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 through LLM Integrations in Google´s internal IDE

---

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 through LLM Integrations in Google´s internal IDE

---

Presented by

## Arij KOUKI

Hosted by



| Company Supervisor | University Supervisor |
|---|---|
|  |  |

Academic Year : 2024/2025

# Acknowledgements

# Table of Contents

# List of Figures

# List of Tables

# 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 : Intelligence Artificielle Générative, Génie Logiciel, 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: Generative AI, Software Engineering, IDE Integration, Code Quality, Developer Productivity**

# 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 environments like those of YouTube engineering teams, developers 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, conducted at Google Zurich, explores the integration of generative artificial intelligence into an internal Integrated Development Environment (IDE) to support software engineers in their day-to-day coding activities. It introduces a Large Language Model (LLM)-powered agent capable of performing in-depth code analysis, detecting complex violations, and offering clear, contextual suggestions for improvement. Embedded directly within the IDE through an extension, this intelligent assistant aims to enhance code quality, reduce technical debt, and support developer productivity at scale.

This report begins by presenting the context and challenges of applying AI in modern software development. It then outlines the design and integration of the AI agent, followed by an evaluation of its contribution to improving engineering workflows within a high-impact, real-world environment.

# Part I
# Project Foundation and Methodology

# 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  Presentation

Founded in 1998, Google LLC is a global leader in technology and innovation. 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 has established 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 I.1** – Google Logo

## 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 I.2** – Overview of some Google products

## 1.3  Focus Area

### 1.3.1  YouTube

Acquired by Google in 2006, YouTube has become 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 I.3** – YouTube Logo

### 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

## 2.1  Project Context

This project was developed within the scope of my host team YoutTube Dev Infra, which focuses on supporting developers by providing tools and extensions that improve their day-to-day workflows. As part of this mission, the team is exploring how artificial intelligence can be leveraged to further assist developers. One concrete idea is to enhance their existing extension by introducing a feature that helps enforce internal best practices. The goal is to explore how large language models (LLMs) can complement traditional approaches, offering developers more intelligent and context-aware guidance directly within the IDE.

In addition, this project is carried out as part of the National Institute of Applied Science and Technology's fifth-year mandatory final project, which is required to obtain the software

engineering degree.

## 2.2   Existing Solutions

To support best practices and maintain code quality, teams currently rely on three main approaches:

- **Code Reviews:** Engineers provide human feedback on design, readability, and best practices during review sessions. Recently, AI-assisted reviews have also been introduced, helping both reviewers and developers by suggesting improvements.

- **Presubmit Checks:** Automated checks executed before code is submitted. These ensure correctness, style consistency, and prevent simple errors.

- **Rule-Based Checks (Work in Progress):** Tools under development to enforce simple, objective rules such as naming conventions or syntax. These checks are intended for measurable guidelines and provide consistent, automated verification.

Together, these mechanisms form the current framework supporting developers in maintaining quality and consistency across the codebase.

## 2.3   Problem Statement

While these approaches provide valuable support, they also leave significant gaps in practice.

- **Delayed Feedback:** The most meaningful guidance on design and best practices typically arrives during code reviews, after development is complete. This often requires rework and slows iteration.

- **Limited Scope of Presubmit Checks:** Presubmits focus on correctness and safety rather than nuanced best practices, leaving developers without proactive guidance in those areas.

- **Surface-Level Coverage of Rule-Based Checks:** Rule-based tools can only enforce simple, objective rules. They are not able to reason about context-dependent or subjective best practices.

As a result, developers lack timely, intelligent support during the actual coding phase, where guidance would be most efficient and impactful.

## 2.4   Proposed Solution

This project introduces an **AI-assisted feedback system integrated directly into the coding workflow**.

- **Real-Time Guidance:** Provide developers with immediate, context-aware suggestions while they are writing code.

- **Framework-Specific Best Practice Enforcement:** Go beyond syntax and correctness by surfacing adherence to YouTube's internal framework guidelines early in the development process.

- **Reduced Review Burden:** Shift part of the best practice enforcement from manual reviews to the authoring stage, making reviews faster and more focused on higher-level insights.

By bringing intelligent, framework-aware feedback closer to the point of coding, the solution aims to reduce back-and-forth during reviews, ensure internal consistency, and accelerate development velocity.

# 3   Work Methodology

## 3.1   Agile Development Approach

Agile practices were adopted to support iterative development and maintain flexibility in responding to evolving requirements. This approach allowed continuous integration of feedback from the host and co-host, ensuring that each increment of work aligned with both technical goals and the broader product vision. Testing, validation, and code reviews were incorporated throughout the process to maintain high quality, while frequent collaboration provided clarity and shared ownership of outcomes. Agile principles also complemented Google's focus on engineering excellence, including rigorous design reviews, thorough testing, robust code reviews, and DevOps-enabled automation.

## 3.2   Kanban Workflow

The dynamic workload of the YouTube Developer Infrastructure (Dev Infra) team, including feature requests, bug fixes, and maintenance tasks, was managed using a Kanban workflow. By visualizing tasks and limiting work in progress, the team could prevent bottlenecks and quickly

shift focus to urgent issues when necessary. Work was structured into stages to maintain clear coordination while allowing the flexibility to adapt priorities as requirements evolved.

The Kanban workflow included the following stages:

- **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 I.4** – Kanban Workflow

## 3.3 Development Process

The project followed an iterative engineering cycle designed to balance thorough planning with incremental delivery. Each stage was supported by structured practices, dedicated tools, and regular collaboration:

- **Ideation and Research:** The project began with an exploration phase to clarify objectives, gather requirements, and investigate potential solution directions. This stage combined independent research with collaborative discussions to assess feasibility and align on priorities.

- **Design and Planning:** A detailed design document was authored to present the technical choices, architectural considerations, and proposed workflow. This document was reviewed by engineers, iteratively refined, and approved. The final work plan was then transferred to the internal task management system, enabling structured tracking and prioritization.

- **Implementation:** Development was performed in small, reviewable increments using Google's internal development environment within the company-wide repository. Each change was submitted accompanied by unit tests, and validated through both manual and AI-assisted code reviews. The implementation combined Python for backend logic, RPCs for inter-service communication, and TypeScript for the frontend.

- **Testing and Validation:** Functionality and reliability were verified continuously. Automated unit tests ensured correctness at the component level, while integration reviews validated the behavior within the larger system.

- **Maintenance and Optimization:** Refactoring, bug fixes, and updates were performed throughout development, particularly as some dependencies evolved or methods became deprecated. This ensured that the solution remained consistent, maintainable, and aligned with evolving standards.

Collaboration was 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 provided timely feedback, clear guidance, and alignment on shared dependencies.

**Figure I.5** – Project Development Iteration Cycle

## 3.4 Project Timeline

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

**Figure I.6** − Project Timeline - Detailed Schedule

# Conclusion

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

# Part II

# Partie 2

# Chapter   II

# Business Understanding and Comparative Analysis

**Summary**

## Introduction

This chapter establishes the theoretical and contextual foundation of the project. It begins with an overview of the software development lifecycle (SDLC) and the role of emerging AI technologies such as Large Language Models (LLMs), Generative AI, and AI agents in modern software engineering. It then presents the state of the art by reviewing existing solutions currently supporting code quality and best practices. Finally, it defines the project requirements, both functional and non-functional, showing how this work addresses existing gaps by integrating AI-driven assistance into a key phase of the SDLC.

# 1   Business and Reasoning

## 1.1   Software Development Life Cycle (SDLC)

The Software Development Life Cycle (SDLC) is a structured framework that defines the process of planning, creating, testing, deploying, and maintaining software systems. It ensures that software development follows a disciplined approach, promoting quality, efficiency, and maintainability while reducing risks and development costs. The SDLC helps teams manage complex projects by providing a roadmap from initial requirements to system retirement.



**Figure II.1** – Software Development Life Cycle Overview

**Key Steps in SDLC**

A typical SDLC consists of seven core phases, each serving a distinct purpose in ensuring that software meets functional and quality requirements:

1. **Planning:** This initial phase defines the project scope, objectives, resources, and timeline. Risk assessment, feasibility studies, and cost-benefit analyses are performed to ensure the project is viable before investing substantial resources.

2. **Requirements Analysis:** Stakeholders, including end-users and business analysts, collaborate to specify functional and non-functional requirements. Clear, documented requirements reduce ambiguity and serve as a reference for design and implementation. Techniques such as use cases, user stories, and requirement specifications are commonly used.

3. **System Design:** Architects and developers define the software architecture, data models, interface designs, and technology stack. This phase produces detailed design documents and prototypes that guide the implementation phase.

4. **Implementation (Coding):** Developers translate design specifications into source code using appropriate programming languages, frameworks, and tools. Best practices such as modularization, code reviews, and version control are critical in this phase.

5. **Testing:** The system undergoes rigorous verification to detect and correct defects. Testing can include unit testing, integration testing, system testing, performance testing, and user acceptance testing (UAT). Automated testing frameworks are increasingly used to ensure efficiency and repeatability.

6. **Deployment:** Once verified, the software is deployed to production environments. Deployment strategies may include phased rollouts, blue-green deployments, or canary releases to minimize risk and disruption.

7. **Maintenance:** After deployment, the software requires ongoing updates, bug fixes, performance improvements, and adaptations to changing requirements or environments. Maintenance ensures long-term system reliability and user satisfaction.

## Common SDLC Models

Various SDLC models provide different approaches to organizing these steps. Each model offers specific advantages and trade-offs, depending on project complexity, team structure, and business goals.

**Waterfall Model**   The Waterfall model is a linear and sequential approach where each phase must be completed before the next begins. Its structure ensures strict documentation and clear progress milestones.

**Figure II.2** – Waterfall SDLC Model

**Advantages:** Simple and easy to manage, clear documentation, predictable timelines.
**Limitations:** Inflexible to changing requirements, late discovery of defects, high risk in complex projects.

**V-Model (Verification and Validation)** The V-Model extends Waterfall by emphasizing verification and validation. Each development phase has a corresponding testing phase to ensure quality at every stage.



**Figure II.3** – V-Model SDLC

**Advantages:** Strong focus on testing, early defect detection.
**Limitations:** Like Waterfall, less flexible with changing requirements.

**Iterative and Incremental Model**   This approach develops software in small iterations, delivering functional increments at each cycle. Feedback from early releases informs subsequent iterations.



**Figure II.4** – Iterative and Incremental SDLC Model

**Advantages:** Flexibility to change, early delivery of usable software, reduced risk.
**Limitations:** Requires careful planning, can increase complexity if iterations are poorly managed.

**Agile Model**   Agile emphasizes collaboration, adaptability, and rapid delivery. Work is organized into sprints, typically 2-4 weeks, producing a potentially shippable product increment. Agile can be implemented through frameworks such as:

- **Scrum:** Defines roles (Product Owner, Scrum Master, Development Team) and ceremonies (Daily Stand-ups, Sprint Planning, Retrospectives) to organize work in sprints.

- **Kanban:** Visualizes workflow and limits work in progress, emphasizing continuous delivery and efficiency rather than fixed-length iterations.

**Figure II.5** – Agile SDLC Model

**Advantages:** Highly flexible, promotes continuous feedback, encourages customer involvement.

**Limitations:** Less predictability in timelines, requires disciplined team practices.

**DevOps**  DevOps is a cultural and technical approach that integrates development and operations teams, promoting continuous integration, continuous delivery (CI/CD), automated testing, and real-time monitoring. By combining SDLC phases with DevOps practices, organizations achieve faster delivery, higher reliability, and improved collaboration between teams.



**Figure II.6** – DevOps Cycle Integration with SDLC

**Advantages:** Accelerates delivery, improves collaboration, enables automated testing and

monitoring, enhances reliability.

**Limitations:** Requires organizational culture change, relies on mature tooling, can be complex to implement in large teams.

**Emerging Enhancements** Modern software development also increasingly relies on AI-powered tools that provide predictive analytics, code suggestions, and real-time quality checks. These enhancements further increase SDLC efficiency and reliability, setting the stage for intelligent, AI-driven development assistance.

## 1.2 Artificial Intelligence in Software Development

**Foundational AI Concepts**

Artificial Intelligence (AI) refers to computational systems capable of performing tasks that typically require human intelligence, such as reasoning, learning, problem-solving, and decision-making. AI encompasses a broad spectrum of techniques and paradigms, each with distinct capabilities and applications.

**Figure II.7** – Hierarchy of AI Technologies

The AI landscape includes several key categories:

- **Symbolic AI:** Systems based on rules and logic to represent knowledge and perform reasoning. These systems excel at tasks with well-defined rules but struggle with ambiguity and complex pattern recognition.

- **Machine Learning (ML):** Algorithms that improve performance on tasks by learning from data rather than being explicitly programmed. ML systems can identify patterns and make predictions based on historical data.

- **Deep Learning (DL):** A subset of ML using artificial neural networks with multiple layers to model complex patterns in data. Deep learning has revolutionized fields like computer vision and natural language processing.

- **Generative AI:** AI systems capable of producing new content, such as text, code, or images, by learning patterns from existing datasets. Unlike traditional AI that classifies or predicts, generative AI creates novel outputs.

### The AI Revolution in Software Engineering

As software development becomes increasingly complex and teams face mounting pressure to deliver high-quality code faster, traditional development approaches are reaching their limits. The emergence of Artificial Intelligence technologies presents a transformative opportunity to address these challenges by bringing intelligent assistance directly into the development workflow.

The integration of AI into software development represents more than just technological advancement—it addresses fundamental business challenges that organizations face in maintaining code quality, reducing technical debt, and scaling development practices across growing teams. This represents a paradigm shift from traditional SDLC approaches to what we can call the AI-Enhanced SDLC—a development lifecycle where intelligent assistance is embedded throughout every phase, creating a more efficient, proactive, and intelligent development process.

This transformation is primarily driven by two key AI technologies: Large Language Models (LLMs) and AI agents, which together enable the intelligent development workflows that characterize the modern AI-enhanced SDLC.

### Large Language Models: The Foundation of Intelligent Development

Large Language Models (LLMs) represent a breakthrough in generative AI technology, trained on massive corpora of natural language and source code. Unlike traditional programming tools that rely on rigid rules and patterns, LLMs understand context, semantics, and intent, making them uniquely suited for complex reasoning tasks across various domains.

**Figure II.8** – A chronological overview of large language models (LLMs) [2018-2024]

The fundamental capabilities of LLMs include:

- **Natural Language Understanding:** Ability to comprehend and interpret human language with nuanced understanding of context, intent, and meaning.

- **Pattern Recognition:** Capacity to identify complex patterns in data, including code structures, design patterns, and domain-specific conventions.

- **Content Generation:** Capability to produce coherent, contextually appropriate text, code, and other content based on learned patterns.

- **Reasoning and Inference:** Ability to perform logical reasoning, make inferences, and solve complex problems through step-by-step analysis.

In the context of software development, these capabilities translate into powerful applications:

- **Contextual Code Analysis:** Understanding code within its broader project context, including relationships between components, dependencies, and business logic.

- **Intelligent Code Generation:** Producing code suggestions, refactoring recommendations, and automated fixes that align with project-specific patterns and requirements.

- **Explanatory Documentation:** Providing clear, human-readable explanations of code behavior, potential issues, and suggested improvements.

- **Semantic Standards Enforcement:** Detecting violations of coding standards and best practices that go beyond simple syntax checking to include design patterns, architectural principles, and domain-specific guidelines.

**AI Agents: Orchestrating Intelligent Development Workflows**

While LLMs provide the foundational intelligence for understanding and generating code, AI agents represent the next evolution—autonomous systems that can reason, plan, and execute complex development tasks. AI agents serve as intelligent orchestrators that bridge the gap between human developers and AI capabilities, providing seamless integration of AI assistance into the development workflow.

**Agent Architecture and Workflow** The architecture of AI agents follows a structured workflow that enables reliable and scalable operation. As illustrated in Figure II.9, the agent orchestration process begins with input processing, where user requests are captured and preprocessed. The core reasoning component then analyzes the input using LLM capabilities to determine the appropriate course of action. This is followed by tool integration, where the agent interfaces with external systems to execute planned actions. Throughout this process, the memory layer maintains context and stores relevant information for future reference.

**AI AGENT**



**Figure II.9** – AI Agent Architecture and Orchestration Workflow

**Core Components of Development-Focused AI Agents** AI agents designed for software development incorporate specialized components tailored to the development workflow:

- **Code Analysis Engine:** Processes source code, understands project context, and identifies violations of coding standards, design patterns, and best practices.

- **Contextual Reasoning:** Utilizes LLMs to understand the broader project context, including dependencies, business logic, and architectural constraints when making recommendations.

- **Development Tool Integration:** Interfaces with IDEs, version control systems, testing frameworks, and other development tools to provide seamless assistance within the existing workflow.

- **Learning and Adaptation:** Continuously learns from codebase patterns, team feedback, and project outcomes to provide increasingly relevant and accurate assistance.

- **Feedback and Explanation System:** Provides clear, actionable explanations for identified issues and suggested improvements, helping developers understand not just what to change, but why.

**Balancing AI Capabilities with Practical Considerations** While AI agents offer transformative potential for software development, their deployment requires careful consideration of both capabilities and limitations:

**Key Capabilities:**

- **Contextual Understanding:** Ability to analyze code within its broader project context, understanding relationships and dependencies

- **Automated Analysis:** Consistent, comprehensive code analysis that scales across large codebases and development teams

- **Intelligent Recommendations:** Context-aware suggestions that go beyond simple rule-based checks to include design patterns and architectural considerations

- **Continuous Learning:** Adaptation to project-specific patterns and team preferences over time

**Practical Limitations:**

- **Computational Cost:** LLM-based analysis requires significant computational resources, impacting response times and operational costs

- **Context Constraints:** Limited ability to process extremely large files or maintain context across very long code sequences

- **Accuracy Considerations:** Potential for misinterpretation or hallucination, requiring human oversight for critical decisions

- **Integration Complexity:** Requires careful integration with existing development tools and workflows

These considerations inform the design of practical AI-assisted development systems that maximize benefits while managing limitations effectively.

## AI-Enhanced SDLC: The Complete Integration

The combination of LLMs and AI agents enables the complete integration of intelligent assistance throughout the software development lifecycle. This AI-enhanced SDLC represents a fundamental transformation from traditional reactive approaches to proactive, intelligent development assistance that addresses the core business challenges identified earlier.

**The Transformative Impact of AI Integration**  The integration of AI into the SDLC addresses key limitations of conventional development practices:

- **From Reactive to Proactive:** Traditional approaches identify issues after they occur, often during code review or testing phases. AI-enhanced development provides real-time feedback during coding, preventing issues before they become embedded in the codebase.

- **From Inconsistent to Scalable:** Human-based quality assurance varies in consistency and availability. AI provides uniform, expert-level analysis that scales across teams and projects without degradation.

- **From Static to Adaptive:** Traditional tools rely on fixed rules and patterns. AI systems learn and adapt to project-specific patterns, team preferences, and evolving best practices.

- **From Isolated to Integrated:** Rather than treating quality assurance as a separate phase, AI integrates intelligent assistance throughout the entire development workflow.

**AI Applications Across the Development Lifecycle** AI technologies provide targeted assistance at each phase of the SDLC, addressing specific challenges and opportunities:

- **Requirements and Planning:** AI analyzes historical project data to predict timelines, identify requirement ambiguities, and suggest optimal resource allocation. Natural language processing helps translate business requirements into technical specifications.

- **Design and Architecture:** AI assists architects by analyzing existing codebases, suggesting architectural patterns, detecting design anti-patterns, and ensuring consistency with organizational standards.

- **Implementation:** This is where AI agents provide the most immediate value, offering real-time code completion, best practice enforcement, bug detection, and refactoring suggestions directly within the development environment.

- **Testing and Quality Assurance:** AI generates comprehensive test cases, prioritizes test execution based on risk analysis, and identifies edge cases that human testers might miss.

- **Deployment and Maintenance:** AI monitors deployment health, predicts potential issues, and automatically suggests optimizations based on usage patterns and performance metrics.

This comprehensive understanding of AI's potential in software development provides the foundation for evaluating current solutions and identifying opportunities for improvement. The next section examines existing approaches to code quality and best practices enforcement, highlighting both their contributions and limitations in addressing the business challenges outlined above.

# 2   State of the Art: Existing Solutions

While the AI-enhanced SDLC presents a vision of intelligent, proactive development assistance, the current reality in most development environments, including our workspace, relies on more traditional approaches to code quality and best practices enforcement. Software engineers currently depend on a combination of human review, automated scripts, and rule-based systems to maintain code quality. These mechanisms vary in scope, effectiveness, and feedback timing, but as our analysis will show, they leave significant gaps in providing the real-time, intelligent assistance that characterizes the AI-enhanced development paradigm.

**Current Feedback Mechanisms**    The existing approaches to code quality and best practices enforcement can be categorized into several key mechanisms:

- **Code Reviews:** Human reviewers examine code for design quality, readability, maintainability, and adherence to standards. While this approach provides high-level, context-aware feedback, it often introduces delays and requires significant effort. Recently, AI-assisted review tools have been developed to accelerate this process by providing initial suggestions or flagging common issues.

- **Presubmit Checks:** Automated scripts that validate code before submission, enforcing style guides, ensuring compilation, and verifying simple correctness and safety constraints. Although fast and reliable, they primarily focus on surface-level checks and do not consider design or contextual issues.

- **Rule-Based Checks:** Systems that enforce coding conventions, naming schemes, and formatting standards. These tools provide consistency and objectivity but cannot reason about complex or context-dependent best practices.

**Limitations**    While valuable, these approaches leave significant gaps, particularly in providing timely and intelligent feedback during the active coding phase. This timeline illustrates where current feedback mechanisms fit into the developer workflow:



**Figure II.10** – Current Developer Workflow Feedback Timeline

During the **Coding Phase**, developers get some support from Rule-Based Checks, but these are limited to simple, objective rules and often miss the nuances of complex best practices.

The most insightful feedback on design and best practices typically comes during **Code Review**, but this happens after the initial development, making changes more disruptive.

Finally, **Presubmit Checks** before submission focus on correctness and safety, not on proactive best practice guidance.

This leaves a significant gap: there's no real-time, intelligent support for adhering to complex best practices while the developer is actively coding, which is the gap our project aims to fill.

| Approach | Strengths | Limitations |
|---|---|---|
| Code Reviews | Context-aware, high-level insights | Feedback delayed, time-consuming |
| Presubmit Checks | Fast, automated validation | Limited scope, mostly syntax and safety |
| Rule-Based Checks | Consistency, objectivity | Cannot handle complex or contextual practices |

**Tableau II.1** – Comparison of available software quality support solutions

**Impact of Delayed Feedback**  Following from the previous analysis, these workflow issues create concrete impacts on development:

- **Technical Debt Accumulation:** Delayed feedback and limited enforcement of best practices lead to accumulating technical debt and inconsistencies in code quality over time.

- **Prolonged Review Cycles:** Review cycles take longer because developers must iterate multiple times to address issues discovered late in the process.

- **Developer Frustration:** The repetitive cycle of late-stage corrections contributes to developer frustration and reduced productivity.

- **Inconsistent Quality Standards:** Without real-time guidance, adherence to best practices varies significantly across team members and projects.

- **Increased Development Costs:** The cost of fixing issues increases exponentially the later they are discovered in the development process.

This analysis clearly demonstrates why we need a solution that brings guidance earlier, directly into the authoring process, addressing the fundamental timing and intelligence gaps in current approaches.

**Trends and Emerging Practices** Modern development increasingly integrates AI-driven assistance, predictive code analysis, automated documentation, and intelligent testing recommendations. These trends aim to provide proactive, context-aware guidance within the IDE, reducing errors and improving developer efficiency. However, as our analysis shows, there remains a critical gap in providing real-time, intelligent feedback during the active coding phase.

The limitations identified in current approaches—delayed feedback, limited scope of automated checks, and inability to handle complex best practices—create a compelling case for integrating LLM-powered assistance directly into the development workflow. This represents the next evolution in development tooling, moving from reactive quality assurance to proactive, intelligent guidance that addresses the fundamental gaps in current solutions.

—

# 3 Project Requirements

To tackle the issue of delayed feedback and limited intelligent assistance, our solution integrates the power of large language models directly into the developer's workflow in the IDE.

The proposed system builds on gaps identified in existing solutions. It aims to provide real-time, AI-driven feedback directly in the developer workflow, while maintaining performance, scalability, and usability.

## 3.1 Functional Requirements

The system must:

- **Detect Framework Violations:** Identify violations of internal YouTube framework best practices and coding standards in real-time during development.

- **Provide Contextual Explanations:** Explain violations in clear, developer-friendly language with context-specific rationale that helps developers understand why certain patterns are problematic.

- **Generate Actionable Fixes:** Suggest specific, actionable fixes leveraging AI-generated solutions tailored to YouTube framework patterns and conventions.

- **Enable Developer Interaction:** Allow developers to accept, reject, or modify AI suggestions, providing full control over the implementation of recommended changes.

- **Maintain Contextual Relevance:** Ensure feedback appears in the appropriate location within the code and remains relevant and accurate as the developer continues coding.

- **Seamless IDE Integration:** Integrate seamlessly with the existing internal IDE, extension, and backend workflow without disrupting the developer's current development process.

## 3.2 Non-Functional Requirements

The system should also meet broader quality criteria:

- **Performance:** Provide fast, near real-time responses to avoid interrupting developer workflow.

- **Scalability:** Efficiently handle large codebases and multiple simultaneous users.

- **Maintainability:** Enable modular updates, addition of new AI models, or coding rules.

- **Reliability:** Ensure robustness in production environments with minimal downtime.

- **Security and Privacy:** Comply with organizational policies, ensuring safe handling of code and data.

- **Usability:** Deliver concise, context-aware, and minimally intrusive feedback.

- **Extensibility:** Easily add new rules, models, or integrations.

| Requirement Type | Description |
|---|---|
| Functional | Framework violation detection, contextual explanations, actionable fixes, developer interaction, contextual relevance, seamless IDE integration |
| Non-Functional | Performance, scalability, maintainability, reliability, security, usability, extensibility |

**Tableau II.2** – Summary of project requirements

**Rationale** These requirements address limitations identified in existing solutions by embedding proactive, context-aware feedback directly in the coding workflow, specifically targeting YouTube framework best practices. This approach enhances developer productivity, reduces framework-specific errors, and supports consistent adherence to internal YouTube development standards.

# Conclusion

This chapter established the business and theoretical foundation for integrating AI into software development workflows. Beginning with an analysis of the Software Development Life Cycle, it identified key challenges in maintaining code quality and consistency across growing development teams. The exploration of AI technologies—particularly Large Language Models and AI agents—revealed their potential for addressing these challenges through proactive, context-aware assistance.

The examination of existing solutions highlighted significant limitations: traditional tools provide limited analysis, and human review offers inconsistent feedback. There remains a critical gap in delivering real-time, intelligent assistance that scales with team growth. The AI-enhanced SDLC paradigm addresses these limitations by embedding intelligent assistance throughout the development process.

The project requirements defined in this chapter focus on real-time, AI-driven feedback for YouTube framework best practices that integrates seamlessly into the development workflow. The proposed solution bridges the gap between AI potential and practical development needs, setting the stage for the detailed system design presented in the following chapter.

# Chapter    III

# System Design and Architecture

**Summary**

## 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 the previous chapter, this chapter details the 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

## 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. The architecture consists of two main components that work together to deliver real-time best practice enforcement:

- **IDE**: The developer's workspace containing the YouTube IDE Extension, which works with the currently open file being analyzed and annotates results directly in the editor.

- **AI Agent Framework**: The processing layer containing the LLM Best Practices Agent, which serves as the core AI processing engine that uses internal LLMs to analyze code and suggest best practices.



**Figure III.1** – High-Level System Architecture

## 1.2  System Workflow

The system operates through a streamlined workflow that begins when a developer triggers analysis via the YouTube IDE Extension. The complete interaction flow is depicted in Figure III.2.

**Figure III.2** – System Interaction Sequence Diagram

The sequence unfolds through the following interactions and responsibilities:

1. **Analysis initiation**: The YouTube developer triggers analysis of the currently open file.

2. **Request submission**: The YouTube IDE Extension submits an analysis request to the agent, identifying the target file.

3. **Processing**: The agent performs best-practices analysis and prepares the resulting findings.

4. **Result delivery**: The agent returns a structured response to the YouTube IDE Extension.

5. **Presentation**: The YouTube IDE Extension presents the best practices violations and suggested fixes within the editor.

This workflow emphasizes decoupling the IDE from heavy AI computation, ensuring that the development environment remains responsive while delegating intensive analysis to the specialized agent framework.

# 2 Components Design

## 2.1 LLM Best Practices Agent

The LLM Best Practices Agent serves as the core intelligence engine of the system, responsible for analyzing code, identifying violations of YouTube framework best practices, and providing actionable feedback to developers. This component represents the convergence of modern AI capabilities with domain-specific software engineering expertise.

### 2.1.1 Agent Architecture Choice

The agent is built using the Executable Agent architecture, provided by our internal AI platform. This architecture offers a structured approach to orchestrating AI-powered workflows, and was selected after careful evaluation of different agent paradigms, considering reliability, performance, and maintainability.

**Executable Agent vs. ReAct Architecture** To motivate the choice of Executable Agent architecture, we contrast it with the more widely known ReAct (Reasoning and Acting) pattern. In our system, the Executable Agent refers to a deterministic, tool-orchestrated workflow where control flow is defined in code rather than delegated to LLM reasoning. For readers outside Google, this is conceptually similar to what Microsoft publicly describes as sequential orchestration [**?** ], though our implementation is independent and internally developed.

- **ReAct Agent**: Uses a reasoning loop where the LLM decides what action to take next, executes it, observes the result, and continues reasoning. This creates a dynamic, LLM-driven execution flow.

- **Executable Agent**: Follows a predefined, deterministic execution flow where the agent orchestrates a sequence of tool calls in a structured manner, with the LLM used primarily for processing within each tool rather than for orchestration decisions.

Figure III.3 illustrates the fundamental difference in control flow between the two approaches.

**Figure III.3** – Comparison of Executable Agent vs. ReAct Agent Architectures

The Executable Agent architecture offers several key advantages for this use case:

- **Deterministic Execution**: Predefined flow ensures predictable behavior and simplifies debugging.

- **Tool Orchestration**: Clean framework for coordinating multiple specialized tools.

- **Error Handling**: Built-in mechanisms for handling failures and graceful degradation.

- **Performance (Latency)**: Efficient execution with low response latency, since orchestration decisions are pre-programmed rather than deferred to open-ended LLM reasoning.

- **Cost Efficiency (Token Usage)**: Reduced token consumption by constraining LLM calls to focused, well-scoped processing steps rather than repeated reasoning loops.

- **Maintainability**: Clear separation of concerns between tools and responsibilities.

### 2.1.2   Core Tools Architecture

To achieve this structured workflow, the agent incorporates five specialized tools, each handling a specific step of the best practices analysis pipeline.

**ReadFileFromWorkspace Tool**   This tool serves as the entry point for code analysis, responsible for retrieving the complete content of the file being analyzed. It handles file system access, ensuring that the agent has access to the full context of the code under examination. The tool is designed to handle multiple file formats and encodings that might arise in different development environments. Design contract: input = file path; output = full file content.

**CodeAnalysisTool**   The core analysis engine that performs violation detection against YouTube framework best practices. This tool uses an internal Large Language Model to perform semantic code analysis, going beyond simple pattern matching to understand intent and context. It can identify violations related to:

- Framework-specific patterns and anti-patterns

- Component structure and organization

- Code complexity and maintainability

- Import patterns and dependencies

- Naming conventions

Design contract: input = file content + convention set; output = list of base violations.

**ViolationExplanationTool**   This tool leverages the internal Large Language Model to generate human-readable explanations for each identified violation, providing developers with clear understanding of why a particular code pattern violates best practices. The explanations are contextual and educational, helping developers not only fix the immediate issue but also understand the underlying principles. Design contract: input = base violation + surrounding code context + convention; output = natural-language explanation.

**CodeFixTool**   This tool uses the same internal Large Language Model to propose actionable fix suggestions for identified violations. It goes beyond simply identifying problems to provide concrete solutions that developers can implement. Design contract: input = base violation + explanation + code context; output = suggested fix. The fixes are designed to be:

- **Safe**: Only modify internal implementation without breaking public APIs

- **Self-contained**: Require no additional changes in other files

- **Contextual**: Take into account the specific code context and framework patterns

- **Educational**: Include comments explaining the reasoning behind the fix

**Finish Tool**  The consolidation component that aggregates all analysis results into a structured output format. This tool ensures that the response is properly formatted and contains all necessary information for the YouTube IDE Extension to display results effectively. It enforces the response schema, deduplicates entries, and merges overlapping ranges when applicable, producing a coherent set of violations, explanations, and fixes.

### 2.1.3  Processing Strategy

The agent employs a balanced processing strategy that weighs performance against reliability. This strategy represents a key architectural decision made after evaluating different processing approaches for handling multiple violations within a single file.

**Processing Strategy Options**  During the design phase, three main processing strategies were considered for handling multiple violations:

- **Fully Sequential Processing**: Each violation is processed one at a time, with explanation generation and fix creation happening sequentially for each violation. This approach ensures maximum reliability and predictability but may result in longer processing times for files with many violations.

- **Fully Parallel Processing**: All violations are processed concurrently, with explanation generation and fix creation happening simultaneously for all violations. This approach maximizes performance but introduces complexity in managing concurrent operations and may face reliability challenges under high load conditions.

- **Hybrid Parallel Processing with Concurrency Limiting**: A balanced approach that processes violations in parallel but with controlled concurrency to maintain system stability. This strategy provides significant performance improvements over sequential processing while ensuring reliable operation even under various load conditions.

The hybrid parallel processing strategy with concurrency limiting was selected as it provides the optimal balance between performance and reliability for production use. The internal workflow is illustrated in the activity diagram below:



**Figure III.4** – Agent Processing Activity Diagram

The workflow demonstrates the key architectural decisions made in the processing strategy:

**Sequential Initial Processing** The agent begins with sequential steps: file reading and violation detection. This ensures that the complete code context is available before any analysis begins, and all violations are identified before parallel processing starts.

**Parallel Violation Processing** When violations are found, the agent employs the hybrid parallel processing strategy with concurrency limiting. Each violation is processed independently, with explanation generation and fix creation happening concurrently across multiple violations while maintaining system stability. A default concurrency limit is applied, and excess tasks are queued to prevent resource saturation. LLM calls use bounded retries with

exponential backoff; on persistent failure, partial results are returned. Failures are isolated per violation—an error in one task does not stop processing of other violations.

**Result Consolidation**   The workflow concludes with result consolidation, ensuring that all parallel processing results are properly aggregated into a coherent response format for the YouTube IDE Extension.

### 2.1.4   Integration with LLM Infrastructure

The agent integrates with an internal AI platform that hosts multiple LLMs. We currently use the latest internal Gemini-based model trained on Google's codebase. The architecture is model-agnostic, enabling seamless adoption of newer Gemini-based models as they are released without requiring changes in the orchestration logic. LLM usage is isolated behind stable interfaces so higher-level logic remains unaffected. Three tools depend on the LLM: CodeAnalysisTool, ViolationExplanationTool, and CodeFixTool. Integration concerns include prompt construction, response parsing, error handling, and token-usage tracking. Token usage monitoring is crucial for both cost management and latency optimization, enabling capacity planning and performance tuning. This design ensures long-term maintainability and benefits from platform improvements without architectural change.

## 2.2   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.

### 2.2.1   Extension Architecture

The YouTube IDE Extension is designed as a lightweight client that orchestrates the interaction between the developer and the AI analysis system. The architecture ensures that the extension remains responsive while delegating computationally intensive analysis to the specialized agent framework.

The extension handles user interface concerns, progress indication, and result presentation, while the heavy lifting of code analysis and AI processing is handled by the dedicated agent infrastructure. This separation of concerns ensures that the development environment remains responsive and that AI capabilities can be scaled independently of the user interface components. The extension handles concurrent analysis requests efficiently, queuing multiple file analyses to prevent resource saturation.

### 2.2.2 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:

- **Editor Action Integration**: A lightbulb icon is strategically placed in the editor title bar, positioned directly above the code being analyzed. This placement ensures maximum visibility and easy discoverability, allowing YouTube developers to quickly identify when AI analysis is available. The icon serves as both a visual indicator and an interactive trigger, making the feature immediately apparent without cluttering the interface.

- **Command Palette Integration**: For developers who prefer keyboard-driven workflows, the YouTube IDE Extension provides a command accessible through the Command Palette: "Analyze File with AI". This approach ensures that the feature is accessible through standard IDE navigation patterns, supporting both mouse-driven and keyboard-driven user interactions.

**Feedback Mechanisms:** The YouTube IDE Extension employs sophisticated feedback mechanisms to present analysis results in a manner that integrates seamlessly with existing development workflows:

- **Progress Notification**: Notification messages show real-time progress during analysis, keeping developers informed of the current status and estimated completion time. This ensures transparency during the AI processing phase and prevents uncertainty about system responsiveness.

- **Diagnostic Integration**: Violations are displayed as native IDE diagnostics, showing detailed explanations of YouTube framework violations. This approach leverages YouTube

developers' familiarity with standard diagnostic patterns while providing enhanced intelligence through AI analysis. The diagnostic messages focus on explaining why a particular code pattern violates best practices.

- **Hover Provider Enhancement**: The hover provider delivers formatted code suggestions when developers interact with diagnostic markers. Since suggested fixes often contain code snippets, the hover interface provides proper formatting and syntax highlighting that cannot be displayed effectively in diagnostic messages. This approach ensures that developers receive well-formatted, actionable code solutions.

### 2.2.3 Handling Stale Diagnostics

A fundamental challenge in integrating AI-powered analysis into real-time development environments is balancing responsiveness with computational efficiency. While developers expect instant feedback on their code changes, invoking a powerful language model on every keystroke would result in prohibitive latency and resource consumption. This creates a classic engineering dilemma: how can we provide developers with immediate feedback that remains relevant and accurate as they continue to modify their code? The YouTube IDE Extension addresses this challenge through a two-tiered feedback system:
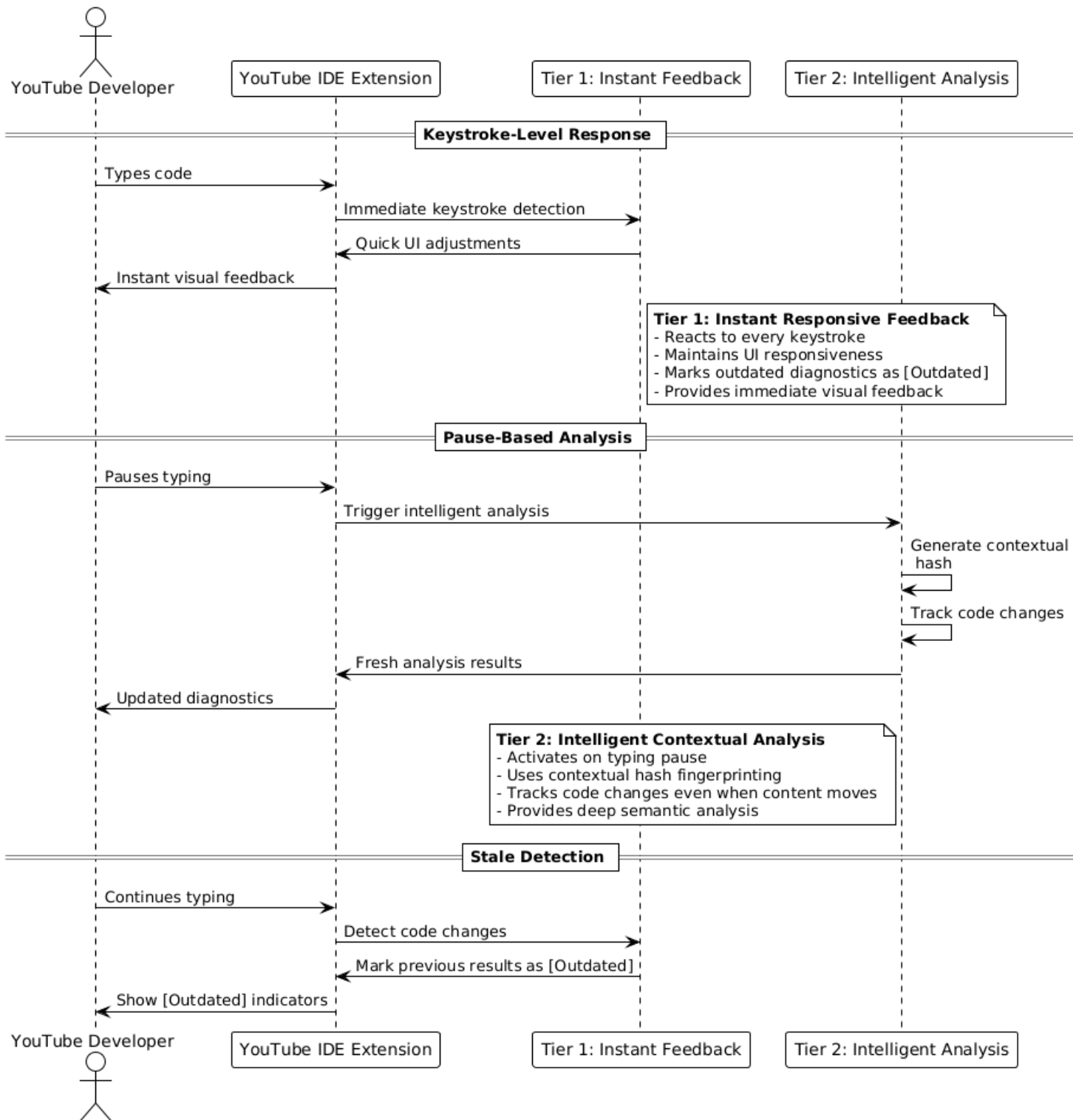
**Figure III.5** – Two-Tiered Feedback System for Handling Stale Diagnostics: Balancing instant responsiveness with intelligent analysis

**Tier 1: Instant Responsive Feedback**  The first tier provides immediate reaction to every keystroke, maintaining UI responsiveness through quick adjustments and intelligent state management. This tier ensures that the interface remains responsive and provides immediate visual feedback, even during rapid code changes. When analysis results become outdated due to code modifications, they are marked as [Outdated] to provide clear indication of their current relevance.

**Tier 2: Intelligent Contextual Analysis**  The second tier activates once the developer pauses typing, employing a debounced contextual hash system, creating a unique fingerprint of the code state. This debounced approach ensures that diagnostics are re-anchored correctly after a sufficient pause in typing, maintaining responsiveness while ensuring accurate positioning. The fingerprint allows the system to track code changes even when content moves within the file, ensuring that analysis results remain relevant and accurate. The intelligent analysis tier provides deep, semantic understanding of the code while maintaining system responsiveness.

**Seamless Integration**  The dual-tier system creates a seamless experience that balances speed with intelligence, providing developers with both immediate feedback and comprehensive analysis. This architecture represents a solution to the classic engineering problem of providing instant feedback without compromising on the depth and accuracy of AI-powered analysis.

### 2.2.4  User Interaction Flow

The complete user interaction flow with the YouTube IDE Extension is illustrated in the following diagram, showing the journey from analysis initiation to fix application:
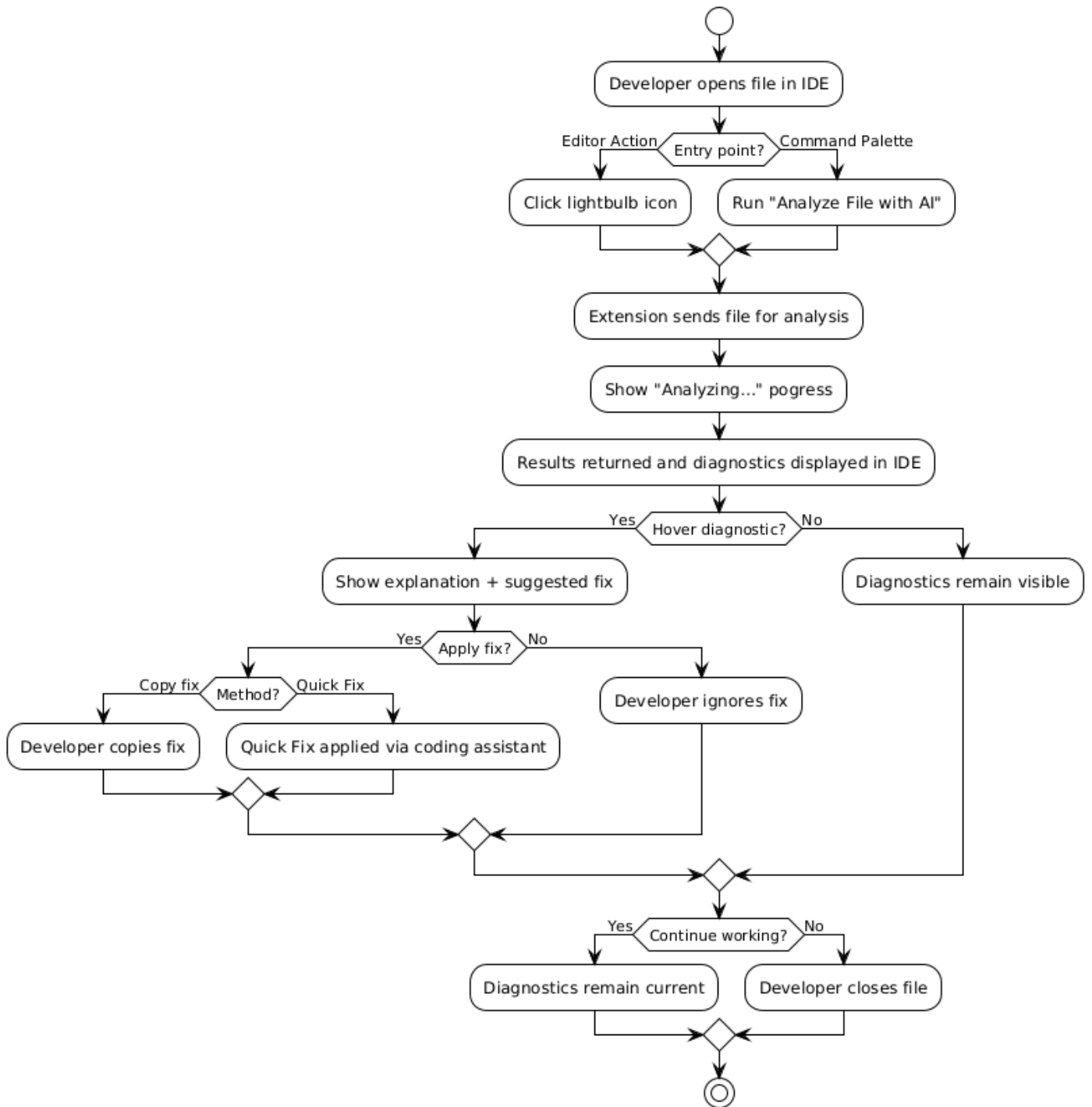
**Figure III.6** – YouTube IDE Extension User Interaction Flow: Complete developer journey from analysis trigger to fix application

The flow demonstrates the key interaction patterns and decision points that YouTube developers encounter when using the extension. The process begins with analysis initiation through

either the editor action or command palette, followed by real-time progress indication during AI processing. Once results are returned, violations are displayed as native IDE diagnostics with explanations. When developers hover over diagnostic markers, they receive both the explanation and formatted code fixes. The fix application process is entirely optional, allowing developers to either copy suggested changes or use the Quick Fix mechanism that calls the IDE's coding assistant. This design ensures that the AI assistance integrates naturally into existing development workflows without requiring any learning curve.

# 3 Data Models and Interfaces

## 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.

**Analysis Request Format**   The YouTube IDE Extension sends analysis requests to the LLM Best Practices Agent using a simple file path string:

- **File Path**: String containing the complete path to the file being analyzed. The file path must follow a specific format that includes the workspace ID and username for proper identification and access control.

This minimal input format ensures that the agent can focus on its core responsibility of code analysis while maintaining security and proper workspace isolation.

**Analysis Response Format**   The agent returns a structured JSON response with the following top-level fields:

- **status**: The overall status of the analysis: "success" or "error"

- **errorDetails**: Details about the error if status is "error", otherwise null

- **violationResults**: An array of violation objects, where each object represents a single violation found in the code

- **tokenUsage**: Optional object containing input_tokens and output_tokens if requested

Each violation object includes:

- **originalViolation**: Contains the data from the initial CodeAnalysisTool pass

  - **conventionId**: The unique identifier for the best practice that was violated (e.g., "local-components-complexity")

  - **range**: The location of the violation in the source file with precise line and character positions (start/end)

- **explanation**: The human-readable explanation of the violation, generated by the ViolationExplanationTool

- **suggestedFix**: A human-readable fix presented in the IDE hover, designed for clear rendering (formatted code snippets, e.g., TypeScript, or concise step-by-step instructions)

A minimal example response:

```
{
  "status": "success",
  "violationResults": [
    {
      "originalViolation": {
        "conventionId": "local-components-complexity",
        "range": {"start": {"line": 15, "character": 0},
                  "end": {"line": 15, "character": 20}}
      },
      "explanation": "This component exceeds the
recommended complexity threshold.",
      "suggestedFix": "Consider breaking this into smaller components."
    }
  ],
  "tokenUsage": {"input_tokens": 150, "output_tokens": 45}
}
```

## 3.2  Convention Data Management

YouTube framework best practices are stored as JSON files in the repository and imported by the agent at runtime. This approach provides version control integration and enables easy maintenance of convention definitions.

**Rationale for JSON-Based Convention Storage**   The system stores YouTube framework best practices as JSON files rather than using a traditional database for several reasons:

- **Read-Only Nature**: Conventions are static during runtime, requiring no dynamic updates, making a simple file sufficient.

- **Version Control**: Storing conventions in the repository allows easy tracking of changes, rollbacks, and collaboration with minimal infrastructure.

- **Simplicity & Portability**: JSON files are lightweight, easy to parse, and can be loaded efficiently at agent startup without introducing database dependencies.

- **Performance Consideration**: Loading all conventions into memory ensures low-latency access for the AI agent during code analysis.

- **Prototype Context**: This design targets a single framework for prototype development, allowing quick iteration. The architecture remains flexible—if multi-framework support or dynamic updates are required in the future, the JSON layer can be replaced with a database or service with minimal impact on other components.

This approach balances maintainability, simplicity, and performance while demonstrating awareness of potential scalability needs.

**Convention Data Structure**   The convention data is organized in a hierarchical format that supports both human readability and programmatic access:

- **Convention Definitions**: Each convention includes:

  - **Unique ID**: Identifier for programmatic reference
  - **Category**: Grouping (e.g., "Component Structure", "Import Patterns", "Naming Conventions")
  - **Description**: Clear explanation of the best practice
  - **Examples**: Code snippets showing correct and incorrect implementations
  - **Rationale**: Explanation of why this practice improves code quality
  - **Severity**: Default severity level for violations

- **Dynamic Retrieval**: Context-specific data inclusion in prompts based on:

  - File type and framework context
  - Previously detected violations in the same file

**Integration Interfaces**   The following minimal, versioned boundaries keep components decoupled. Detailed request/response fields are defined earlier in Input/Output Specification.

- **IDE Extension Interface**: AnalysisRequest/AnalysisResponse contract between the YouTube IDE Extension and the LLM Best Practices Agent. Requests include the file path (workspace + username); responses follow the Analysis Response Format. This boundary lets the UI evolve independently from the agent.

- **Convention Access**: The agent loads a local JSON array of convention definitions at startup (read-only). This simple import avoids network dependencies and provides direct access to IDs, descriptions, examples, and severity.

- **Monitoring Interface**: Captures token usage from LLM method calls and aggregate counts of violations by category for observability. Latency and accuracy measurements are discussed in the Evaluation chapter.

## 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 diagnostics, 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.

# Conclusion and Perspectives

C'est l'une des parties les plus importantes et pourtant les plus négligées du rapport. Ce qu'on ne veut pas voir ici, c'est combien ce stage vous a été bénéfique, comment il vous a appris à vous intégrer, à connaître le monde du travail, etc.

Franchement, personne n'en a rien à faire, du moins dans cette partie. Pour cela, vous avez les remerciements et les dédicaces, vous pourrez vous y exprimer à souhait.

La conclusion, c'est très simple : c'est d'abord le résumé de ce que vous avez raconté dans le rapport : vous reprenez votre contribution, en y ajoutant ici les outils que vous avez utilisé, votre manière de procéder. Vous pouvez même mettre les difficultés rencontrées. En deuxième lieu, on y met les perspectives du travail : ce qu'on pourrait ajouter à votre application, comment on pourrait l'améliorer.

# Appendix : Miscellaneous remarks

- Un rapport doit toujours être bien numéroté;

- De préférence, ne pas utiliser plus que deux couleurs, ni un caractère fantaisiste;

- Essayer de toujours garder votre rapport sobre et professionnel;

- Ne jamais utiliser de je ni de on, mais toujours le nous (même si tu as tout fait tout seul);

- Si on n'a pas de paragraphe 1.2, ne pas mettre de 1.1;

- TOUJOURS, TOUJOURS faire relire votre rapport à quelqu'un d'autre (de préférence qui n'est pas du domaine) pour vous corriger les fautes d'orthographe et de français;

- Toujours valoriser votre travail : votre contribution doit être bien claire et mise en évidence;

- Dans chaque chapitre, on doit trouver une introduction et une conclusion;

- Ayez toujours un fil conducteur dans votre rapport. Il faut que le lecteur suive un raisonnement bien clair, et trouve la relation entre les différentes parties;

- Il faut toujours que les abréviations soient définies au moins la première fois où elles sont utilisées. Si vous en avez beaucoup, utilisez un glossaire.

- Vous avez tendance, en décrivant l'environnement matériel, à parler de votre ordinateur, sur lequel vous avez développé : ceci est inutile. Dans cette partie, on ne cite que le matériel qui a une influence sur votre application. Que vous l'ayez développé sur Windows Vista ou sur Ubuntu n'a aucune importance;

- Ne jamais mettre de titres en fin de page;

- Essayer toujours d'utiliser des termes français, et éviter l'anglicisme. Si certains termes sont plus connus en anglais, donner leur équivalent en français la première fois que vous les utilisez, puis utilisez le mot anglais, mais en italique;

- Éviter les phrases trop longues : clair et concis, c'est la règle générale !

Rappelez vous que votre rapport est le visage de votre travail : un mauvais rapport peut éclipser de l'excellent travail. Alors prêtez-y l'attention nécessaire.