

ZUCCHETTI

Specifica Tecnica

Gini

Integrazione di LLM in Git lato versionamento

Leonardo Trolese

Indice

1	Introduzione	3
1.1	Scopo del documento	3
1.2	Obiettivi del progetto	3
2	Tecnologie	3
2.1	Stack tecnologico	3
2.2	Linguaggi utilizzati	4
2.3	Framework e librerie	4
3	Architettura Logica	4
3.1	1. CLI (Interfaccia a Riga di Comando)	5
3.2	2. Application Layer	5
3.3	3. Dominio	5
3.4	4. Infrastructure	5
3.5	5. Configurazione	5
3.6	6. Scripts	5
3.7	Relazioni tra i moduli	5
3.8	Diagramma logico	6
4	Architettura di Sistema	6
4.1	Panoramica Generale	6
4.2	Modulo 1 – Generazione ed esecuzione di comandi Git	6
4.2.1	Struttura e responsabilità	7
4.2.2	Pattern utilizzati	7
4.2.3	Estendibilità	8
4.2.4	Sintesi del flusso	8
4.3	Modulo 2 – Generazione di messaggi di commit	8
4.3.1	Struttura e responsabilità	8
4.3.2	Pattern utilizzati	9
4.3.3	Estendibilità	9
4.3.4	Sintesi del flusso	9

1 Introduzione

1.1 Scopo del documento

Questo documento descrive l'architettura logica e di sistema del software sviluppato nell'ambito del progetto di stage presso Zucchetti S.p.A., svolto dallo studente Leonardo Trolese. La Specifica Tecnica è rivolta a fornire una visione strutturata e dettagliata delle componenti del sistema, delle interazioni tra i moduli e delle scelte tecnologiche adottate. Essa costituisce un riferimento per lo sviluppo, il testing, la manutenzione e l'estensione futura del plugin realizzato.

1.2 Obiettivi del progetto

L'obiettivo principale del progetto è la realizzazione di un plugin per il software Git, che introduca in quest'ultimo la possibilità di gestire e generare testo in linguaggio naturale, in accompagnamento ad alcune delle principali funzionalità del software. Le principali features implementate dal plugin sono:

- generazione di comandi Git per CLI a partire da query in linguaggio naturale che descrivano l'intenzione dell'utente;
- generazione automatica di commenti di commit a seguito di analisi del codice modificato;
- l'analisi dell'impatto delle modifiche correnti sul repository remoto;
- utilizzo di LLM per la risoluzione automatizzata di conflitti di merge.

Le attività saranno accompagnate da sessioni di test e dalla produzione di documentazione tecnica, al fine di assicurare l'affidabilità e la riusabilità della soluzione finale.

2 Tecnologie

Questa sezione descrive le tecnologie impiegate nello sviluppo del progetto, con particolare attenzione agli strumenti software, ai linguaggi di programmazione e alle librerie che compongono l'ecosistema applicativo. La scelta di queste tecnologie è stata guidata da criteri di compatibilità, efficienza, supporto alla modularità e disponibilità di risorse open source.

2.1 Stack tecnologico

Il progetto è stato sviluppato all'interno di un ambiente Python, con il supporto di strumenti e librerie open source orientati all'elaborazione del linguaggio naturale, alla generazione di codice e alla manipolazione di repository Git.

Il software si affida al database vettoriale ChromaDB per la memorizzazione del corpo di intenti utilizzato in fase di Rag per la raccolta di esempi coerenti con la query utente, da fornire come contesto all'LLM.

Oltre a Python come base di sviluppo e ChromaDB come database, il sistema si affida a un software esterno per l'interazione con LLM. Trattasi di Ollama, utilizzato per l'interazione con LLM, e fruibile in python mediante l'apposita libreria Ollama Python.

2.2 Linguaggi utilizzati

Nel corso del progetto sono stati utilizzati i seguenti linguaggi:

- **Python:** linguaggio principale impiegato per lo sviluppo del plugin e delle sue componenti interne. La scelta è motivata dalla vasta disponibilità di librerie per l'IA, l'elaborazione del linguaggio naturale e l'interazione con Git.
- **JSON:** formato di serializzazione leggero utilizzato per la configurazione dell'applicazione e la definizione di dati strutturati, come input di esempio o risposte attese.

2.3 Framework e librerie

Il sistema integra numerosi pacchetti Python e strumenti software open source, ciascuno con un ruolo specifico all'interno del flusso operativo:

- **chromadb:** libreria Python che gestisce l'interazione con l'omonimo database vettoriale utilizzato per la memorizzazione e la ricerca di rappresentazioni semantiche (embedding) di documenti testuali. È un componente chiave per la fase di retrieval nel contesto RAG.
- **rank_bm25:** implementazione dell'algoritmo BM25, impiegata per il ranking dei documenti recuperati sulla base della rilevanza testuale. Viene utilizzata in combinazione con le ricerche vettoriali per aumentare la precisione.
- **ollama:** libreria Python utilizzata per interagire con l'omonimo pensato per l'utilizzo locale di modelli LLM. Consente di eseguire inferenze senza dipendere da servizi cloud esterni, garantendo maggiore controllo e privacy.
- **GitPython:** libreria che permette di interagire in modo programmatico con repository Git. È fondamentale per l'automazione delle operazioni di commit, merge e generazione di commenti.
- **attrs:** utilizzata per la definizione concisa e robusta di classi dati in Python, facilitando la scrittura di codice manutenibile e ben strutturato.
- **PyMuPDF:** libreria opzionale utilizzata per il parsing e la gestione di file PDF, utile in eventuali casi di estrazione o analisi di contenuti documentali.
- **colorama:** impiegata per la gestione della colorazione dell'output su terminale, migliorando l'esperienza utente durante l'utilizzo del plugin in modalità CLI.
- **pytest, pytest-cov:** strumenti dedicati al testing automatico delle componenti software. Consentono l'esecuzione di test unitari e la valutazione della copertura del codice, a supporto della qualità del progetto.

L'adozione di queste tecnologie ha permesso di costruire un sistema robusto, flessibile e facilmente estendibile, in linea con gli obiettivi di progetto e le pratiche moderne di sviluppo software.

3 Architettura Logica

L'architettura logica del software nella cartella `src` segue un approccio modulare e a livelli, ispirato ai principi della Clean Architecture e del Domain-Driven Design. I principali moduli individuati sono:

3.1 1. CLI (Interfaccia a Riga di Comando)

La cartella `cli/` contiene i punti di ingresso dell'applicazione (`main.py`), la gestione dell'input/output utente (`user_io.py`) e il parser dei comandi (`command_parser.py`). Questo livello si occupa esclusivamente dell'interazione con l'utente e della raccolta delle richieste, delegando la logica applicativa ai moduli sottostanti.

3.2 2. Application Layer

La cartella `application/` implementa la logica di orchestrazione delle funzionalità principali tramite handler specializzati:

- `cmd_conversation_handler.py`, `commit_conversation_handler.py`, `commit_impact_handler.py`, `merge_conflict_handler.py`:

gestiscono i diversi flussi conversazionali e le operazioni core (generazione comandi, commit, analisi impatto, risoluzione conflitti).

- Sottocartelle come `rag/` contengono la logica per la costruzione del contesto tramite RAG (Retrieval-Augmented Generation).

3.3 3. Dominio

La cartella `domain/` definisce le strutture dati centrali e i prompt per l'interazione con il modello LLM:

- `chat.py`, `prompts.py`, `response_structure.py`: rappresentano le entità, i messaggi e i formati di risposta usati in tutto

il sistema.

3.4 4. Infrastructure

La cartella `infrastructure/` fornisce servizi di basso livello e componenti tecnici:

- `embedding/`: pipeline di embedding, chunking e gestione RAG.
- `git_service/`: astrazioni per l'interazione con Git (diff, repository, annotazioni).
- `llm/`: client per la comunicazione con modelli LLM (es. Ollama).
- `database/`: gestione del database vettoriale per il retrieval.
- File contenente gli esempi che costituiscono il corpo del sistema di RAG (`examples.json`).

3.5 5. Configurazione

La cartella `config/` contiene la configurazione centralizzata del sistema (`config.py`, `config.json`), permettendo la parametrizzazione dei modelli e delle lingue.

3.6 6. Scripts

La cartella `scripts/` ospita utility e script di supporto che non fanno strettamente parte del software eseguibile, ma sono stati utilizzati nel corso della produzione di tale software.

3.7 Relazioni tra i moduli

- Il livello CLI riceve input dall'utente e invoca gli handler dell'application layer.

- Gli handler orchestrano le operazioni, interagendo con il dominio per la gestione dei dati e con l'infrastruttura per

i servizi tecnici (Git, embedding, LLM).

- La configurazione è accessibile trasversalmente da tutti i moduli.
- L'infrastruttura è progettata per essere facilmente sostituibile o estendibile (es. nuovi modelli LLM, database, chunker).

3.8 Diagramma logico

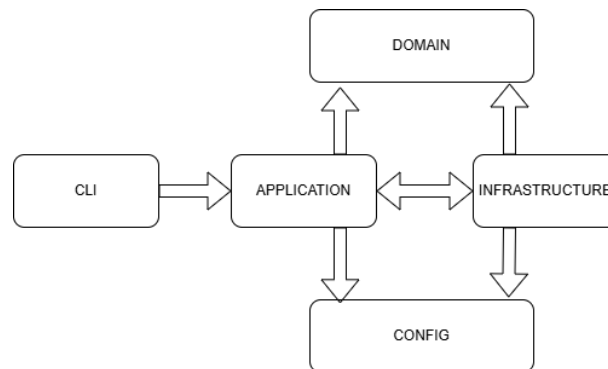


Figure 1: Diagramma logico dell'architettura del software

Questa architettura garantisce separazione delle responsabilità, testabilità e facilità di estensione, rendendo il sistema robusto e manutenibile.

4 Architettura di Sistema

4.1 Panoramica Generale

Il sistema descritto è un plugin progettato per assistere gli utenti nell'utilizzo avanzato di Git tramite un'interfaccia a riga di comando intelligente, potenziata da modelli di linguaggio (LLM) e tecniche di ricerca ibrida realizzata mediante Retrieval-Augmented Generation (RAG) e algoritmo di ranking BM25. L'architettura è pensata per garantire modularità, estendibilità e facilità di manutenzione, ed è stata realizzata applicando best practices come dependency injection, e pattern specifici per la risoluzione di problemi ricorrenti riscontrati durante la realizzazione del software.

Il plugin implementa quattro funzionalità principali, ciascuna orchestrata da moduli dedicati, che coprono la generazione di comandi Git, la creazione di messaggi di commit, l'analisi dell'impatto delle modifiche e la risoluzione automatica dei conflitti.

4.2 Modulo 1 – Generazione ed esecuzione di comandi Git

Il modulo di generazione comandi è progettato secondo un'architettura orientata ai pattern Command e Adapter, con particolare attenzione all'estendibilità e alla separazione delle responsabilità.

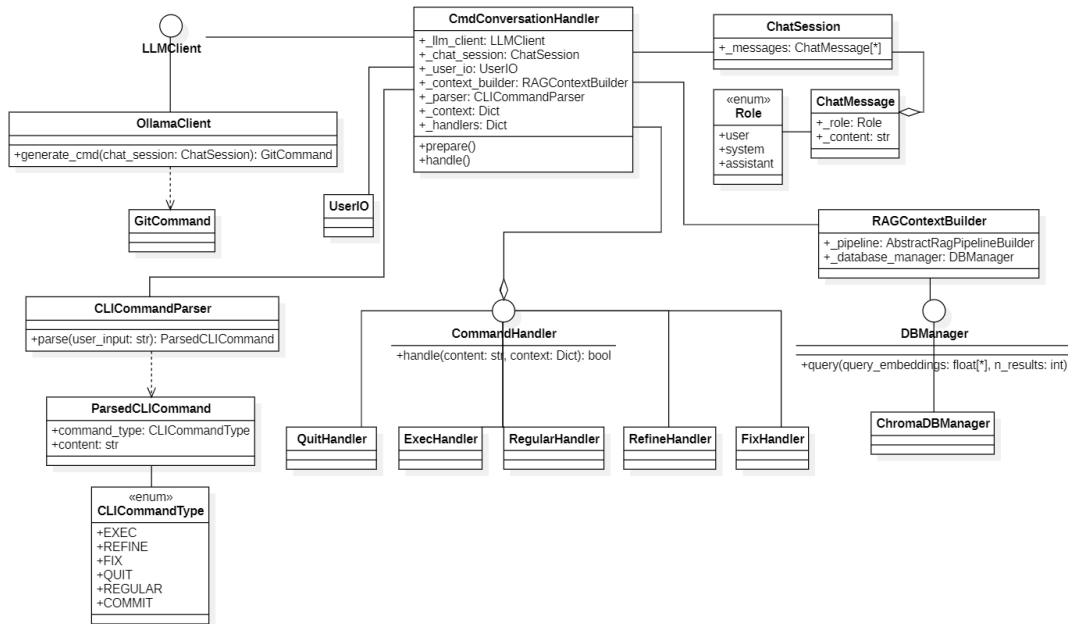


Figure 2: Diagramma logico dell'architettura del software

4.2.1 Struttura e responsabilità

Il cuore del modulo è rappresentato dalla classe `CmdConversationHandler`, che funge da *orchestratore*: riceve l'input dell'utente tramite `UserIO`, costruisce il contesto necessario (con `RAGContextBuilder`), e delega la gestione del comando al parser (`CLICommandParser`). Quest'ultimo trasforma l'input in un oggetto `ParsedCLICommand`, identificando il tipo di comando richiesto.

La gestione effettiva dei comandi è affidata a una gerarchia di handler che implementano il pattern *Command*: la classe astratta `CommandHandler` definisce l'interfaccia comune (`handle(content, context)`), mentre le sottoclassi concrete (`ExecHandler`, `RegularHandler`, `RefineHandler`, `FixHandler`, `QuitHandler`) implementano la logica specifica per ciascun tipo di comando. Questo approccio consente di aggiungere facilmente nuovi tipi di comando semplicemente introducendo nuovi handler, senza modificare la logica esistente.

Il recupero del contesto è invece demandato alla classe `RAGContextBuilder` che attraverso l'attributo `_pipeline: AbstractRagPipelineBuilder`, estrae dal database gli intenti più coerenti con la query dell'utente, andando poi a raffinare la ricerca attraverso l'applicazione, su di essi, dell'algoritmo di ranking BM25. Nell'insieme questa componente realizza una ricerca ibrida.

4.2.2 Pattern utilizzati

- **Command Pattern:** Ogni comando dell'utente viene incapsulato in un oggetto handler dedicato, che può essere eseguito in modo indipendente. Questo favorisce l'estensibilità e la separazione delle responsabilità.
- **Adapter Pattern:** L'interazione con componenti esterni o potenzialmente variabili (come il modello LLM) avviene tramite interfacce astratte (`LLMClient`). Ad esempio, `OllamaClient` implementa l'interfaccia `LLMClient`, o `ChromaDBManager` che implementa `DBManager`, permettendo di sostituire facilmente il backend LLM senza impattare il resto del sistema.

4.2.3 Estendibilità

L'architettura è pensata per facilitare l'aggiunta di nuove funzionalità: per supportare un nuovo tipo di comando, è sufficiente implementare una nuova sottoclasse di `CommandHandler` e registrarla nell'orchestratore. Analogamente, l'uso di interfacce astratte per i servizi esterni (LLM, database, ecc.) permette di adattare il sistema a nuove tecnologie o provider con modifiche minime.

4.2.4 Sintesi del flusso

1. L'utente inserisce una query tramite CLI.
2. `CmdConversationHandler` raccoglie l'input, e lo passa al parser.
3. `CLICommandParser` restituisce un oggetto `ParsedCLICommand` che identifica il tipo di comando.
4. L'orchestratore seleziona l'handler appropriato e ne invoca il metodo `handle`.
5. Gli handler possono interagire con il modello LLM (tramite l'interfaccia `LLMClient`), con la pipeline di ricerca ibrida e con altri servizi secondo necessità.
6. Il risultato viene restituito all'utente tramite `UserIO`.

Questa organizzazione garantisce modularità, testabilità e facilità di evoluzione del modulo di generazione

4.3 Modulo 2 – Generazione di messaggi di commit

4.3.1 Struttura e responsabilità

Il modulo per la generazione dei messaggi di commit è centrato sulla classe `CommitConversationHandler`, che svolge il ruolo di orchestratore del flusso conversazionale. Essa riceve l'input utente tramite `UserIO`, mantiene lo stato della conversazione (`ChatSession`) e utilizza il parser (`CLICommandParser`) per interpretare i comandi. L'interazione con il modello LLM avviene tramite l'interfaccia `LLMClient`, la cui implementazione concreta (`OllamaClient`) si occupa di generare la risposta (`LLMResponse`). La risposta del modello può essere di due tipi, distinti dal campo `mode`:

- **COMMIT:** contiene un oggetto `CommitMsg` con titolo e corpo del messaggio di commit generato.
- **QUESTION:** contiene una lista di domande (`Questions`) che il sistema pone all'utente per chiarire o arricchire il contesto prima di generare il messaggio di commit definitivo.

La gestione delle azioni conseguenti al comando è affidata a una gerarchia di handler che estendono `CommandHandler`, tra cui `CommitHandler` e `RefineCommitHandler`, responsabili rispettivamente della generazione e della raffinazione dei messaggi di commit.

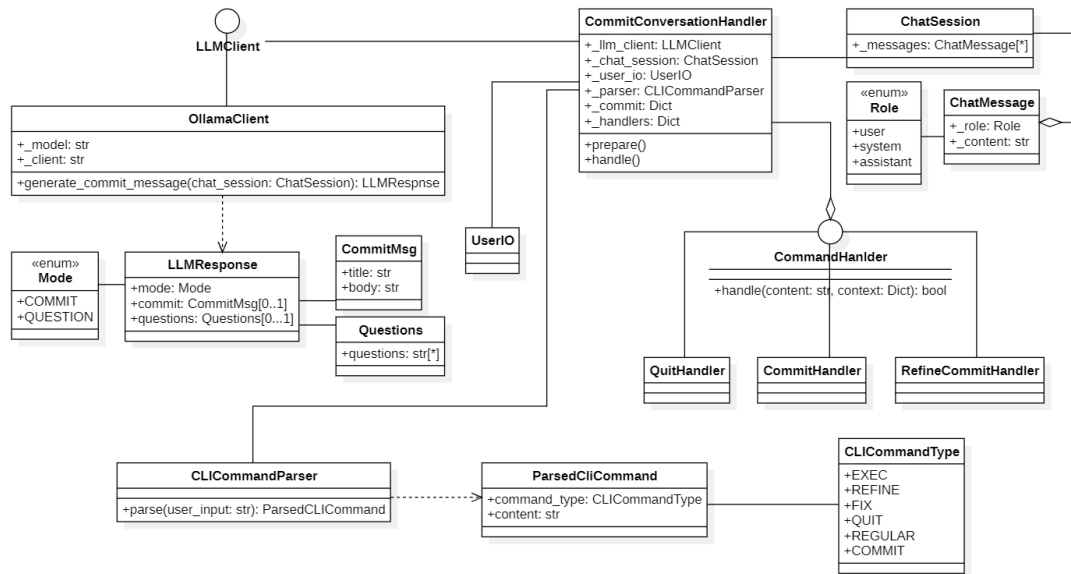


Figure 3: Diagramma logico dell'architettura del software

4.3.2 Pattern utilizzati

- **Command Pattern:** ogni tipo di comando (commit, refine, quit, ecc.) è gestito da un handler dedicato

che implementa una logica specifica, facilitando la separazione delle responsabilità e l'estendibilità del sistema.

4.3.3 Estendibilità

La struttura modulare consente di aggiungere facilmente nuovi tipi di handler per gestire ulteriori modalità di interazione o nuove tipologie di comandi, senza modificare la logica esistente. La separazione tra orchestratore, parser, handler e interfaccia verso il modello LLM garantisce una facile manutenibilità e possibilità di evoluzione del modulo.

4.3.4 Sintesi del flusso

1. L'utente inserisce un comando che avvia la generazione del messaggio di commit tramite CLI.
2. CommitConversationHandler raccoglie l'input, aggiorna la sessione e lo passa al parser.
3. Il comando viene interpretato e gestito dall'handler appropriato (CommitHandler, RefineCommitHandler, ecc.).
4. Il CommitConversationHandler interagisce con il modello LLM per generare una risposta (LLMResponse).
5. Se la risposta è di tipo *COMMIT*, viene proposto un messaggio di commit all'utente.
6. Se la risposta è di tipo *QUESTION*, vengono poste domande di chiarimento all'utente e il ciclo riprende fino a ottenere tutte le informazioni necessarie.
7. Il risultato finale viene restituito all'utente tramite UserIO.