	24
2	Components Design	25
3.2.1	LLM Best Practices Agent	25
3.2.1.1	Agent Architecture Options	26
3.2.1.2	Processing Strategy Options	27
3.2.1.3	Core Tools Architecture	27
3.2.1.4	Integration with LLM Infrastructure	28
3.2.2	Evaluation Framework for Architecture Decisions	28
3.2.2.1	Test Suite Framework	28
3.2.2.2	LLM-as-a-Judge Methodology	29
3.2.2.3	Key Metrics	29
3.2.2.4	Candidate Architectures and Strategies	29
3.2.3	YouTube IDE Extension	30
3.2.3.1	Extension Architecture	30

3.2.3.2	User-Triggered vs. Automatic Analysis Design Decision	30
3.2.3.3	User Interface Design	31
3.2.3.4	User Interaction Flow	32
3	Data Models and Interfaces	34
3.3.1	Input/Output Specification	34
3.3.2	Convention Data Management	34
IV Implementation		37
1	Working Environment	37
4.1.1	Development Infrastructure	37
4.1.1.1	Internal IDE	37
4.1.1.2	Internal RPC Playground	37
4.1.1.3	Google Colab	38
4.1.1.4	Internal Repository Integration	38
4.1.2	Project Management and Documentation	38
4.1.2.1	Internal Version Control	38
4.1.2.2	Internal Code Review Platform	38
4.1.2.3	Internal Project Management System	39
4.1.2.4	Google Docs	39
2	Technologies	39
4.2.1	Backend Technologies (AI Agent)	39
4.2.1.1	Python Programming Language	39
4.2.1.2	YouTube DevInfra Agent Framework	40
4.2.1.3	Internal AI Platform	40
4.2.1.4	LLM Libraries and Frameworks	40
4.2.2	Frontend Technologies (IDE Extension)	40
4.2.2.1	TypeScript Programming Language	40
4.2.2.2	VS Code Extension API	41
3	Realization	41
4.3.1	Agent Realization	41
4.3.1.1	Evaluation Results and Architecture Decision	41
4.3.1.2	Agent Implementation Overview	43
4.3.1.3	Core Tools Implementation	44
4.3.1.4	Processing Workflow	46
4.3.1.5	Resilience and Error Handling	47

Table des Matières

4.3.1.6	Convention Data Management	47
4.3.2	Extension Integration	48
4.3.2.1	Extension Architecture	48
4.3.2.2	User Interaction	50
4.3.2.3	Stale Diagnostics Handling	52
4.3.2.4	Resilience and Communication	55
Conclusion and Perspectives		56
References		57
Appendix : Miscellaneous remarks		59

List of Figures

1.1	Google Logo	2
1.2	Overview of some Google products	3
1.3	YouTube Logo	3
1.4	Kanban Workflow	6
1.5	Project Development Cycle	7
1.6	Project Timeline - Detailed Schedule	9
2.1	Hierarchy of AI Technologies	12
2.2	Chronological Overview of Large Language Models (LLMs)	13
2.3	AI Agent Architecture and Orchestration Workflow	14
2.4	Current Developer Workflow Feedback Timeline	18
2.5	System Use Case Diagram	20
3.1	High-Level System Architecture	24
3.2	System Interaction Sequence Diagram	25
3.3	Comparison of Executable Agent vs. ReAct Agent Architectures	26
3.4	YouTube IDE Extension User Interaction Flow: Complete developer journey from analysis trigger to fix application	33
4.1	Latency Performance Across 12 Test Cases	42
4.2	Cost Analysis Across 12 Test Cases	42
4.3	Accuracy Performance Across 12 Test Cases	43
4.4	Agent Processing Activity Diagram	47
4.5	System Architecture: IDE Extension, Proxy, and Backend Communication Flow	49
4.6	VS Code Interface: Editor Actions (Illustrative)	50
4.7	VS Code Command Palette (Illustrative)	51
4.8	VS Code Notification Interface (Illustrative)	51
4.9	VS Code Diagnostics Interface (Illustrative)	52
4.10	Stale Diagnostics Handling: Two-Tiered System (Illustrative)	54

List of Tables

2.1	Summary of Project Requirements	22
3.1	Comparison of User-Triggered vs. Automatic Analysis Approaches	31
4.1	Single-Violation File Performance Comparison	42

List of Acronyms

AI	Artificial Intelligence
API	Application Programming Interface
CI/CD	Continuous Integration/Continuous Deployment
DL	Deep Learning
FIFO	First In, First Out
IDE	Integrated Development Environment
JSON	JavaScript Object Notation
LLM	Large Language Model
ML	Machine Learning
QA	Quality Assurance
RPC	Remote Procedure Call
SDLC	Software Development Life Cycle
VS	Visual Studio

General Introduction

The software engineering industry is experiencing a major shift driven by the rapid evolution of artificial intelligence and the growing demand for scalable, high-quality code development practices. As development teams grow and systems become more complex, ensuring consistent code quality and adherence to best practices presents an ongoing challenge, especially in large organizations managing massive codebases across distributed teams.

Traditional static analysis tools and linters, while helpful, often fall short when it comes to identifying nuanced or subjective coding issues that depend on context or internal guidelines. In fast-paced development environments, engineers need intelligent, responsive tools that go beyond rule-based checks to provide deeper insights and real-time guidance, without adding friction to their daily workflows.

This graduation project explores the integration of generative artificial intelligence into modern Integrated Development Environments (IDEs) to support software engineers in their day-to-day coding activities. The research investigates how Large Language Models (LLMs) can be leveraged to perform in-depth code analysis, detect complex violations, and offer clear, contextual suggestions for improvement. By embedding intelligent assistance directly within the development workflow, this work aims to enhance code quality, reduce technical debt, and support developer productivity at scale.

This report begins by establishing the theoretical foundation and examining the current state of AI-assisted development tools. It then presents the design and implementation of an intelligent code analysis system, followed by an evaluation of its effectiveness in real-world development scenarios. The work contributes to the broader understanding of how artificial intelligence can transform software engineering practices and improve developer experience.

Chapter I

Project Overview

Introduction

This opening chapter establishes the foundation of the graduation project by introducing the host company and defining the project's scope. We examine the organizational context, outline the main objectives and challenges, and present the methodological framework that guides the development process.

1 Host Company: Google

1.1.1 Presentation

Founded in 1998, Google LLC is a global leader in technology and innovation [1]. As a subsidiary of Alphabet Inc., Google's mission is to organize the world's information and make it universally accessible and useful. Guided by values such as innovation, accessibility, sustainability, and user trust, Google establishes itself as one of the most influential companies shaping the digital era. Its culture emphasizes collaboration, diversity, inclusion, and impact-driven engineering, enabling continuous leadership in research and product development.

Figure 1.1 shows the Google logo.



Figure 1.1 – Google Logo

1.1.2 Products and services

Google offers a broad ecosystem of products and services that touch nearly every aspect of digital life. Among its flagship consumer products are Google Search, Maps, Gmail, Chrome, and the Android operating system, serving billions of users daily.

Beyond consumer services, Google develops enterprise and cloud-based solutions such as Google Cloud Platform and Google Workspace, as well as advanced AI systems like Vertex AI. The company also invests in hardware, including Pixel devices, Nest smart home products, and ChromeOS.

These products reflect Google's commitment to connecting people, improving productivity, and driving digital transformation worldwide.

Figure 1.2 provides an overview of some Google products.



Figure 1.2 – Overview of some Google products

1.1.3 Focus Area

1.1.3.1 YouTube

Acquired by Google in 2006 [2], YouTube becomes the world's leading video-sharing platform, serving more than two billion logged-in users monthly. It empowers individuals to create, share, and discover video content globally while sustaining a vibrant creator economy. From a technical perspective, YouTube integrates video processing, recommendation systems, live streaming, advertising, and trust and safety to deliver a seamless experience across devices.

Figure 1.3 displays the YouTube logo.



Figure 1.3 – YouTube Logo

1.1.3.2 YouTube Developer Infrastructure Team

Within YouTube's engineering organization, the Developer Infrastructure (Dev Infra) team supports thousands of engineers building the platform. The team develops tooling, automation,

and guidelines that improve efficiency, reliability, and consistency in software development. By maintaining developer velocity and quality at scale, the Dev Infra team contributes directly to YouTube's ability to innovate and grow.

2 Project Overview

1.2.1 Project Context

This project develops within the scope of a development infrastructure team dedicated to supporting developers by providing tools and extensions that enhance their daily workflows. As part of this mission, the team explores how artificial intelligence can be leveraged to assist developers in maintaining code quality and adhering to best practices. The goal is to investigate how large language models (LLMs) can complement traditional approaches by offering more intelligent and context-aware guidance directly within the IDE.

In parallel, this work also constitutes the mandatory end-of-studies-internship project required for obtaining the software engineering degree at the National Institute of Applied Science and Technology, providing both academic and practical significance.

1.2.2 Problem Statement

In large-scale software development environments, maintaining uniform adherence to **internal framework-specific guidelines** across multiple teams is crucial for ensuring code quality, consistency, and long-term maintainability. While modern development environments provide assistance for general programming practices or widely used frameworks, they lack intelligent support for the nuanced, evolving rules of internal frameworks. As a result, developers often receive feedback on internal best practices only during code reviews, after significant effort has already been invested. This delayed feedback cycle leads to inefficiencies such as rework, slower iterations, and frustration among teams who must refactor code that was previously considered complete. The absence of real-time, context-aware guidance tailored to internal frameworks leaves developers navigating complex design decisions without adequate support, leading to technical debt, inconsistent quality, and higher onboarding complexity. Addressing this gap requires solutions that proactively enforce internal framework best practices during the coding phase, providing timely and context-specific feedback directly within the IDE.

1.2.3 Proposed Solution

This project introduces an **AI-assisted feedback system integrated directly into the coding workflow**. The solution is designed to address the challenges outlined above through three key capabilities:

- **Shift-Left Feedback:** Provide developers with earlier, context-aware guidance during the coding phase, ensuring that issues are detected and addressed well before code reviews.
- **Framework-Specific Best Practice Enforcement:** Surface adherence to internal framework guidelines early in the development process, going beyond syntax and correctness.
- **Reduced Review Burden:** Shift part of the best practice enforcement from manual reviews to the authoring stage, allowing reviews to focus on higher-level insights.

By integrating intelligent, framework-aware feedback directly into the coding workflow, this solution aims to minimize rework, improve adherence to internal standards, and accelerate development velocity.

3 Work Methodology

1.3.1 Agile Development Approach

Agile software development, as defined by Beck et al. [3], emphasizes "individuals and interactions over processes and tools, working software over comprehensive documentation, customer collaboration over contract negotiation, and responding to change over following a plan." This methodology is adopted to support iterative development and maintain flexibility in responding to evolving requirements. According to Martin [4], agile practices enable continuous integration of feedback, ensuring that each increment of work aligns with both technical goals and the broader product vision. Testing, validation, and code reviews are incorporated throughout the process to maintain high quality, while frequent collaboration provides clarity and shared ownership of outcomes. Agile principles complement the focus on engineering excellence, including rigorous design reviews, thorough testing, robust code reviews, and DevOps-enabled automation.

1.3.2 Kanban Workflow

Kanban, as described by Anderson [5], is "a method for managing knowledge work with an emphasis on just-in-time delivery while not overloading the team members." The dynamic workload of the development infrastructure team, including feature requests, bug fixes, and maintenance tasks, is managed using this Kanban workflow. According to Kniberg and Skarin [6], Kanban enables teams to visualize tasks and limit work in progress, preventing bottlenecks and allowing quick focus shifts to urgent issues when necessary. Work is structured into stages to maintain clear coordination while allowing the flexibility to adapt priorities as requirements evolve.

The Kanban workflow includes the following stages, as illustrated in Figure 1.4:

- **Backlog:** Prioritized collection of feature requests, enhancements, and bug fixes.
- **Research & Design:** Assessment of technical feasibility and preparation of design specifications.
- **Development:** Implementation and integration of features into the system.
- **Review & Testing:** Code review, unit tests, and integration tests to ensure quality and correctness.
- **Deployment:** Release of validated features to developer environments.

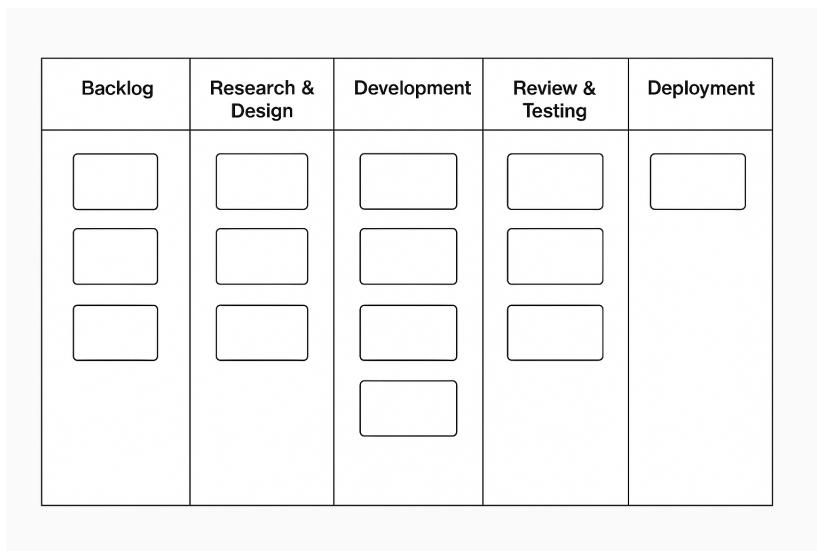


Figure 1.4 – Kanban Workflow

1.3.3 Development Process

The project follows an iterative engineering cycle designed to balance thorough planning with incremental delivery, as illustrated in Figure 1.5. This cycle guides the work through structured stages, supported by dedicated tools and regular collaboration:

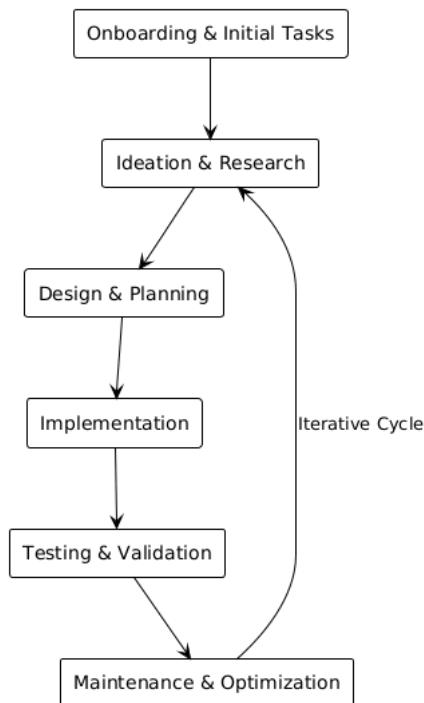


Figure 1.5 – Project Development Cycle

The main stages, as depicted in the figure, include:

- **Onboarding and Initial Tasks:** The project begins with a structured onboarding phase, where familiarization with internal tools, repositories, and coding standards combines with the resolution of assigned bugs. This phase ensures a smooth transition into the team's workflow and provides early practical experience.
- **Ideation and Research:** Following onboarding, the project enters an exploration phase to clarify objectives, gather requirements, and investigate potential solution directions. Independent research complements collaborative discussions to assess feasibility and align priorities.
- **Design and Planning:** A detailed design document is authored to present technical choices, architectural considerations, and the proposed workflow. The document un-

dergoes iterative review by engineers, and the final approved plan is transferred to the internal task management system for structured tracking and prioritization.

- **Implementation:** Development is performed in small, reviewable increments using the internal development environment within the company-wide repository. Each change is submitted with unit tests and validated through manual and AI-assisted code reviews.
- **Testing and Validation:** Functionality and reliability are verified continuously. Automated unit tests ensure correctness at the component level, while integration reviews and structured evaluations validate the behavior within the larger system.
- **Maintenance and Optimization:** Refactoring, bug fixes, and updates are performed throughout development, particularly as dependencies evolve or methods become deprecated. This ensures that the solution remains consistent, maintainable, and aligned with evolving standards.

Collaboration is supported through a structured communication rhythm, combining regular syncs with the host and co-host, weekly team meetings, and occasional cross-team discussions. This cadence provides timely feedback, clear guidance, and alignment on shared dependencies throughout the iterative cycle.

1.3.4 Project Timeline

The project spans four months (May 5 – September 5, 2025). Work is scheduled based on business priorities and technical dependencies. Early weeks focus on research and design, followed by implementation, testing, and iterative refinement.

Figure 1.6 shows the detailed project timeline.

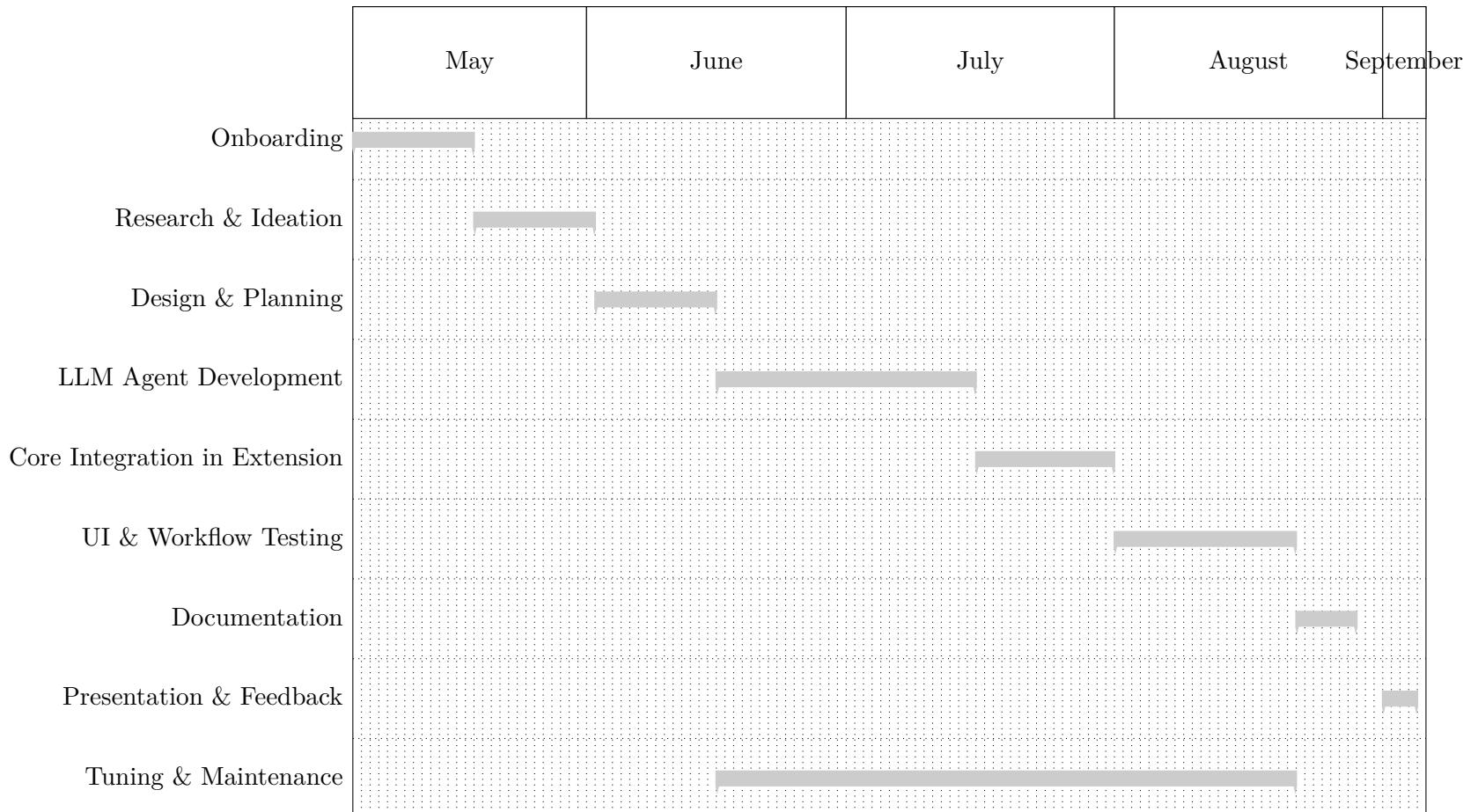


Figure 1.6 – Project Timeline - Detailed Schedule

Conclusion

In summary, this chapter outlines the project's context by presenting the host company, defining the problem, and clarifying the main objectives and challenges. It also describes the methodology chosen to guide the work, which provides the basis for the technical developments detailed in the following chapters.

Chapter II

Business understanding and project requirements

Introduction

This chapter establishes the theoretical and contextual foundation of the project. It begins with an overview of AI technologies such as Large Language Models (LLMs), Generative AI, and AI agents, highlighting their relevance to modern software engineering. Next, it presents the state of the art by reviewing existing solutions for code quality and best practices. Finally, it defines the project requirements, showing how this work addresses gaps by integrating AI-driven assistance into key development phases.

1 Foundations and Motivation

2.1.1 Foundational AI Concepts

Artificial Intelligence (AI) refers to computational systems capable of performing tasks that typically require human intelligence, including reasoning, learning, problem-solving, and decision-making [7]. AI encompasses a wide range of techniques with distinct capabilities and applications. Figure 2.1 illustrates a hierarchy of AI technologies.

II.1 Foundations and Motivation

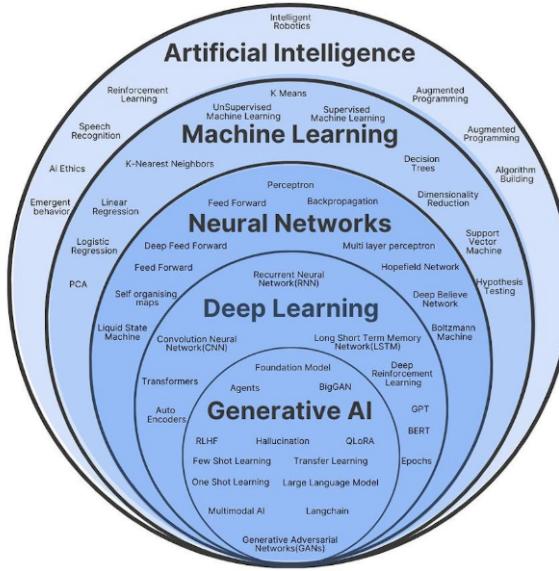


Figure 2.1 – Hierarchy of AI Technologies

Key categories include:

- **Symbolic AI:** Rule-based reasoning and knowledge representation systems.
- **Machine Learning (ML):** Algorithms that learn from data to improve task performance [8].
- **Deep Learning (DL):** Neural networks with multiple layers capable of modeling complex patterns [9].
- **Generative AI:** Systems that produce new content, such as text or code, by learning patterns from existing datasets [10].

2.1.1.1 Large Language Models (LLMs)

LLMs represent a breakthrough in generative AI, trained on extensive natural language and code corpora [11]. Unlike traditional tools, LLMs understand context, semantics, and intent, enabling complex reasoning across domains. Figure 2.2 provides a chronological overview of LLM development from 2018–2024.

II.1 Foundations and Motivation

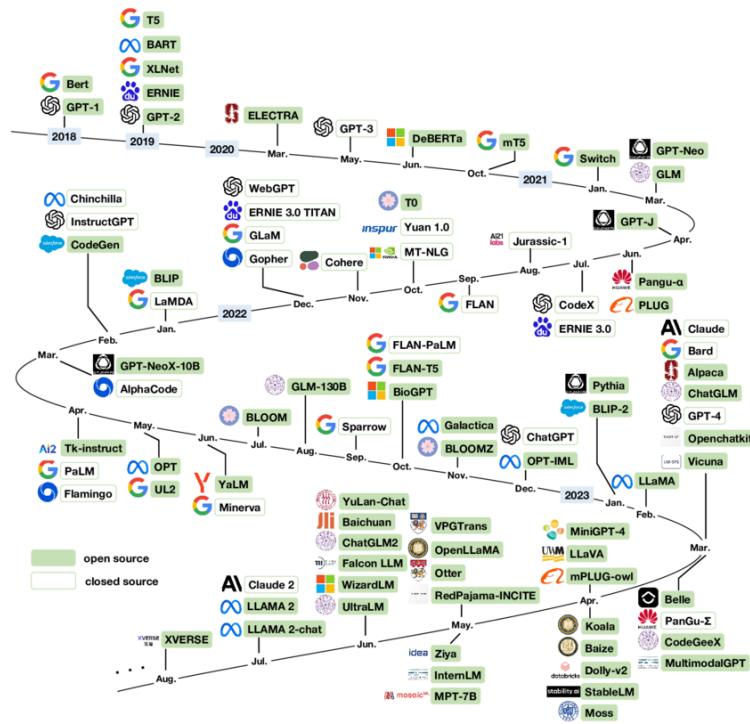


Figure 2.2 – Chronological Overview of Large Language Models (LLMs)

LLM capabilities include:

- Natural Language Understanding
- Pattern Recognition
- Content Generation
- Reasoning and Inference

In software development, these capabilities translate into:

- Contextual Code Analysis
- Intelligent Code Generation
- Explanatory Documentation
- Semantic Standards Enforcement

2.1.1.2 AI Agents

AI agents build on LLMs by autonomously reasoning, planning, and executing development tasks. They orchestrate LLM capabilities within workflows, providing seamless integration into IDEs, testing frameworks, and version control systems. Figure 2.3 illustrates a typical AI agent architecture and orchestration workflow.

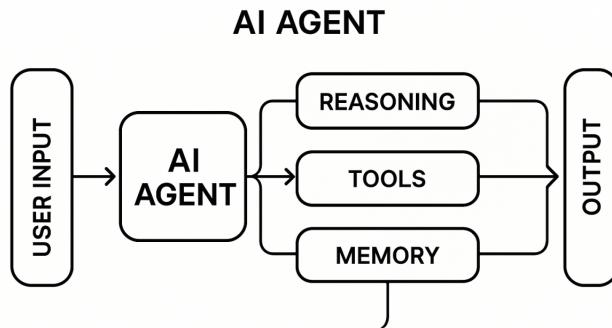


Figure 2.3 – AI Agent Architecture and Orchestration Workflow

Core components include:

- Code Analysis Engine
- Contextual Reasoning
- Development Tool Integration
- Contextual Adaptation via system prompts
- Feedback and Explanation System

Balancing AI Capabilities and Practical Considerations While LLM-powered systems provide strong *capabilities*—contextual understanding, automated analysis, and project-specific adaptation—they also present *practical constraints*, including computational cost, context window limits, accuracy variability, and integration complexity. The system design explicitly navigates these trade-offs to deliver reliable, low-latency guidance directly inside the IDE.

2.1.2 The AI Revolution in Software Engineering

Modern software development faces increasing complexity. Traditional practices often struggle to maintain code quality while meeting deadlines. AI technologies offer transformative opportunities by embedding intelligent assistance directly into the development workflow.

II.1 Foundations and Motivation

Industry adoption highlights this impact: surveys show AI-generated code accounts for a growing portion of development output in major organizations [12, 13]. AI integration addresses critical challenges such as maintaining code quality, reducing technical debt, and scaling development practices across teams.

2.1.2.1 AI-Enhanced SDLC

The AI-enhanced Software Development Life Cycle (SDLC) embeds intelligent assistance throughout all phases, directly influencing planning, design, coding, testing, and deployment.

Its transformative impact on development practices includes:

- **From Reactive to Proactive:** Feedback is delivered continuously during coding and design, preventing issues before they require refactoring.
- **From Inconsistent to Scalable:** AI agents provide uniform, expert-level guidance across teams and codebases.
- **From Static to Adaptive:** AI adapts to project-specific patterns, team preferences, and evolving best practices.
- **From Isolated to Integrated:** Guidance is embedded within workflows rather than treated as a separate phase, bridging planning, implementation, testing, and maintenance.

2.1.2.2 AI Applications Across the SDLC

AI can support specific phases of the SDLC as follows:

- **Requirements and Planning:** Estimate timelines, identify ambiguities, and translate requirements into technical specifications.
- **Design and Architecture:** Recommend patterns, detect anti-patterns, and ensure compliance with organizational standards.
- **Implementation:** Context-aware code completion, best practice enforcement, bug detection, and refactoring guidance.
- **Testing and Quality Assurance:** Generate test cases, identify edge cases, and prioritize test execution.
- **Deployment and Maintenance:** Monitor performance, predict issues, and recommend optimizations.

II.2 State of the Art and Existing Solutions

Overall, the AI-enhanced SDLC shifts development from reactive, fragmented practices to a proactive, adaptive, and integrated approach. This aligns directly with the goals of this project: delivering real-time, framework-specific guidance within the IDE to improve developer efficiency and maintain code quality.

2 State of the Art and Existing Solutions

Building on the foundations above, this section examines both market-available solutions and environment-specific approaches to understand the current landscape of code quality enforcement and developer assistance, identifying gaps that the proposed solution aims to fill.

2.2.1 State of the Art

The software development market offers a variety of AI-powered tools and platforms designed to enhance code quality and developer productivity. These solutions represent the current state of the art in intelligent development assistance.

2.2.1.1 AI-Powered Code Analysis Tools

Key categories include:

- **AI Code Assistants:** GitHub Copilot, Amazon CodeWhisperer, and Tabnine provide AI-powered code completion and generation, helping developers write code more efficiently.
- **AI-Powered IDEs:** Tools like Cursor and Claude Code integrate AI assistance directly into the coding workflow, offering context-aware code generation, refactoring, and intelligent suggestions.
- **Static Analysis Platforms:** Solutions such as SonarQube, CodeClimate, and DeepCode offer automated code quality analysis with AI-enhanced pattern detection.
- **AI Code Review Tools:** Platforms like PullRequest.com and CodeRabbit provide AI-assisted code review, offering automated suggestions and quality assessments.

2.2.1.2 Capabilities of Market Solutions

Typical features include:

II.2 State of the Art and Existing Solutions

- **General Code Analysis:** Broad pattern recognition and quality assessment across multiple languages and frameworks.
- **AI-Powered Suggestions:** Recommendations for code improvements, refactoring, and general best practices.
- **IDE Integration:** Seamless integration with widely used environments such as VS Code, IntelliJ, and Eclipse.

2.2.2 Existing Environment-Specific Approaches

Within our development environment, software engineers rely on a combination of modern AI-powered tools and traditional mechanisms to maintain code quality.

2.2.2.1 Current Feedback Mechanisms

- **Code Reviews:** Human reviewers provide context-aware feedback on design quality, readability, maintainability, and adherence to standards. Feedback is high-level but often delayed and resource-intensive.
- **Presubmit Checks:** Automated scripts enforce style guides, compilation correctness, and basic safety constraints. They are fast but primarily focus on surface-level checks.
- **Coding Assistant:** Internal IDE features AI-powered code completion, generation, and basic suggestions.
- **Linters:** IDE-integrated linters provide real-time style and deprecation warnings but do not enforce complex best practices.
- **Rule-Based Checks:** Enforce coding conventions and naming schemes consistently but cannot reason about complex or context-dependent practices.

2.2.2.2 Gap Analysis: Market vs. Environment Solutions

While both market solutions and environment-specific approaches provide valuable capabilities, they leave significant gaps in enforcing internal framework-specific best practices during the coding phase.

II.2 State of the Art and Existing Solutions

Limitations of Market Solutions

- **Generic Analysis:** Lack deep understanding of internal framework-specific conventions.
- **External Dependency:** May require sharing internal code, posing security or privacy concerns.
- **Limited Customization:** Cannot easily enforce project-specific practices or architectural patterns.

Limitations of Environment Solutions

- Coding Assistant and linters focus on general guidance, not framework-specific practices.
- Rule-based and presubmit checks are reactive rather than proactive, offering feedback after coding rather than during authoring.

Figure 2.4 illustrates how current mechanisms provide feedback at different points in the workflow, highlighting the gap during active coding.

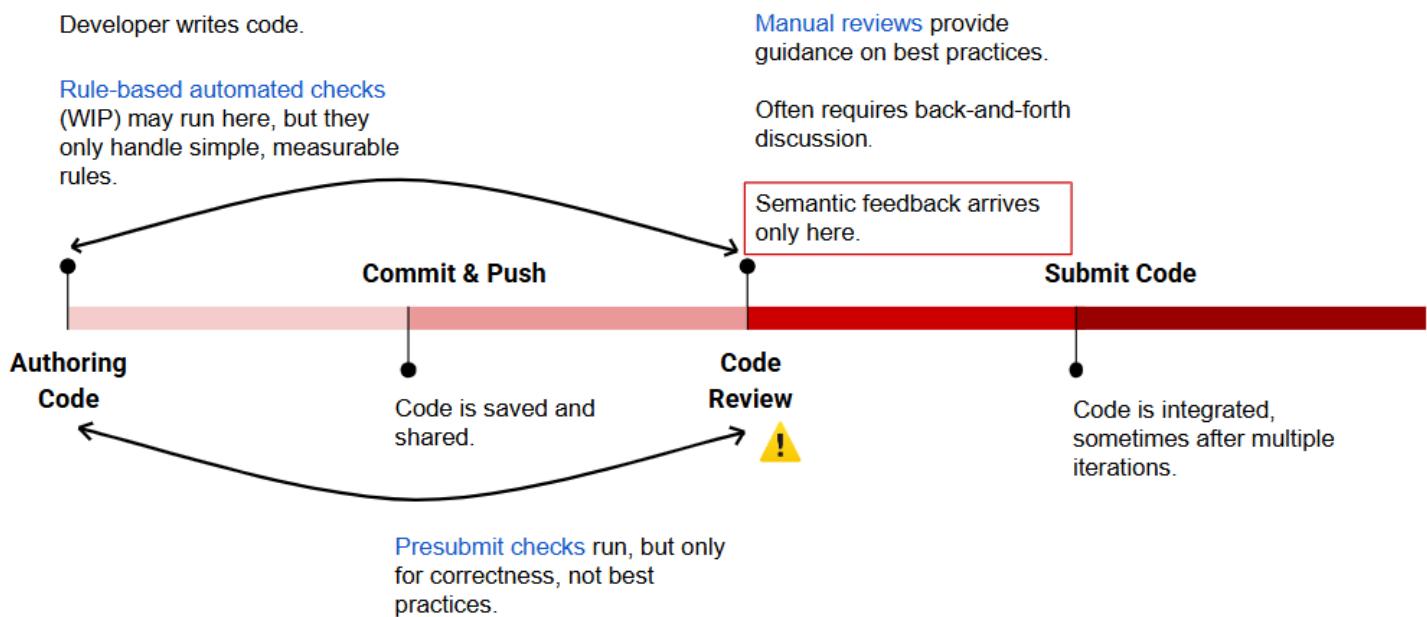


Figure 2.4 – Current Developer Workflow Feedback Timeline

2.2.2.3 Impact of Delayed Feedback

The current workflow introduces several challenges:

- **Technical Debt Accumulation:** Delayed feedback leads to inconsistencies and long-term maintenance costs.
- **Prolonged Review Cycles:** Multiple iterations are required to fix issues discovered late.
- **Developer Frustration:** Repetitive corrections decrease productivity.
- **Inconsistent Quality Standards:** Lack of real-time guidance leads to variable adherence across teams.
- **Increased Costs:** Late discovery of issues exponentially increases remediation effort.

2.2.3 Opportunity for Framework-Specific AI Solutions

The analysis reveals a clear opportunity for integrating framework-specific AI solutions. Combining the intelligence of market tools with the specificity required for internal frameworks, these solutions can deliver real-time, context-aware guidance directly in the coding phase.

Such a solution addresses the observed gap:

- Real-time adherence to internal framework best practices.
- Integration directly into the developer workflow.
- Proactive feedback that reduces review effort and technical debt.

This motivates the development of an LLM-powered, IDE-integrated assistant that provides intelligent, framework-aware support, complementing existing tools while filling the critical gap in real-time guidance.

3 Project Requirements

We now translate the identified opportunity into concrete requirements. The proposed solution integrates LLM-powered assistance directly into the developer workflow within the IDE to provide real-time, actionable guidance while maintaining performance, usability, and scalability.

2.3.1 Motivation for Use Case Analysis

Understanding how developers will interact with the system is critical for designing effective functionality. Use case analysis provides a structured approach to capture user-system interactions, identify key actors, and define the boundaries of the system. This ensures that the system delivers targeted support precisely where it is needed during the active coding phase.

2.3.1.1 System Use Cases

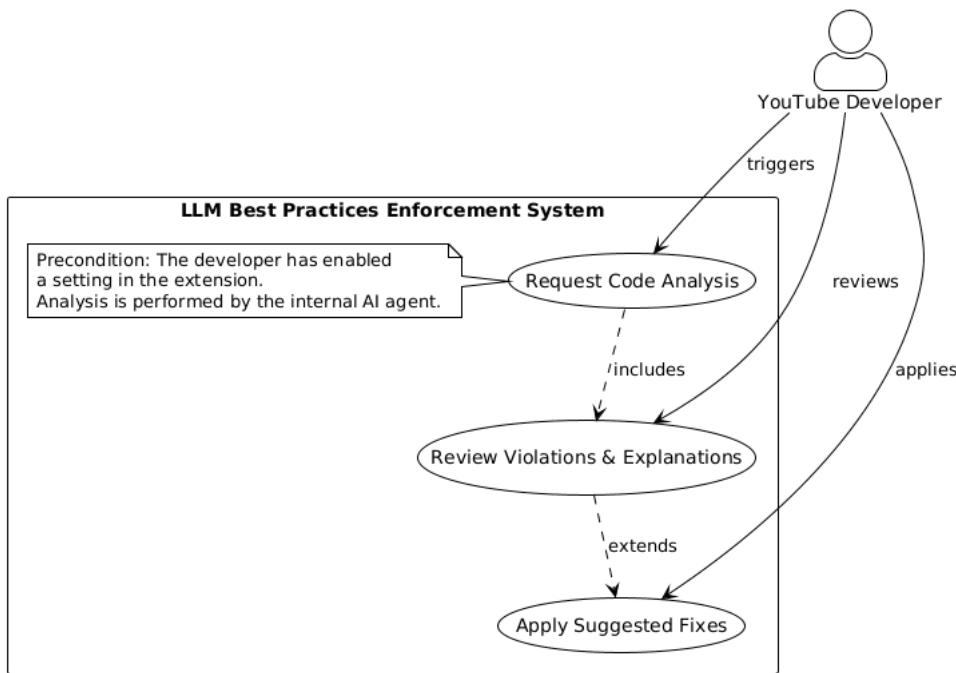


Figure 2.5 – System Use Case Diagram

Figure 2.5 illustrates the core functionality of the system from the perspective of YouTube developers. The primary actor is the **YouTube Developer**, who interacts with the system through three main use cases:

1. **Request Code Analysis:** Trigger real-time evaluation of code for adherence to framework best practices.
2. **Review Violations and Explanations:** Inspect flagged violations, with clear, context-specific explanations provided by the system.
3. **Apply Suggested Fixes (Optional):** Accept, modify, or reject actionable AI-generated suggestions to improve code quality.

II.3 Project Requirements

This design ensures that developers remain in control of their workflow while benefiting from comprehensive, intelligent feedback. Optional scenarios, such as applying fixes, allow flexibility and respect developer autonomy. The use case diagram focuses on core interactions, leaving authentication and administrative workflows out of scope.

2.3.2 Functional Requirements

Based on the use cases and the gaps identified in existing tools, the system must provide the following core functionalities:

- **Detect Framework Violations:** Identify violations of internal YouTube framework best practices in real-time.
- **Provide Contextual Explanations:** Deliver developer-friendly explanations that clarify why a particular pattern is problematic.
- **Generate Actionable Fixes:** Suggest concrete AI-driven solutions aligned with framework conventions.
- **Enable Developer Interaction:** Allow developers to accept, reject, or modify suggested fixes, maintaining workflow control.
- **Maintain Contextual Relevance:** Ensure feedback remains accurate and properly anchored as code evolves.
- **Seamless IDE Integration:** Embed within the existing internal IDE without disrupting normal development processes.

2.3.3 Non-Functional Requirements

In addition to core functionality, the system must satisfy broader quality criteria:

- **Performance:** Provide near real-time feedback to avoid interrupting workflow.
- **Scalability:** Efficiently handle large codebases and multiple simultaneous users.
- **Maintainability:** Support modular updates, addition of new rules, and AI model improvements.
- **Reliability:** Operate robustly in production with minimal downtime.

II.3 Project Requirements

- **Security and Privacy:** Comply with organizational policies, safeguarding code and data.
- **Usability:** Deliver concise, context-aware, and minimally intrusive feedback.
- **Extensibility:** Allow easy addition of new rules, models, or integrations.

Table 2.1 summarizes functional and non-functional requirements.

Table 2.1 – Summary of Project Requirements

Requirement Type	Description
Functional	Detection, explanations, fixes, interaction, context, IDE integration
Non-Functional	Performance, scalability, maintainability, reliability, security, usability, extensibility

These requirements directly address the gaps identified in both market and environment-specific solutions. By embedding proactive, context-aware feedback into the coding workflow, the system:

- Reduces framework-specific errors during development.
- Improves adherence to YouTube internal standards.
- Enhances developer productivity by providing guidance in real-time.

The use case analysis connects these requirements to real developer interactions, ensuring that the system's functionality aligns with actual workflow needs and supports effective adoption.

Conclusion

This chapter formalizes the project requirements, bridging the gap between the problem analysis and system design. It introduces the rationale for capturing developer interactions via use cases, defines functional and non-functional requirements, and establishes concrete acceptance criteria. Together, these elements provide a solid foundation for the subsequent system design chapter.

Chapter III

System Design and Architecture

Introduction

This chapter presents the system design and architecture of the LLM-powered best practices enforcement system. Building on the business understanding and requirements established in Chapter , it details technical design decisions, architectural patterns, and system components that enable real-time, intelligent feedback for YouTube framework development.

The design follows traditional software engineering principles while incorporating modern AI technologies. This chapter covers the overall system architecture, component design, data models, and integration patterns that form the foundation of the implemented solution.

1 System Architecture Overview

3.1.1 High-Level Architecture

The system architecture is designed to integrate seamlessly into the developer's existing workflow while providing intelligent, context-aware feedback. It consists of two main components:

- **IDE:** The developer's workspace, including the YouTube IDE Extension, which works with the currently open file and annotates results directly in the editor.
- **AI Agent Framework:** The processing layer containing the LLM Best Practices Agent, which performs code analysis and generates best practice suggestions.

III.1 System Architecture Overview

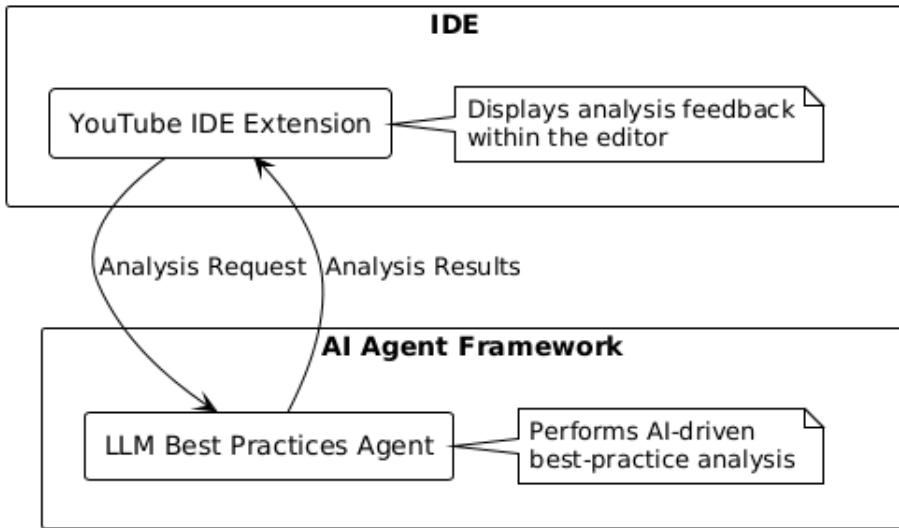


Figure 3.1 – High-Level System Architecture

Figure 3.1 illustrates the separation between the IDE extension and the AI processing backend. The IDE provides immediate access to analysis capabilities, while the AI Agent Framework handles computationally intensive tasks. This separation allows independent scaling of AI capabilities without impacting IDE responsiveness.

3.1.2 System Workflow

The system operates through a streamlined workflow beginning when a developer triggers analysis via the IDE Extension (Figure 3.2).

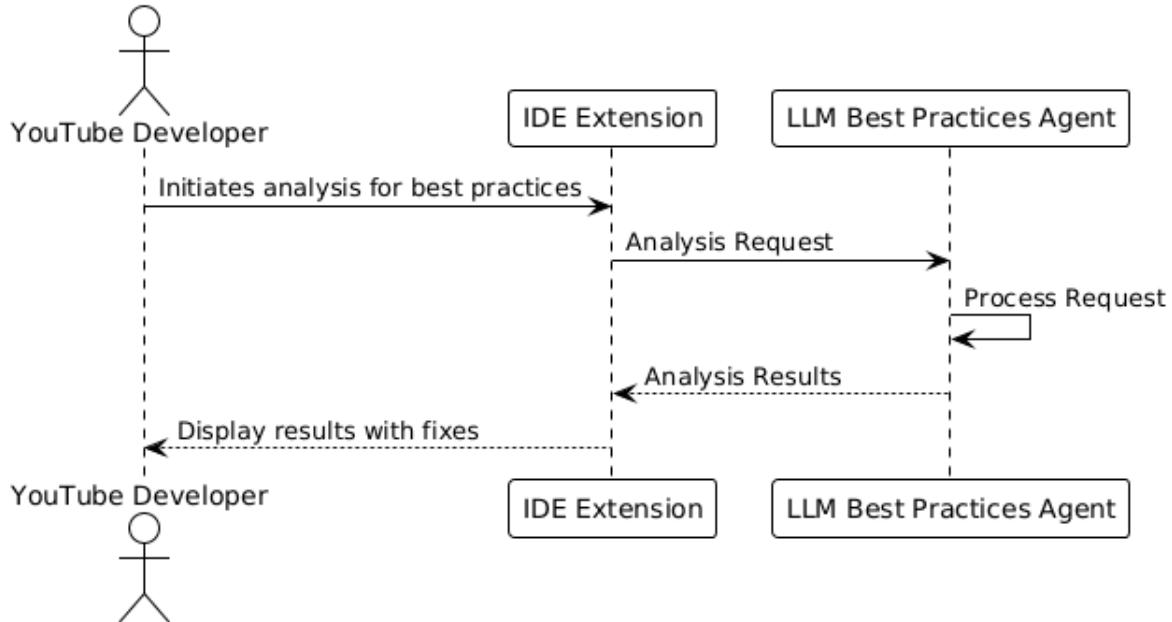


Figure 3.2 – System Interaction Sequence Diagram

The primary workflow includes:

1. **Analysis Initiation:** The developer triggers analysis on the open file.
2. **Request Submission:** The IDE Extension submits the request to the AI agent, identifying the target file.
3. **Processing:** The agent performs best-practices analysis.
4. **Result Delivery:** The agent returns structured findings.
5. **Presentation:** The IDE Extension displays violations and suggested fixes in the editor.

This workflow decouples the IDE from heavy computation, maintaining responsiveness while delegating analysis to the specialized agent framework.

2 Components Design

3.2.1 LLM Best Practices Agent

The LLM Best Practices Agent is the core intelligence engine of the system. Its responsibilities include analyzing code, detecting violations of internal YouTube framework best practices, providing human-readable explanations, and suggesting actionable fixes. This component represents the intersection of AI capabilities with domain-specific software engineering expertise.

3.2.1.1 Agent Architecture Options

We explored multiple architectural paradigms for the agent, drawing inspiration from established AI agent patterns [14]. Figure 3.3 illustrates the fundamental differences in control flow between the two main options.

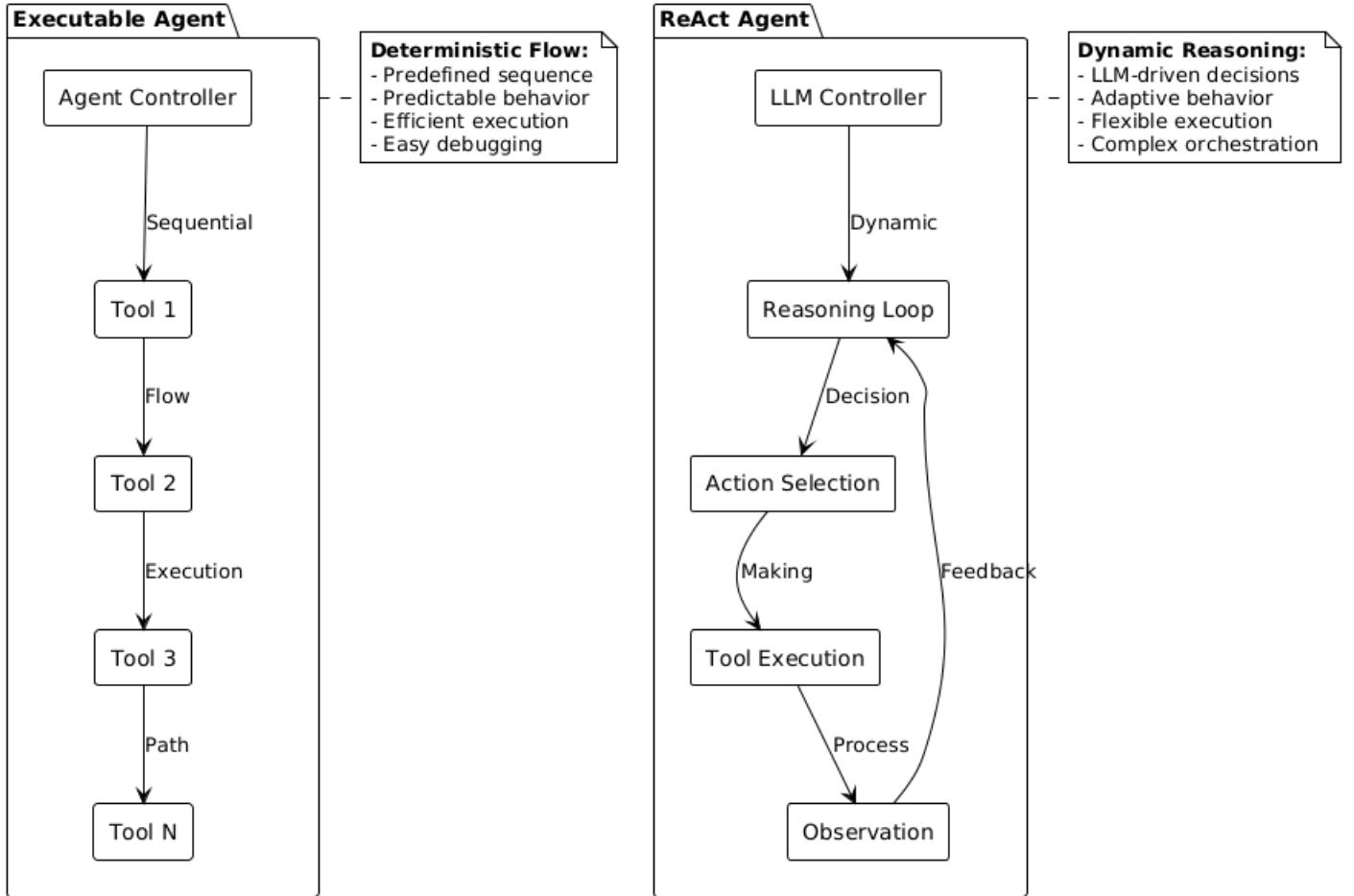


Figure 3.3 – Comparison of Executable Agent vs. ReAct Agent Architectures

- **ReAct Agent:** This architecture integrates *reasoning and acting* loops. The LLM decides the next action, executes it, observes outcomes, and continues iterating. While highly flexible and capable of dynamic decision-making, ReAct agents introduce variability in execution flow, increased token consumption, and difficulty in guaranteeing deterministic outcomes.
- **Executable Agent:** The Executable Agent pattern [14] orchestrates a deterministic, predefined workflow. Control flow is explicitly encoded, with the LLM primarily responsible for domain-specific processing within each tool rather than orchestrating overall

behavior. This approach provides predictable behavior, easier debugging, structured tool integration, and efficient token usage.

Each option has trade-offs:

- *ReAct*: High flexibility, adaptive reasoning, but non-deterministic and potentially higher latency/cost.
- *Executable*: Predictable, maintainable, supports modular tool orchestration, but less flexible for emergent reasoning scenarios.

3.2.1.2 Processing Strategy Options

For handling multiple violations in a single file, we considered three strategies:

- **Fully Sequential**: Processes violations one at a time, ensuring maximal reliability and deterministic outcomes. However, this introduces latency for large files with many violations.
- **Fully Parallel**: Processes all violations concurrently to maximize throughput. While performant, this increases complexity in managing concurrency, resource contention, and error handling.
- **Hybrid Approach**: Limits parallel execution to a controlled number of concurrent violations. This balances throughput and stability, leveraging parallelism without overloading resources or compromising reliability.

3.2.1.3 Core Tools Architecture

To maintain modularity and separation of concerns, the agent workflow is decomposed into specialized tools, each responsible for a specific step of the best practices enforcement pipeline:

- **File Reading Tool**: Retrieves complete file content and metadata. Ensures full context for downstream analysis.
- **Code Analysis Tool**: Performs semantic and structural analysis, detecting violations against internal framework rules. Integrates static analysis heuristics with LLM-based reasoning.
- **Violation Explanation Tool**: Converts raw violations into human-readable explanations, providing actionable insight and contextual rationale for developers.

III.2 Components Design

- **Code Fix Tool:** Suggests concrete corrective actions for each violation. Balances automated recommendations with developer flexibility, allowing acceptance, modification, or rejection.
- **Result Consolidation Tool:** Aggregates outputs from all tools into a structured response consumable by the IDE extension, ensuring clear and contextually accurate presentation.

3.2.1.4 Integration with LLM Infrastructure

The agent interfaces with an internal AI platform hosting multiple LLM models. Key design considerations include:

- **Model-Agnostic Orchestration:** The architecture supports seamless adoption of new models without modifying the agent workflow.
- **Stable Interfaces:** Each tool interacts with the LLM via encapsulated, stable APIs to ensure maintainability.
- **Long-Term Flexibility:** New analysis rules or models can be integrated without architectural changes.

3.2.2 Evaluation Framework for Architecture Decisions

Given the complexity and multiple viable options for agent architecture and processing strategies, empirical evaluation is essential. Decisions cannot be based solely on intuition; performance, reliability, and quality need to be measured under realistic conditions.

3.2.2.1 Test Suite Framework

A set of 12 test cases was designed to cover a representative variety of scenarios, including:

- Files with few versus many violations
- Simple versus complex violation patterns
- Edge cases where violations overlap or conflict
- Different programming constructs and language features

This diversity ensures that evaluations measure robustness across real-world conditions.

3.2.2.2 LLM-as-a-Judge Methodology

To assess output quality, we adopted an *LLM-as-a-Judge* methodology. Instead of relying on brittle string matching, which can fail to recognize semantically correct outputs expressed differently, a separate LLM acts as an impartial evaluator:

- Compares agent outputs against gold-standard solutions.
- Evaluates semantic correctness, completeness, and clarity.
- Scores explanations, suggested fixes, and overall alignment with framework best practices.

This approach enables robust evaluation of both ReAct and Executable workflows while capturing the nuances of human-readable explanations and context-aware fixes.

3.2.2.3 Key Metrics

The evaluation framework measures three primary dimensions:

- **Accuracy:** Semantic correctness of analysis results and suggested fixes, as judged by the LLM evaluator.
- **Latency:** End-to-end response time, ensuring real-time usability in the IDE.
- **Cost:** Token consumption, API usage, and compute resource requirements.

3.2.2.4 Candidate Architectures and Strategies

The evaluation framework allows systematic comparison of all considered options:

- **Sequential Executable Agent**
- **Parallel Executable Agent**
- **ReAct Agent**

Each candidate is assessed using the test suite and LLM-as-a-Judge methodology, producing data-driven insights into trade-offs between determinism, throughput, accuracy, and operational cost. This combination of structured evaluation and modular design ensures that architecture and processing strategies can be selected with confidence, balancing performance, reliability, maintainability, and developer usability.

3.2.3 YouTube IDE Extension

The YouTube IDE Extension serves as the user-facing interface that seamlessly integrates the LLM Best Practices Agent into YouTube developers' daily workflow. Since YouTube developers are the primary target audience for this system, the YouTube IDE Extension was chosen as the natural entry point, leveraging their existing development environment and workflow patterns. This component is designed to provide intelligent, context-aware feedback while maintaining the responsiveness and familiarity that developers expect from their development environment. The feature becomes available when developers enable a user setting in the extension, and entry points only appear for files that belong to the internal YouTube framework for which we enforce best practices.

3.2.3.1 Extension Architecture

The YouTube IDE Extension serves as a lightweight client that orchestrates the interaction between developers and the AI analysis system. The architecture ensures responsiveness by delegating computationally intensive analysis to the specialized agent framework while handling user interface concerns, progress indication, and result presentation locally.

3.2.3.2 User-Triggered vs. Automatic Analysis Design Decision

A fundamental design decision for the IDE extension was whether to implement user-triggered analysis or automatic analysis. This choice significantly impacts user experience, system performance, and resource utilization.

The primary motivation for user-triggered analysis stems from the need to maintain developer productivity and system efficiency. LLM analysis is computationally expensive and resource-intensive, making continuous analysis impractical for maintaining IDE responsiveness. User-triggered analysis ensures that analysis occurs only when developers specifically request it, providing contextually relevant feedback at optimal moments without interrupting their workflow. This approach aligns with developer expectations of having control over their development environment while ensuring that computational resources are used efficiently.

Table 3.1 summarizes the comparison of user-triggered vs. automatic analysis approaches:

III.2 Components Design

Table 3.1 – Comparison of User-Triggered vs. Automatic Analysis Approaches

Criteria	User-Triggered	Automatic
User Control	✓	
Resource Efficiency	✓	
IDE Performance	✓	
Contextual Timing	✓	
Discoverability		✓
Always Current		✓

Based on this analysis, the user-triggered approach was selected as it provides superior resource management, user control, and system performance. The trade-offs in discoverability and stale state management are addressed through intuitive UI design and comprehensive feedback mechanisms.

Stale State Challenge The user-triggered approach introduces a fundamental design challenge: maintaining the relevance and accuracy of analysis results as developers continue modifying their code. This challenge requires balancing system responsiveness with result accuracy, ensuring that feedback remains useful throughout the development process.

3.2.3.3 User Interface Design

The YouTube IDE Extension is designed around two core interaction patterns: entry points for initiating analysis and feedback mechanisms for presenting results.

Entry Points The YouTube IDE Extension provides multiple entry points to ensure accessibility and discoverability for different user preferences and workflows:

- **Visual Interface Integration:** Visual indicators are integrated into the development environment to provide immediate visibility of AI analysis availability, ensuring maximum discoverability while maintaining a clean interface.
- **Command-Based Access:** For developers who prefer keyboard-driven workflows, the extension provides command-based access through standard IDE navigation patterns, supporting both mouse-driven and keyboard-driven user interactions.

Feedback Mechanisms The YouTube IDE Extension employs three core feedback mechanisms designed to integrate seamlessly with existing development workflows:

III.2 Components Design

- **Progress Indication:** Real-time status updates during analysis processing to maintain developer awareness and system transparency.
- **Violation Display:** Presentation of analysis results using familiar interface patterns that leverage developers' existing knowledge of standard feedback mechanisms.
- **Contextual Suggestions:** Interactive code solutions that appear when developers interact with violation markers, providing actionable recommendations.

3.2.3.4 User Interaction Flow

The user interaction flow with the YouTube IDE Extension follows a structured pattern from analysis initiation to optional fix application:

III.2 Components Design

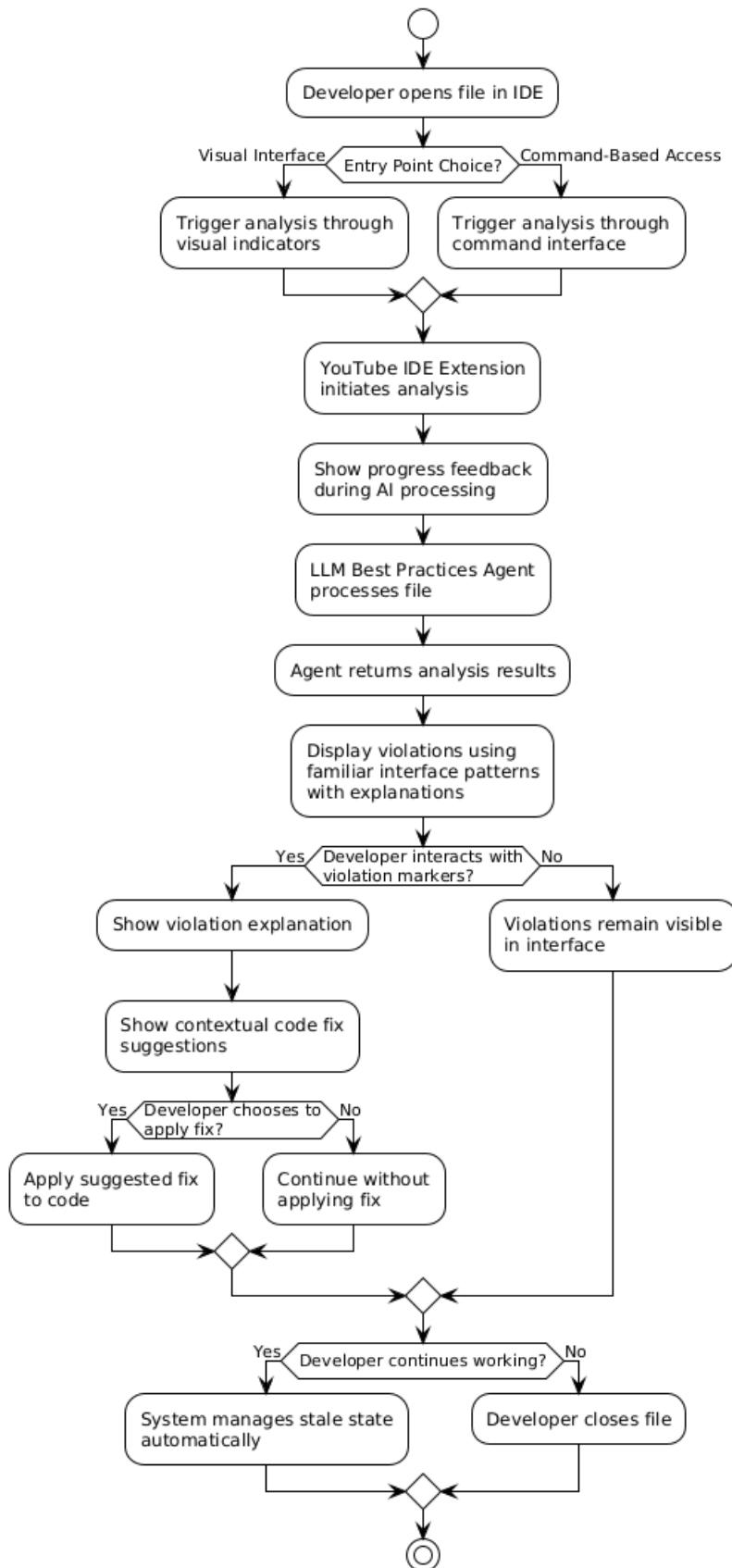


Figure 3.4 – YouTube IDE Extension User Interaction Flow: Complete developer journey from analysis trigger to fix application

The flow demonstrates the core interaction pattern: developers initiate analysis through visual or command interfaces, receive progress feedback during processing, and interact with displayed violations to access explanations and optional fixes. This design ensures developers maintain control over their workflow while providing comprehensive feedback when requested.

3 Data Models and Interfaces

3.3.1 Input/Output Specification

The system's data models define the contracts between components, ensuring consistent communication and data exchange throughout the analysis pipeline. These interfaces establish clear boundaries between the IDE extension and the AI agent, enabling independent evolution of each component.

Analysis Request Format The YouTube IDE Extension sends analysis requests to the LLM Best Practices Agent using a minimal input format that identifies the target file for analysis. This design choice ensures that the agent can focus on its core responsibility of code analysis while maintaining security and proper workspace isolation. The request format includes the file path.

Analysis Response Format The agent returns a structured response containing the analysis results, error information, and optional metadata. The response format includes status information indicating success or failure, violation details with explanations and suggested fixes, and usage statistics for monitoring purposes. This standardized format ensures that the IDE extension can consistently process and display results regardless of the underlying analysis complexity.

3.3.2 Convention Data Management

The system's convention data model defines how YouTube framework best practices are structured, stored, and accessed throughout the analysis pipeline.

Convention Data Structure The convention data model captures best practice definitions as structured objects that support efficient programmatic access and analysis. Each convention definition includes essential metadata such as unique identifiers, descriptions, correct examples, and incorrect examples. This structure enables rapid lookup and context-specific retrieval

III.3 Data Models and Interfaces

during code analysis, with the design optimized for constant-time access patterns required by the agent’s processing pipeline.

Storage Architecture Decision The system employs an in-memory storage approach using structured Python objects rather than external file-based or database storage. This design decision balances several architectural considerations:

- **Data Characteristics:** Conventions are static during runtime, requiring no dynamic updates, making in-memory storage appropriate for performance-critical analysis.
- **Performance Requirements:** In-memory access ensures ultra-low-latency retrieval for real-time analysis, eliminating I/O overhead during agent execution.
- **Simplicity & Reliability:** Structured Python objects provide type safety and eliminate parsing overhead while ensuring data integrity.
- **Resource Efficiency:** Conventions are loaded once at startup, minimizing runtime resource consumption and avoiding repeated file system access.
- **Architectural Flexibility:** The design allows future migration to external storage if multi-framework support or dynamic updates become necessary.

Integration Interfaces The system defines minimal, versioned boundaries that keep components decoupled:

- **IDE Extension Interface:** Contract between the YouTube IDE Extension and the LLM Best Practices Agent that defines the analysis request and response format, enabling independent evolution of UI and agent components.
- **Convention Access Interface:** Defines how the agent accesses convention definitions through in-memory lookup mechanisms, providing a stable boundary for data retrieval without external dependencies.
- **Monitoring Interface:** Captures usage statistics and performance metrics for system observability and evaluation purposes.

Conclusion

The architecture of the LLM Best Practices Agent is defined by three central design decisions. First, the adoption of the Executable Agent paradigm ensures deterministic execution, structured tool orchestration, and predictable performance, avoiding the drawbacks of open-ended reasoning loops. Second, the hybrid parallel–sequential processing strategy balances efficiency with interpretability, allowing analyses to scale while maintaining transparency in intermediate results. Finally, the IDE extension provides a seamless developer experience, integrating feedback, explanations, and fixes directly into familiar workflows while managing state and staleness automatically. Together, these pillars create a robust, cost-efficient, and developer-friendly system for embedding AI-driven best practice enforcement into the coding environment.

Chapter IV

Implementation

Introduction

This chapter presents the evaluation results that informed the final architecture choice, then details how we implemented that choice. We report what the evaluation revealed and, based on those results, what we chose and how it is realized in the system.

The chapter is structured as follows: first, we present the evaluation results and architecture decision; then, we describe the working environment and technology stack; finally, we detail the implementation of both the AI agent backend and the IDE extension frontend.

1 Working Environment

4.1.1 Development Infrastructure

The implementation leverages Google’s internal development infrastructure, providing support for large-scale software development. This infrastructure ensures security, scalability, and integration with existing YouTube development workflows.

4.1.1.1 Internal IDE

The development environment utilizes Google’s internal IDE, which provides a development experience similar to Visual Studio Code but optimized for Google’s internal infrastructure and security requirements.

4.1.1.2 Internal RPC Playground

The RPC Playground is Google’s internal tool for testing and debugging Remote Procedure Call (RPC) services [15]. This tool serves as a playground for sending RPC requests and was essential for developing and testing the communication protocol between the AI Agent and the YouTube IDE Extension.

4.1.1.3 Google Colab

Google Colab was used during early prototyping to iterate on prompt design, tool orchestration, and Executable Agent behaviors before production hardening. Colab provided hosted notebooks with on-demand compute (including GPUs/TPUs) and easy sharing for rapid experiments [16, 17]. It was part of the development environment rather than the deployed technology stack.

4.1.1.4 Internal Repository Integration

All code is stored and versioned within Google’s internal repository system, enabling proper code review processes and collaboration.

4.1.2 Project Management and Documentation

4.1.2.1 Internal Version Control

The project utilizes Google’s internal version control system, which provides Git-like functionality while ensuring compliance with internal security and access control requirements. Git is a fast, scalable, distributed version control system designed to handle everything from small to very large projects with speed and efficiency [18].

4.1.2.2 Internal Code Review Platform

Google’s internal code review platform provides code review capabilities, ensuring code quality and knowledge sharing across development teams.

The code review platform includes:

- **Automated Review Suggestions:** AI-powered suggestions for code improvements, best practices, and potential issues.
- **Collaborative Review Process:** Tools for managing review workflows, assigning reviewers, and tracking review progress.
- **Integration with CI/CD:** Automatic triggering of builds and tests when code changes are submitted for review [19].

4.1.2.3 Internal Project Management System

The project management system provides project tracking, task management, and collaboration capabilities similar to Jira but optimized for Google's internal workflows. JIRA is a flexible issue tracking system that provides project management capabilities [20].

4.1.2.4 Google Docs

Google Docs was used for authoring and reviewing design documents, leveraging the internal built-in Approvals workflow to formalize stakeholder sign-off. The review process combined live comments, suggestions, and targeted approvals to ensure traceable decisions.

These project workflows ensured fast iteration, early detection of issues, and compliance with Google's security and code quality standards.

2 Technologies

This section presents the concrete technologies used to implement the system. We distinguish between industry-standard tools (e.g., Python, TypeScript, JSON) and internal platforms operated within Google (e.g., YouTube DevInfra Agent Framework, internal AI platform).

4.2.1 Backend Technologies (AI Agent)

4.2.1.1 Python Programming Language

Python serves as the primary programming language for the AI agent framework implementation. Python is a high-level, interpreted programming language known for its simplicity, readability, and library ecosystem [21]. The language's dynamic typing and support for artificial intelligence and machine learning libraries make it suitable for AI agent development [22].

Python's advantages for this implementation include:

- **AI/ML Ecosystem:** Extensive libraries for machine learning, natural language processing, and AI development.
- **Asynchronous Programming:** Built-in support for asynchronous programming patterns essential for handling concurrent requests.
- **JSON Processing:** Native support for JSON serialization and deserialization required for API communication.

4.2.1.2 YouTube DevInfra Agent Framework

The implementation utilizes the YouTube DevInfra Agent Framework — the serving infrastructure designed by YouTube for YouTube Developer Infrastructure agents. It underpins all YouTube agents, providing standardized execution, deployment, and operational primitives for LLM-powered applications.

This framework was selected because it natively supports the *Executable Agent* pattern and supplies built-in tool orchestration, workflow control, and production lifecycle management. Using it aligns the implementation with DevInfra standards and avoids rebuilding common agent infrastructure, allowing focus on best-practices enforcement logic.

4.2.1.3 Internal AI Platform

The system integrates with Google’s internal AI platform, which hosts internal LLM models trained on Google’s codebase, ensuring organization-specific knowledge and compliance with internal security requirements.

4.2.1.4 LLM Libraries and Frameworks

The internal agent framework utilizes several specialized libraries and frameworks for LLM interaction and agent development:

- **LLM Interaction Libraries:** Specialized libraries for communicating with internal LLM models, monitoring token usage, and optimizing API calls.
- **Agent Orchestration Libraries:** Libraries that provide the Executable Agent pattern implementation, tool registration, and workflow management.
- **Prompt Engineering Libraries:** Frameworks for constructing, optimizing, and managing prompts.

4.2.2 Frontend Technologies (IDE Extension)

4.2.2.1 TypeScript Programming Language

TypeScript is used for the frontend development of the YouTube IDE Extension. TypeScript is a strongly typed superset of JavaScript that compiles to plain JavaScript [23]. The language provides static type checking, which helps prevent runtime errors and improves code maintainability in large-scale applications.

TypeScript was chosen because we are building a feature inside an existing extension implemented in TypeScript, ensuring direct compatibility and reuse. Its static typing and interfaces improve maintainability and reduce runtime errors in complex UI state and service interactions. It also integrates seamlessly with VS Code API and other development environments.

4.2.2.2 VS Code Extension API

The system integrates with Visual Studio Code through its Extension API. Visual Studio Code is a source-code editor developed by Microsoft, built on the Electron framework [24]. The VS Code Extension API provides capabilities for extending the editor's functionality.

Given that the internal IDE is VS Code-like, adopting the VS Code Extension API is the natural choice: it is natively supported within the environment (requiring no additional infrastructure), exposes the command, user-interface, and configuration interfaces required by the best-practices enforcement feature, and ensures compatibility with the existing extension ecosystem. In practice, the API provides the integration points necessary to implement the designed interaction flow without introducing custom runtime scaffolding.

3 Realization

The realization phase translates architectural decisions into a concrete implementation of the system. It consists of two primary components: (1) the **AI Agent backend**, responsible for code analysis and violation handling, and (2) the **IDE Extension frontend**, responsible for developer interaction. This section first discusses the evaluation results that informed our architectural choice, then presents the detailed implementation of the agent and its integration with the IDE extension.

4.3.1 Agent Realization

4.3.1.1 Evaluation Results and Architecture Decision

To determine the most suitable agent design, we evaluated three candidate architectures: the Parallel Executable, Sequential Executable, and ReAct Agent. These were assessed against 12 representative test cases covering a range of violation patterns and code complexities.

For a simple single-violation file (TC8_PROPS_MISSING_EXPORT), the results are shown in Table 4.1. Parallel outperformed Sequential by $\sim 33\%$ in latency, while both were vastly more efficient than ReAct in token usage.

IV.3 Realization

Table 4.1 – Single-Violation File Performance Comparison

Architecture	Latency	Input Tokens	Output Tokens
Parallel Executable	4.1s	5,157	331
Sequential Executable	6.1s	5,157	331
ReAct Agent	15.3s	37,989	1,523

Scaling to the full evaluation suite (Figures 4.1–4.3) revealed deeper trade-offs. Parallel maintained faster throughput but suffered reliability issues such as `OverLimitException` and `DEADLINE_EXCEEDED`. Sequential was slower but fully reliable with identical accuracy (~91.7%). ReAct proved infeasible due to extreme token consumption and lower accuracy.



Figure 4.1 – Latency Performance Across 12 Test Cases

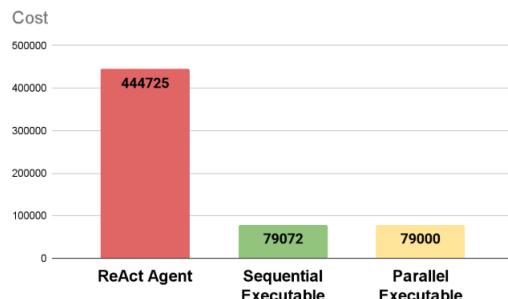


Figure 4.2 – Cost Analysis Across 12 Test Cases

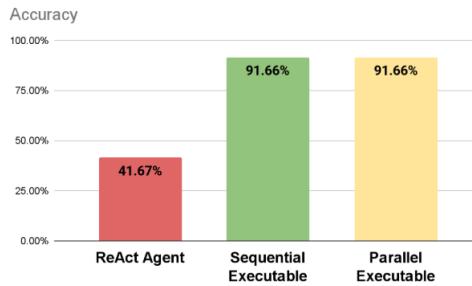


Figure 4.3 – Accuracy Performance Across 12 Test Cases

From these results, the final decision was to adopt a **Parallel Executable architecture with concurrency limiting**. This preserves most of the performance advantage of Parallel while ensuring reliability through semaphore-based concurrency control, preventing resource saturation. This hybrid approach combines determinism, throughput, and robustness, making it suitable for production use in developer IDEs.

4.3.1.2 Agent Implementation Overview

The AI Agent implements the chosen architecture using Python and integrates with Google’s internal AI platform hosting Gemini-based LLMs. The design is class-based and modular, ensuring extensibility, observability, and testability. The agent exposes a single `execute()` entry point, orchestrating a pipeline of specialized tools that handle file reading, analysis, explanation generation, fix generation, and result consolidation.

Configuration and Initialization At startup, the agent performs several critical setup tasks:

- **Model Selection:** Configured to use the latest Gemini-based model trained on internal Google code.
- **Convention Loading:** Best practices are retrieved from repository-stored JSON files and cached in memory.
- **Tool Registration:** The five specialized tools are instantiated and registered into a deterministic workflow.

Implementation Structure The implementation follows a layered orchestration pattern:

- A central agent class orchestrates the pipeline via dependency-injected tools.

- Each tool is encapsulated as a class with clear `run()` contracts and typed inputs/outputs.
- Metrics and token usage are logged per tool, enabling fine-grained observability.
- Separation of concerns allows tools to be replaced or extended without modifying orchestration logic.

4.3.1.3 Core Tools Implementation

The agent implements five specialized tools, each corresponding to one stage of the analysis pipeline. Tools are designed as independent classes that expose a public `run()` method, which enforces a clear input/output contract and raises typed errors when failures occur. This contract-based design makes the system modular, testable, and resilient to partial failures.

ReadFileFromWorkspace Tool Contract: input = file path; output = file content (string). The ReadFileFromWorkspace tool is responsible for retrieving the contents of the developer's source file. Implementation details include:

- **File system access:** Uses Python's built-in file handling with UTF-8 as the default encoding, while detecting and recovering from alternative encodings.
- **Error handling:** Raises typed errors such as `FileNotFoundException`, `PermissionDeniedError`, and `EncodingException`. These errors are logged and surfaced to the developer in structured form.
- **Robustness:** Handles large files by streaming content, ensuring memory efficiency.

CodeAnalysisTool Contract: input = file content; output = list of base violations. This tool forms the analysis core of the system. It detects framework violations by orchestrating LLM queries enriched with contextual information.

- **Prompt construction:** Implements a template system with slots for file type, code snippet, and dynamically filtered convention definitions.
- **Context injection:** Selects relevant conventions from the cache and injects them into the prompt to guide the model.
- **Response parsing:** Uses strict JSON schema validation with recovery mechanisms for malformed LLM responses.
- **Performance optimizations:** Employs caching of templates and convention definitions, reducing repeated token usage and lowering latency.

ViolationExplanationTool Contract: input = base violation; output = natural-language explanation. The ViolationExplanationTool generates educational explanations that help developers understand not only what the violation is, but why it matters.

- **Contextualization:** Pulls in violation metadata (rule violated, line number, code snippet) to ground explanations in concrete evidence.
- **Generation style:** Prompts the LLM to balance technical accuracy with readability, avoiding overly generic statements.
- **Implementation detail:** Each explanation is post-processed for clarity, removing redundant phrasing and enforcing concise output.

CodeFixTool Contract: input = violation + explanation; output = code fix (annotated snippet). This tool provides actionable, safe, and educational fixes.

- **Safety:** Fixes are constrained to local code changes, avoiding edits that break APIs or dependencies.
- **Self-containment:** Ensures that each fix can be applied without requiring modifications in other files.
- **Contextual adaptation:** Uses file content and violation metadata to adapt fixes to the surrounding code style.
- **Educational emphasis:** Each fix includes explanatory comments describing why the change is required.
- **Error handling:** Raises a `FixGenerationError` if the LLM output cannot be parsed or validated as compilable code.

Finish Tool Contract: input = violations + explanations + fixes; output = structured response for IDE extension. The Finish Tool acts as the consolidation component, producing a structured JSON-like response that the IDE extension can directly render.

- **Aggregation:** Combines violations, explanations, and fixes into a unified structure keyed by violation ID.
- **Deduplication:** Identifies overlapping or redundant violations and merges them to reduce noise.

- **Validation:** Ensures schema compliance so that the IDE can reliably parse and render results.
- **Formatting:** Applies consistent formatting (line numbers, code blocks, explanations) to improve readability in the UI.

Collectively, these tools form a deterministic pipeline coordinated by the agent's `execute()` method. Each tool adheres to explicit contracts, which improves modularity, observability (per-tool token usage is logged), and long-term maintainability.

4.3.1.4 Processing Workflow

The hybrid strategy is realized in three stages:

1. **Sequential Preprocessing:** Full file content is read and all violations are identified.
2. **Parallel Processing:** Explanation and fix generation tasks are executed concurrently with semaphore-based limits on in-flight LLM calls.
3. **Result Consolidation:** Outputs are aggregated, validated, and prepared for the IDE extension.

Figure 4.4 illustrates this workflow.

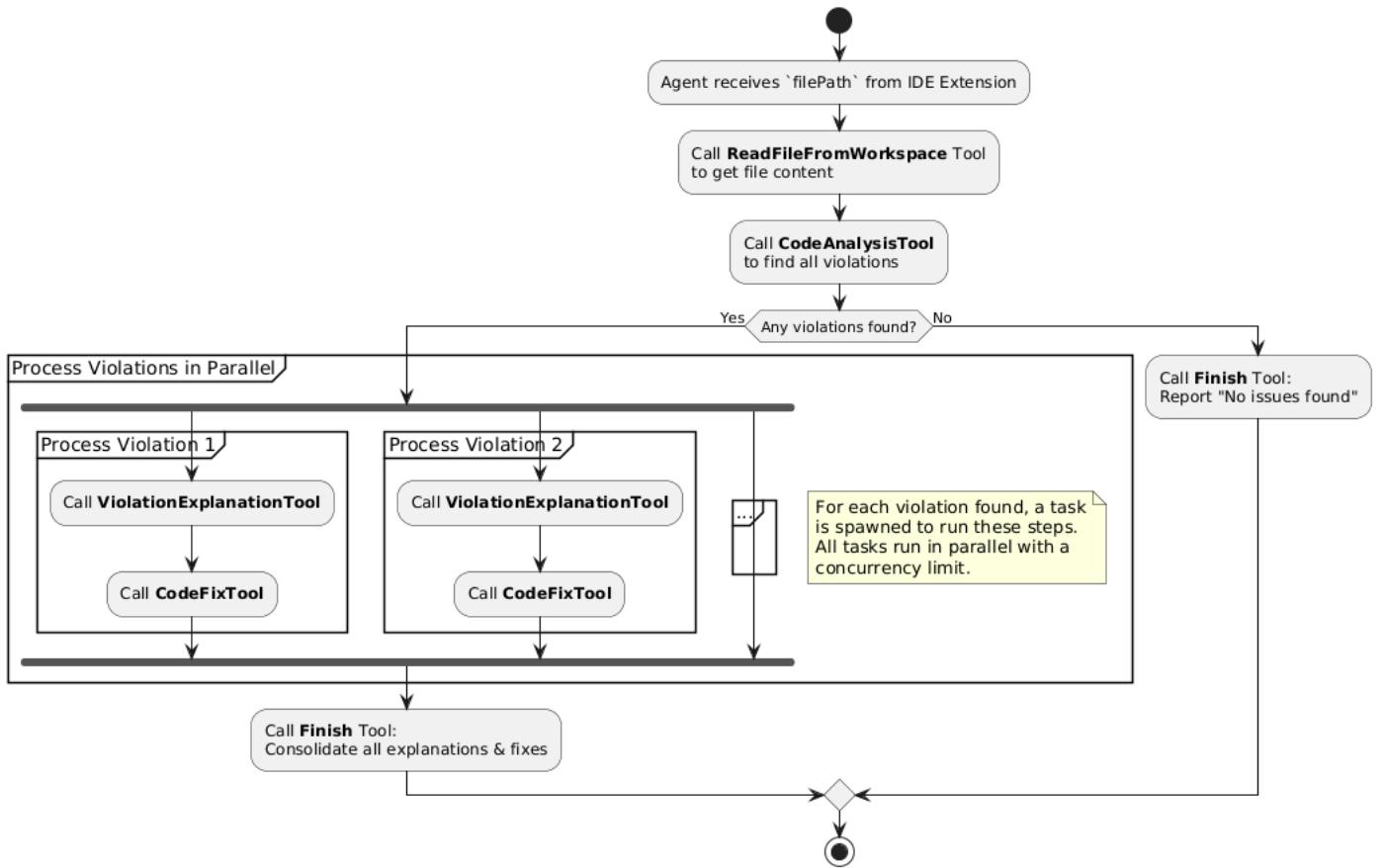


Figure 4.4 – Agent Processing Activity Diagram

4.3.1.5 Resilience and Error Handling

The system ensures graceful degradation under failures:

- Independent processing of violations isolates failures to individual tasks.
- Retry with exponential backoff mitigates transient network or LLM errors.
- Typed error classes facilitate debugging and developer support.
- Partial results are preserved, ensuring developers always receive usable feedback.

4.3.1.6 Convention Data Management

The convention data management system implements loading, caching, and retrieval mechanisms for YouTube framework best practices. The conventions are stored as Python objects in

an array, each containing a unique identifier, description, correct example, and incorrect example. For efficient runtime access, the system constructs an in-memory map keyed by convention ID, allowing the agent tools to retrieve only the relevant convention on demand.

Loading and Initialization At startup, all convention objects are loaded into memory from the Python array. A dictionary (map) is created with convention IDs as keys and convention objects as values, providing constant-time access for subsequent tool invocations.

Memory Caching and Access This in-memory caching strategy ensures low-latency access during code analysis:

- **Efficient Lookup:** Tools retrieve conventions by ID from the map, avoiding iteration over the full array.
- **Dynamic Selection:** Only conventions relevant to the current file type and analysis context are queried, minimizing unnecessary data processing.
- **Lightweight and Fast:** The cache resides entirely in memory, requiring no external services, and supports rapid retrieval during concurrent tool executions.

4.3.2 Extension Integration

4.3.2.1 Extension Architecture

The IDE Extension implements a layered architecture that integrates with the overall system through three distinct layers: the Extension layer containing user-facing components, a Proxy layer for authentication and request routing, and the Backend layer hosting AI agent services.

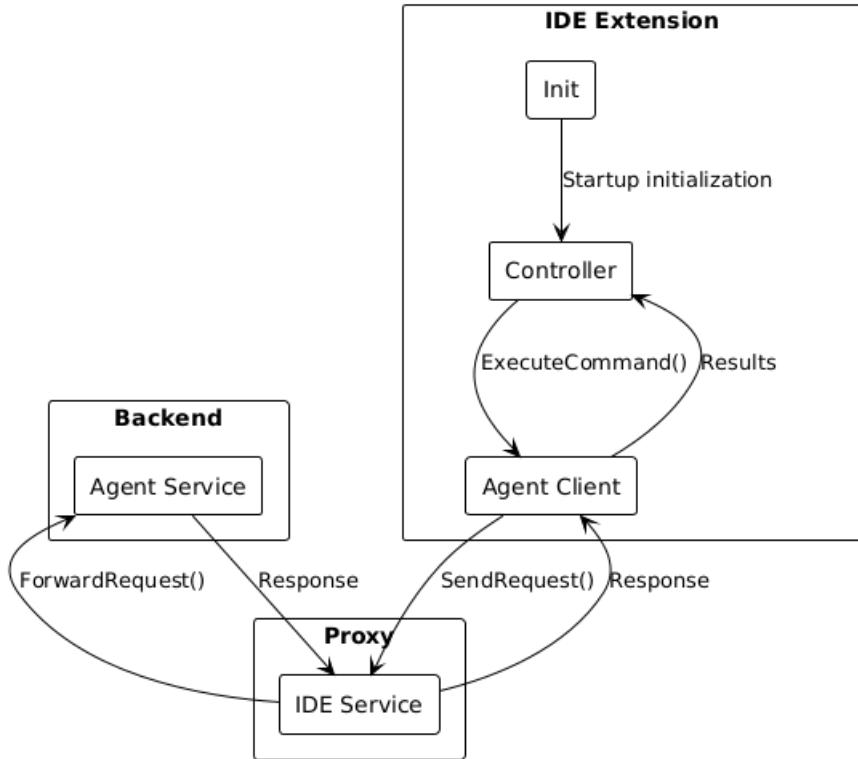


Figure 4.5 – System Architecture: IDE Extension, Proxy, and Backend Communication Flow

The architecture follows a clear request-response pattern where user interactions trigger analysis requests that flow through the IDE Service proxy for authentication and authorization, then to the Agent Service in the backend for processing. Responses follow the same path in reverse, ensuring secure and authenticated communication throughout the entire pipeline while maintaining clear separation of responsibilities between layers.

Extension Components The extension's internal architecture consists of three core components that work together to provide seamless integration with the development environment, as illustrated in Figure 4.5.

Init Component: Handles extension initialization, reading user settings, registering commands and editor actions, and performing health checks. The component ensures proper setup of all dependencies.

Controller: Centralizes all UI-related state and orchestrates interactions between components. The Controller processes commands, manages notifications, routes analysis results, renders diagnostics and hover-based suggestions, and coordinates stale-state transitions to ensure consistent behavior across all entry points.

AgentClient: Manages communication with the backend services through the proxy layer.

The component handles request formatting, implements retry logic with exponential backoff, manages timeouts, and processes responses from the AI agent.

4.3.2.2 User Interaction

The feature is controlled through a dedicated **user setting**. This setting appears as a simple checkbox: when enabled, the feature becomes available in the IDE; when disabled, it is entirely hidden from the interface. This ensures that developers can opt in seamlessly without cluttering the environment for those who do not use the feature.

The primary entry point for triggering analysis is an **Editor Action** integrated into the file title bar (Figure 4.6). This placement ensures high visibility and aligns naturally with the developer's workflow when working on individual files.

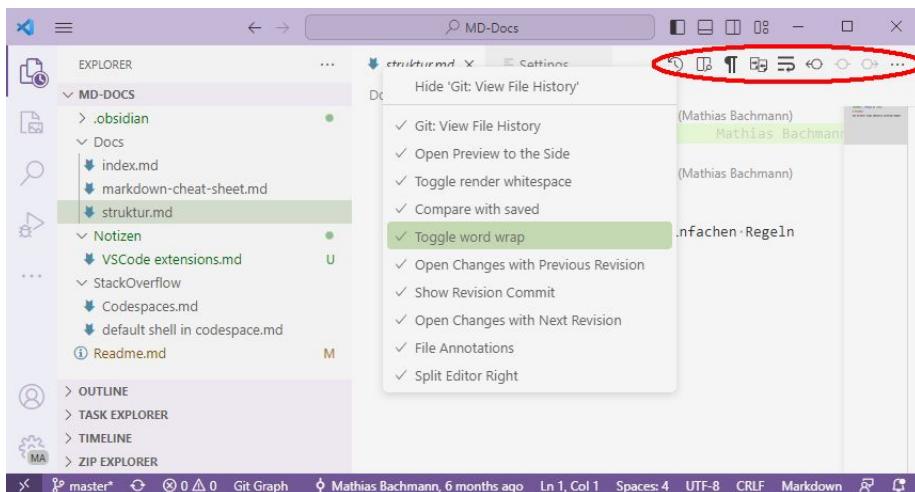


Figure 4.6 – VS Code Interface: Editor Actions (Illustrative).

An additional entry point is provided through the **Command Palette**, which can be invoked using **Ctrl+Shift+P** (or **Cmd+Shift+P** on macOS). This pathway makes the feature equally accessible to developers who prefer keyboard-driven workflows and ensures discoverability for new users exploring available commands (Figure 4.7).

IV.3 Realization

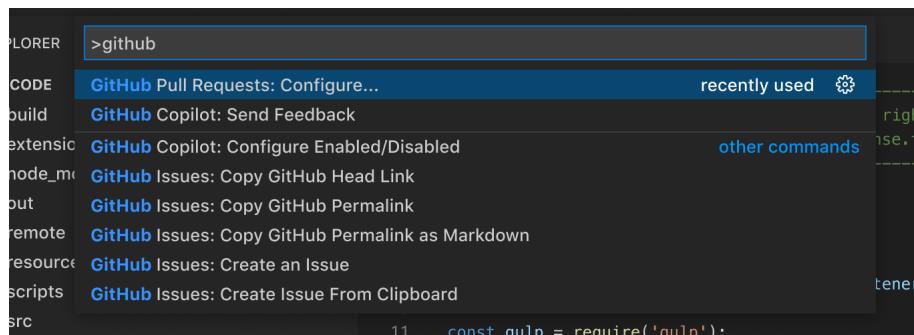


Figure 4.7 – VS Code Command Palette (Illustrative).

Once analysis is triggered, the extension provides immediate feedback via **VS Code notifications** (Figure 4.8). These notifications confirm that a request has been received, update progress, and display clear error messages if issues occur. This ensures transparent communication throughout the request lifecycle.

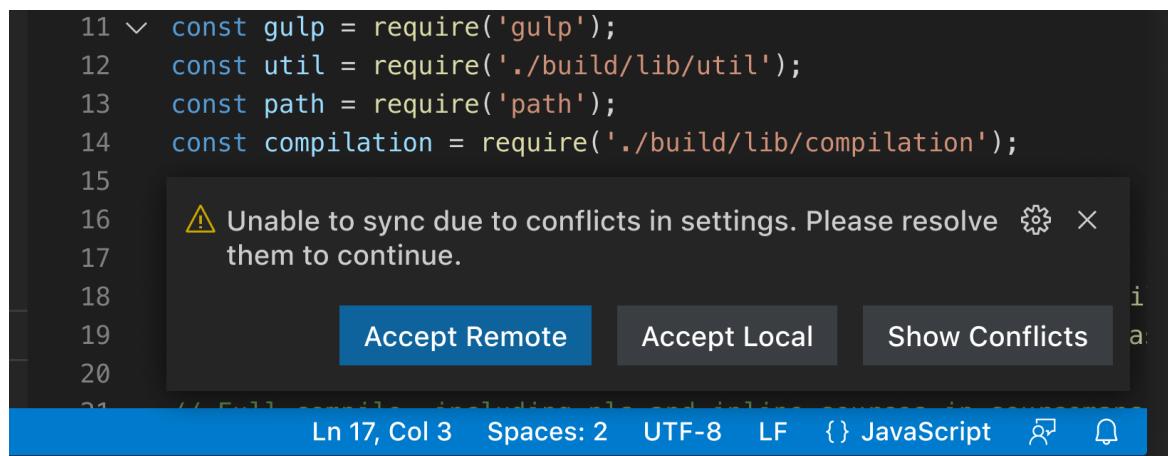


Figure 4.8 – VS Code Notification Interface (Illustrative).

The results of the analysis are surfaced through VS Code's native **diagnostic system**. Violations appear in the **Problems panel** and are underlined directly in the editor, marking the exact range of code that violates a best practice (Figure 4.9). Hovering over the highlighted code reveals the diagnostic explanation, helping developers quickly understand the issue in context.

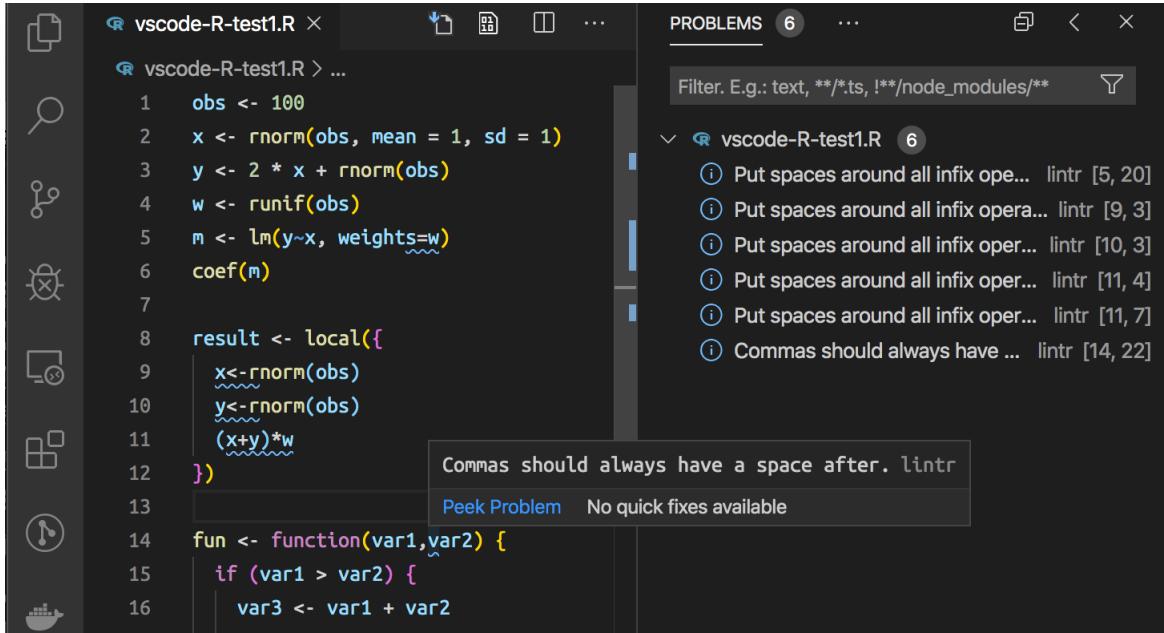


Figure 4.9 – VS Code Diagnostics Interface (Illustrative).

For more detailed guidance, an integrated **hover provider** presents formatted fix suggestions directly within the editor. This approach allows fixes to be displayed with proper syntax highlighting and inline code snippets, offering a clear and actionable path to resolution without leaving the development workflow.

4.3.2.3 Stale Diagnostics Handling

One of the most challenging aspects of IDE integration is maintaining diagnostic accuracy as developers continuously modify their code. The extension implements a sophisticated two-tiered system that balances immediate responsiveness with accurate analysis results.

The Challenge Traditional diagnostic systems struggle with code that changes rapidly, often displaying outdated information that confuses developers and reduces trust in the tool. The challenge is to provide immediate visual feedback while ensuring that diagnostics remain accurate and relevant to the current code state.

Two-Tiered Solution The extension handles stale diagnostics using a two-tiered approach that separates immediate responsiveness from precise re-anchoring:

Tier 1 - Instant Adjustment: Provides immediate feedback on every keystroke. Diagnostics are shifted based on simple text edits and marked as [Outdated] if the flagged code itself is edited. This ensures high **responsiveness** without impacting performance.

IV.3 Realization

Tier 2 - Debounced Re-anchoring: Activates after a 1-second pause in typing to improve diagnostic **accuracy**. The process involves:

1. **Fingerprint:** Creates a contextual hash from the code and its surrounding context to identify exact matches.
2. **Scan with Regex:** Finds all possible text matches across the document.
3. **Re-anchor:** Moves the diagnostic to the correct new location based on the highest match score.

This two-tiered strategy balances responsiveness with accuracy, ensuring developers receive timely feedback while maintaining the integrity of diagnostics even during active code editing.

IV.3 Realization

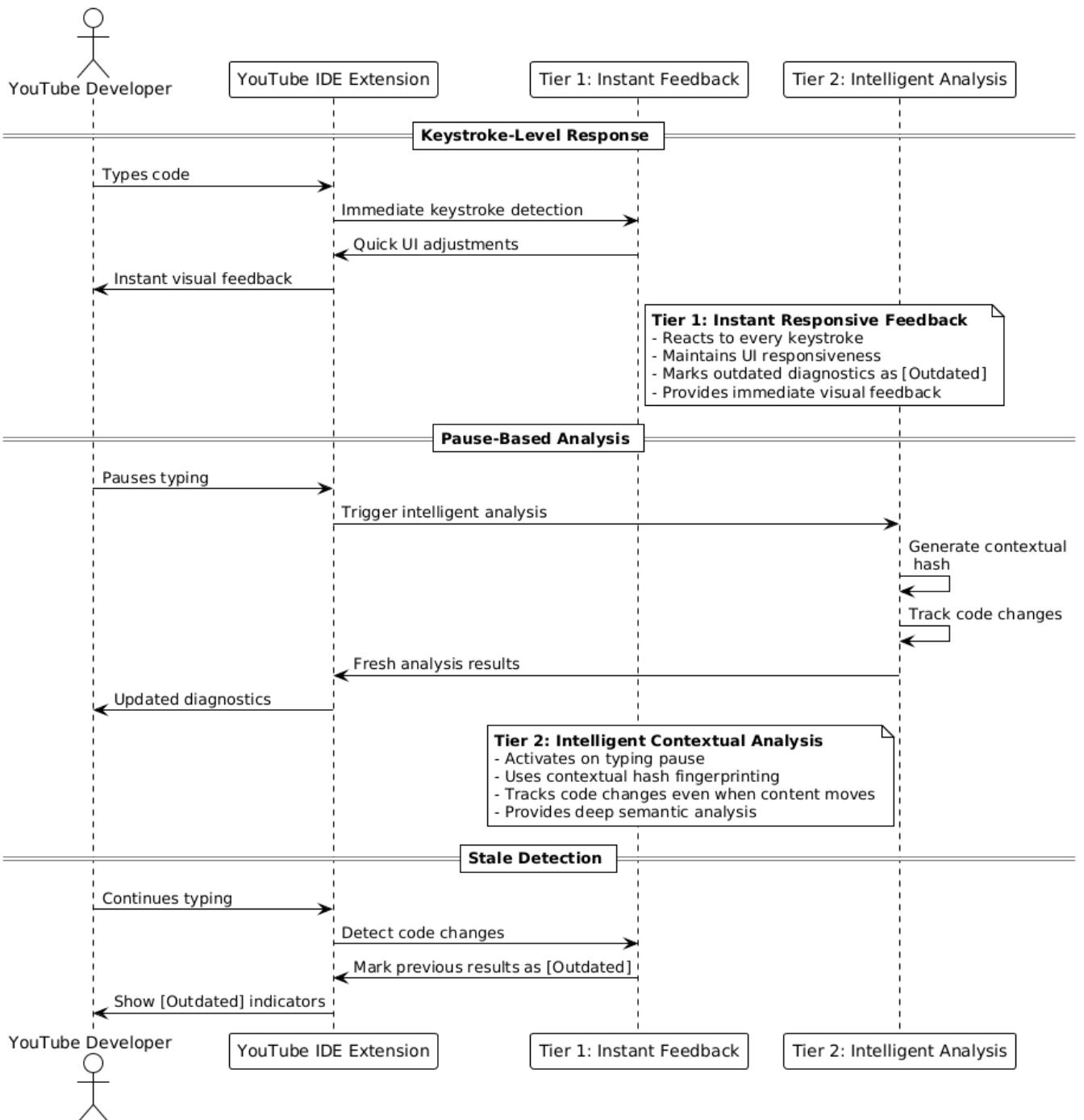


Figure 4.10 – Stale Diagnostics Handling: Two-Tiered System (Illustrative)

This approach ensures that developers receive immediate visual feedback while maintaining diagnostic accuracy through intelligent analysis timing and state management, representing a significant technical contribution to IDE integration challenges.

4.3.2.4 Resilience and Communication

The extension implements comprehensive error handling and resilience mechanisms to ensure reliable operation in production environments. The communication system uses internal RPC infrastructure with JSON payloads for debugging and cross-language compatibility.

Error handling follows fault-tolerant design principles with multi-level error isolation, ensuring that failures in one component do not cascade to others. The system implements specific error types for different failure scenarios, including network communication errors, service unavailability, and timeout errors, with detailed error information and suggested recovery actions.

Retry mechanisms with exponential backoff and bounded attempts ensure operation under transient failure conditions, while timeout management prevents indefinite waiting periods and maintains responsive user experience. The implementation includes request validation and response parsing to prevent communication errors and ensure data integrity throughout the analysis pipeline.

Conclusion

This chapter presented the evaluation results, the resulting architecture decision (Parallel Executable with concurrency limiting), and the complete implementation: environment, technologies, and realization across backend and IDE. The system is production-oriented, balancing speed with stability and cost.

Conclusion and Perspectives

This project has developed and implemented an intelligent assistance system integrated into YouTube’s internal development environment, aimed at enforcing framework-specific best practices in real time. The system addresses a key challenge in large-scale software engineering: the absence of immediate, contextual feedback on internal frameworks. Unlike existing tools that focus on syntax or public libraries, our solution delivers intelligent, context-aware guidance directly within the developer workflow, reducing technical debt and improving code consistency.

The main contribution consists of an AI agent based on Large Language Models (LLMs) that orchestrates five specialized tools for code analysis, explanation, and fix generation, coupled with a YouTube IDE extension offering intuitive entry points and a two-tier diagnostic state management system that balances responsiveness with precision. Development leveraged Python for the agent, TypeScript and the VS Code API for the extension, and Google’s internal AI platform for LLM integration, following an agile Kanban-based workflow.

Design decisions throughout the project were guided by empirical testing and data-driven insights, ensuring that chosen approaches balanced performance, reliability, and usability. Significant technical challenges were addressed, including stale diagnostic handling through a two-tier system, concurrency control for stable performance, and robust semantic evaluation using an LLM-as-a-Judge methodology.

Looking ahead, the next step is implementing a Tiered Analysis Approach that combines a lightweight, rule-based linter for fast detection of simple issues with the LLM agent reserved for complex, subjective cases. This hybrid strategy promises faster, cheaper, and more effective feedback, further enhancing the developer experience. A teammate is already building the linter, and integration plans are in place.

Overall, this project demonstrates the feasibility and value of embedding AI agents into development environments for enforcing framework-specific best practices. It provides a strong foundation for future improvements and represents a meaningful contribution to advancing intelligent developer tools at YouTube.

References

- [1] GOOGLE. Google company overview. <https://about.google/>, (2024). Accessed: 2024-12-19.
- [2] YOUTUBE. Youtube platform statistics. <https://www.youtube.com/about/>, (2024). Accessed: 2024-12-19.
- [3] KENT BECK, MIKE BEEDLE, ARIE VAN BENNEKUM, ALISTAIR COCKBURN, WARD CUNNINGHAM, MARTIN FOWLER, JAMES GRENNING, JIM HIGHSMITH, ANDREW HUNT, BRIAN MARICK, ET AL. *Manifesto for agile software development*. Agile Alliance (2001).
- [4] ROBERT C MARTIN. *Agile software development: principles, patterns, and practices*. Prentice Hall PTR (2003).
- [5] DAVID J ANDERSON. *Kanban: Successful evolutionary change for your technology business*. Blue Hole Press (2010).
- [6] HENRIK KNIBERG AND MATTIAS SKARIN. *Kanban and Scrum-making the most of both*. Lulu. com (2011).
- [7] STUART RUSSELL AND PETER NORVIG. *Artificial Intelligence: A Modern Approach*. Pearson Education, 4th edition (2024).
- [8] CHRISTOPHER M BISHOP. *Pattern Recognition and Machine Learning*. Springer, 1st edition (2006).
- [9] IAN GOODFELLOW, YOSHUA BENGIO, AND AARON COURVILLE. *Deep Learning*. MIT Press (2016).
- [10] RISHI BOMMASANI, DREW A HUDSON, EHSAN ADELI, RUSS ALTMAN, SIMRAN ARORA, SYDNEY VON ARX, MICHAEL S BERNSTEIN, JEANNETTE BOHG, ANTOINE BOSSELUT, EMMA BRUNSKILL, ET AL. *Generative artificial intelligence: An overview*. arXiv preprint arXiv:2301.04226 (2023).
- [11] TOM BROWN, BENJAMIN MANN, NICK RYDER, MELANIE SUBBIAH, JARED D KAPLAN, PRAFULLA DHARIWAL, ARVIND NEELAKANTAN, PRANAV SHYAM, GIRISH SASTRY, AMANDA ASKELL, ET AL. *Large language models: A breakthrough in natural language processing*. Advances in Neural Information Processing Systems **33**, 1877–1901 (2020).
- [12] GOOGLE. Google's ai code generation statistics. Google I/O 2024 Keynote, (2024). Accessed: 2024-12-19.

REFERENCES

- [13] GOOGLE. Google developer survey: Ai in software development. Google Developer Blog, (2024). Accessed: 2024-12-19.
- [14] MICROSOFT AZURE ARCHITECTURE CENTER. Ai agent design patterns. <https://learn.microsoft.com/en-us/azure/architecture/ai-ml/guide/ai-agent-design-patterns>, (2024). Accessed: 2025-09-07.
- [15] SUN MICROSYSTEMS. Rpc: Remote procedure call protocol specification, (1984).
- [16] GOOGLE RESEARCH. Google colab. <https://colab.research.google.com/>, (2017). Accessed: 2024-12-19.
- [17] THOMAS KLUYVER, BENJAMIN RAGAN-KELLEY, FERNANDO PÉREZ, BRIAN GRANGER, MATTHIAS BUSSONNIER, JONATHAN FREDERIC, KYLE KELLEY, JESSICA HAMRICK, JASON GROUT, SYLVAIN CORLAY, ET AL. *Jupyter notebooks—a publishing format for reproducible computational workflows*. Positioning and Power in Academic Publishing: Players, Agents and Agendas pages 87–90 (2014).
- [18] LINUS TORVALDS AND JUNIO HAMANO. *Git: A fast, scalable, distributed version control system*. Proceedings of the Linux Symposium **1**, 3–10 (2005).
- [19] MOJTABA SHAHIN, MUHAMMAD ALI BABAR, AND LIMING ZHU. *Continuous integration and continuous deployment: A systematic review*. ACM Computing Surveys **49**(1), 1–45 (2016).
- [20] ATLASSIAN. *Jira: A flexible issue tracking system*. Atlassian Documentation (2002).
- [21] GUIDO VAN ROSSUM AND FRED L DRAKE JR. *Python reference manual*. Amsterdam: CWI (Centre for Mathematics and Computer Science) (1995).
- [22] FABIAN PEDREGOSA, GAËL VAROQUAUX, ALEXANDRE GRAMFORT, VINCENT MICHEL, BERTRAND THIRION, OLIVIER GRISSEL, MATHIEU BLONDEL, PETER PRETTENHOFER, RON WEISS, VINCENT DUBOURG, ET AL. *Scikit-learn: Machine learning in python*. Journal of machine learning research **12**(Oct), 2825–2830 (2011).
- [23] GAVIN BIERNAN, MARTIN ABADI, AND MADS TORGERSEN. Understanding typescript. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications*, pages 257–268 (2014).
- [24] FERNANDO CASTOR, EDUARDO FIGUEIREDO, NÉLIO CACHO, BRENO SENA, LEONARDO TEIXEIRA, GUSTAVO PINTO, AND BRENO FONSECA. Visual studio code: A new way of developing web applications. In *Proceedings of the 2016 ACM International Conference on Object Oriented Programming Systems Languages & Applications*, pages 1–2 (2016).

Appendix : Miscellaneous remarks

فرض أفضل الممارسات مع تكامل نموذج اللغة الكبير وبيئة التطوير المتكاملة

أنجز هذا المشروع في شركة جوجل زيورخ ضمن متطلبات الدبلوم الوطني للهندسة في هندسة البرمجيات. يدرس دمج الذكاء الاصطناعي التوليدي في تطوير البرمجيات عبر وكيل يعتمد على نموذج اللغة الكبير (LLM) داخل بيئة التطوير المتكاملة (IDE).

يحلّ الوكيل الشيفرة لاكتشاف انتهاكات معقدة قد تغفلها أدوات التحليل الثابت، وينتج شروط واضحة واقتراحات عملية تساعد مطوري يوتوب على رفع جودة الشيفرة والالتزام بالمعايير الفضلى. وباندماجه في سير العمل عن طريق إضافة للبيئة، يعزز الإنتاجية ويقلل الديون التقنية دون تعطيل تجربة التطوير.

الكلمات المفتاحية: هندسة برمجيات، ذكاء اصطناعي توليدي، تكامل IDE، جودة الشيفرة، إنتاجية المطورين

Application des meilleures pratiques avec l'intégration LLM-IDE

Ce projet, réalisé chez Google Zurich, intègre un agent basé sur un LLM au sein d'un IDE interne pour soutenir le développement logiciel.

L'agent analyse le code, fournit des explications claires et des suggestions actionnables, et s'intègre au flux de travail par une extension IDE. Il aide les contributeurs YouTube à maintenir une haute qualité de code et à réduire la dette technique sans perturber l'activité.

Mots-clés : Génie logiciel, IA générative, Intégration IDE, Qualité du code, Productivité

Enforcing Best Practices with LLM-IDE Integration

This project, conducted at Google Zurich, embeds an LLM-powered agent inside an internal IDE to support software development.

The agent analyzes code to detect complex violations missed by traditional static tools, producing clear explanations and actionable fixes. Integrated via an IDE extension, it boosts productivity and reduces technical debt without disrupting the developer workflow.

Keywords: Software Engineering, Generative AI, IDE Integration, Code Quality, Productivity

