

Corso di Laurea in Ingegneria e Scienze Informatiche

Integrazione di RAG e LLM nello Sviluppo del Software

Tesi di laurea in:
PROGRAMMAZIONE AD OGGETTI

Relatore

Prof. Viroli Mirko

Candidato

Bollini Simone

Correlatori

Dott. Aguzzi Gianluca

Dott. Farabegoli Nicolas

Abstract

I Large Language Model (LLM) addestrati per sviluppare il codice sono oggi altamente efficaci e in grado di generare soluzioni utili e funzionanti. L'addestramento fatto sui modelli è però su fonti e soluzioni generali, questo non dà quindi la possibilità al modello di proporre soluzioni su misura per una specifica richiesta utilizzando casistiche già create dal programmatore o dalla propria azienda per casi simili. Da questo nasce l'esigenza di addestrare il modello per personalizzare le soluzioni proposte, contestualizzandole alla propria realtà aziendale e al proprio stile nel programmare. Il LLM non conosce le librerie interne dell'azienda, i pattern di programmazione adottati e quindi le risposte ottenute sono troppo generiche. Per rispondere a questa esigenza entra in gioco la Retrieval-Augmented Generation (RAG) ovvero il processo di ottimizzazione dell'output di un LLM, per permettergli di fruire di un'ulteriore base di conoscenza personalizzata, unica e privata. L'output di questa **matrice di conoscenza** si inserisce e completa la query inviata al LLM, estendendo la base di informazioni sulla quale genererà l'output con la risposta. Questa tesi approfondisce questi concetti e sperimenta l'integrazione di un RAG con un LLM con lo scopo di ottenere dal LLM risposte personalizzate che solo con la conoscenza del LLM anche se estremamente performante e preparato sarebbe stato impossibile ottenere.

*A Giulia e ai miei figli, il dono più grande.
A tutta la mia famiglia.*

Grazie a tutti voi.

Contents

Abstract	iii
1 Introduzione	1
1.1 Essere programmatori nel 2025	1
2 Addestrare un LLM per la Generazione del Codice	5
2.1 Raccolta e Preparazione dei Dati	5
2.2 Pre-Addestramento	6
2.3 Fine-Tuning	7
2.4 Pre-Addestramento vs Fine-Tuning	8
2.4.1 Pre-Addestramento	9
2.4.2 Fine-Tuning	9
2.5 Architettura del Modello	9
2.6 Valutazione e Ottimizzazione	10
2.6.1 Metriche di Valutazione	10
2.6.2 Tecniche di Ottimizzazione	10
3 RAG	11
3.1 Introduzione	11
3.2 Funzionamento	12
3.2.1 Creazione degli Embedding	12
3.2.2 Fase 1: Function Calling	13
3.2.3 Fase 2: Recupero delle Informazioni	13
3.2.4 Fase 3: Aumento del Prompt	13
3.3 Gestione dell'Aggiornamento dei Dati	13
4 Caso Studio: Implementazione di un Sistema RAG per lo Sviluppo del Software	15
4.1 Obiettivi del Caso Studio	15
4.2 Architettura del Sistema	16
4.3 Software Utilizzati	16

CONTENTS

4.3.1	Ollama	16
4.3.2	LLM	16
4.3.3	BGE-M3	17
4.4	Dataset	17
4.5	Implementazione	17
4.5.1	Modello embedding	17
4.5.2	Creazione Rag	17
4.5.3	AI Agent	17
5	Conclusioni	19
5.1	Risultati Ottenuti	19
5.2	Impatto sullo Sviluppo Software	19
5.3	Sfide e Prospettive Future	19
		21
	Bibliography	21

List of Figures

3.1	Flusso di una request ad un LLM integrato con un RAG	12
4.1	Visualizzazione degli embeddings con TSNE	18

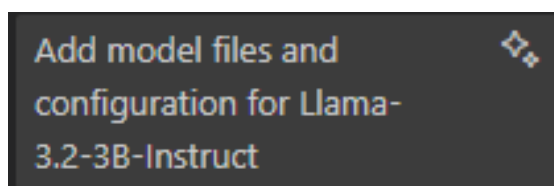
LIST OF FIGURES

Chapter 1

Introduzione

1.1 Essere programmatori nel 2025

Sono disponibili tantissimi (IDE) per lo sviluppo del codice uno di questi è **Visual Studio Code**, mentre **Github** può essere lo strumento utilizzato condividere progetti per lavoro in maniera collaborativa. Se richiesta memoria GPU per piccoli progetti accademici è disponibile **COLAB** che permette di eseguire in remoto codice offrendo anche gratuitamente utilizzo di GPU. Questi esempi mostrano una panoramica di strumenti vasta, complessa e in rapita evoluzione, con un frequente cambio di software per realizzare un programma. Un esempio d'utilizzo con gli strumenti sopra elencato potrebbe essere la realizzazione iniziale del progetto in locale utilizzando Visual Studio Code per poi riportare il tutto su GitHub. In un secondo momento il codice viene ripreso e aperto su Colab dove a sua volta il programma viene modificato ed infine rieseguito il Push sul progetto radice presente su GitHub. Ora nel 2025, la cosa che accomuna questi strumenti, è l'implementazione al loro interno di funzioni che basate sull'IA, in grado di completare il codice, suggerire correzioni e creare documentazione pertinente. Un esempio semplice ma che offre già un'idea della vastità e della potenza di queste funzioni è l'utility di **Github Copilot** 'Generate Commit Message with Copilot' che propone il testo da utilizzare come descrizione di un commit, ho provato a riscontrare quanto fosse contestualizzato e coerente con quanto aggiornato e ho ottenuto il seguente risultato:



Nel mio caso quanto proposto era corretto ed ho quindi eseguito il Commit con la descrizione proposta. Quanto è riuscito a fare Copilot è strabiliante, in pochi istanti ha analizzato il contesto ritornando come output una risposta semplice ma coerente rispetto a quanto cambiato. L'uso di questi strumenti sta rendendo il lavoro molto più dinamico e veloce, riducendo le interruzioni nel cercare soluzioni o per trovare le giuste parole per descrivere quanto fatto.



L'intelligenza artificiale sta rivoluzionando il modo in cui il software viene sviluppato, strumenti come Copilot utilizzando tutto il loro potenziale, possono creare la spina dorsale di un progetto in pochi secondi lasciando al programmatore il compito di verificare e correggere solo in parte il codice proposto. In progetti complessi questo non riduce il ruolo del programmatore, anzi lo eleva a compiti di precisione e ad alto valore aggiunto lasciando la stesura di parti del codice semplici e ripetitive al software stesso. Sapere cosa chiedere e formulare correttamente le domande al LLM è fondamentale, esplicitando nel dettaglio con parole chiave mirate come deve essere realizzato il codice. Altro compito complesso per il

programmatore è non farsi troppo ammaliare dalle soluzioni proposte perché non sempre necessarie per quanto richiesto oppure diverse da quanto già conosciuto per realizzare una determinata funzione. Questo nuovo modo di lavorare per mette di conoscere nuove soluzioni ma comporta test e tempo non sempre disponibile, il programmatore deve sempre avere il controllo del progetto accettando generazione del codice automatica solo dove consapevole di quanto proposto e del suo impatto anche in casi di revisione e manutenzione futuri. L'ultimo miglio da percorrere per sfruttare questi strumenti è la personalizzazione delle risposte del LLM, per ottenere risposte coerenti con quanto già realizzato e conosciuto, per fare questo entra in gioco la RAG.

Chapter 2

Addestrare un LLM per la Generazione del Codice

L'addestramento di LLM per la generazione di codice di programmazione richiede una serie di passaggi metodici e risorse computazionali significative. Conoscere questo processo è utile per la successiva integrazione con la RAG. La procedura si divide nelle selle seguenti fasi:

2.1 Raccolta e Preparazione dei Dati

La qualità e la quantità dei dati per l'addestramento è di primaria importanza per preparare un modello alla generazione di codice in maniera efficace. È quindi essenziale utilizzare per il training codice sorgente proveniente da molteplici fonti tra cui codice sorgente, file Readme, documentazione tecnica, commenti nel codice, pagine Wiki, API e discussioni su forum specializzati in programmazione. In rete è possibile trovare diverso materiale open source tra cui dataset già etichettati. Alcuni dataset hanno un valore altissimo, per tutelare il costo per produrli per certi dataset è previsto il diritto d'autore. I dati si dividono in due tipologie:

- **Dati Strutturati:** seguono un formato specifico e predefinito.
- **Dati non Strutturati:** non sono organizzati e sono quindi più difficili da interpretare dal modello.

La raccolta di dati va visionata con cura, se non si conosce la provenienza del codice è possibile che contenga bug o codice opsoleto che possono essere trasmessi al modello. I dati raccolti devono essere quindi puliti e pre-processati per rimuovere errori e informazioni non pertinenti, garantendo così un dataset di alta qualità per l'addestramento. Sui dataset viene utilizzato un tokenizer specializzato che riconosce costrutti di programmazione come keyword, operatori e strutture sintattiche.

2.2 Pre-Addestramento

Il pre-addestramento di un LLM da utilizzare per la generazione di codice richiede un approccio specifico. A differenza del pre-addestramento generico, utilizzando i dataset precedentemente preparati il modello impara a:

- Predire il completamento del codice
- Comprendere la struttura sintattica dei linguaggi di programmazione
- Riconoscere pattern comuni nel codice
- Identificare le relazioni tra diversi blocchi di codice

Un esempio pratico di pre-addestramento può essere implementato utilizzando la libreria transformers [WDS⁺20, FGT⁺20]:

```
1 from transformers import RobertaConfig, RobertaTokenizerFast
2
3 # Configurazione del modello per il codice
4 config = RobertaConfig(
5     vocab_size=50000, # Dimensione del vocabolario
6     max_position_embeddings=514, # Lunghezza massima sequenza
7     num_attention_heads=12, # Teste di attenzione
8     num_hidden_layers=6, # Strati nascosti
9     type_vocab_size=1 # Tipo di vocabolario
10 )
11
12 # Tokenizer specializzato per il codice
13 tokenizer = RobertaTokenizerFast.from_pretrained(
14     "microsoft/codebert-base",
15     max_length=512,
16     truncation=True,
```



```
17 padding=True
18 )
```

Durante questa fase, il modello sviluppa una comprensione profonda della sintassi e della semantica del codice, che verrà poi raffinata durante il fine-tuning per compiti specifici di generazione del codice.

2.3 Fine-Tuning

Il fine-tuning è la fase in cui il modello viene specializzato per la generazione di codice, documentazione e risposta a quesiti specifici del contesto di programmazione. Durante questa fase, il modello affina le sue capacità attraverso:

- **Dataset Specializzati:** Utilizzo di dataset contenenti:
 - Coppie di descrizioni-implementazioni
 - Documentazione tecnica e commenti
 - Esempi di bug fixing e refactoring
- **Tecniche di Apprendimento:**
 - **Apprendimento Supervisionato:** Training su coppie input-output predefinite
 - **Apprendimento per Rinforzo:** Ottimizzazione basata su feedback e metriche di qualità
 - **Few-shot Learning:** Adattamento a nuovi contesti con pochi esempi

Un esempio pratico di fine-tuning può essere implementato utilizzando la libreria transformers [WDS⁺20]:

```
1 from transformers import Trainer, TrainingArguments
2 from datasets import load_dataset
3
4 # Caricamento del dataset per il fine-tuning
5 dataset = load_dataset("code_search_net", "python")
6
7 # Configurazione del training
8 training_args = TrainingArguments(
```

```
9     output_dir="./results",
10     num_train_epochs=3,
11     per_device_train_batch_size=8,
12     per_device_eval_batch_size=8,
13     warmup_steps=500,
14     weight_decay=0.01,
15     logging_dir="./logs",
16     logging_steps=10,
17     evaluation_strategy="epoch"
18 )
19
20 # Inizializzazione del trainer
21 trainer = Trainer(
22     model=model,                # Modello pre-addestrato
23     args=training_args,        # Argomenti di training
24     train_dataset=dataset["train"],
25     eval_dataset=dataset["validation"],
26     tokenizer=tokenizer,        # Tokenizer specializzato per il codice
27 )
28
29 # Avvio del fine-tuning
30 trainer.train()
```

Durante il fine-tuning, il modello sviluppa capacità specifiche come:

- Generazione di codice a partire da descrizioni in linguaggio naturale
- Completamento intelligente del codice basato sul contesto
- Creazione di documentazione tecnica
- Identificazione e correzione di bug
- Refactoring del codice seguendo best practices

Il processo di fine-tuning richiede un attento bilanciamento tra:

- **Overfitting:** Evitare che il modello memorizzi i dati di training
- **Generalizzazione:** Mantenere la capacità di adattarsi a nuovi contesti
- **Prestazioni:** Ottimizzare la velocità e la qualità delle risposte

2.4 Pre-Addestramento vs Fine-Tuning

È importante comprendere la distinzione tra queste due fasi dell'addestramento:

2.4.1 Pre-Addestramento

Il pre-addestramento è la fase iniziale dove il modello:

- Acquisisce una comprensione **generale** del linguaggio di programmazione
- Viene addestrato su **grandi quantità** di codice sorgente generico
- Impara le strutture base e la sintassi del linguaggio
- Non è ancora specializzato per compiti specifici

2.4.2 Fine-Tuning

Il fine-tuning è invece la fase di specializzazione dove il modello:

- Si adatta a un **dominio specifico** o a compiti particolari
- Utilizza dataset più piccoli ma **mirati**
- Affina le conoscenze per generare codice per specifici casi d'uso

Analogia: Si può paragonare a:

- Pre-addestramento: Imparare la grammatica e il vocabolario di base di una lingua
- Fine-tuning: Specializzarsi nel linguaggio tecnico di un settore specifico

2.5 Architettura del Modello

Gli LLM utilizzano tipicamente architetture basate su trasformatori, che sono particolarmente efficaci nell'elaborazione di sequenze di dati, come il testo e il codice. I trasformatori utilizzano meccanismi di auto-attenzione per valutare l'importanza di diversi elementi in una sequenza, permettendo al modello di comprendere le relazioni tra parole o token. Questa capacità è fondamentale nella generazione del codice, poiché le dipendenze tra variabili e funzioni possono estendersi su ampie sezioni del codice, richiedendo al modello di considerare un ampio contesto per trovare le risposte corrette.

2.6 Valutazione e Ottimizzazione

Una volta addestrato, il modello deve essere rigorosamente valutato utilizzando metriche specifiche per la generazione di codice, come la correttezza sintattica, la funzionalità e l'efficienza del codice prodotto. I risultati della valutazione possono essere utilizzati per ulteriori ottimizzazioni, come aggiustamenti dei pesi del modello, modifiche all'architettura o includere dati di addestramento aggiuntivi per affrontare eventuali carenze.

2.6.1 Metriche di Valutazione

- **Correttezza Sintattica:** Verifica che il codice generato sia sintatticamente corretto.
- **Funzionalità:** Verifica che il codice generato realizzi la funzionalità desiderata.
- **Efficienza:** Valuta le prestazioni del codice in termini di tempo di esecuzione e utilizzo delle risorse.

2.6.2 Tecniche di Ottimizzazione

- **Aggiustamento dei Pesi:** Modifica dei pesi del modello per migliorare le prestazioni.
- **Modifiche all'Architettura:** Introduzione di nuove componenti o modifiche a quelle esistenti.
- **Integrazione di Dati Aggiuntivi:** Utilizzo di ulteriori dati di addestramento per migliorare le prestazioni.

Chapter 3

RAG

3.1 Introduzione

Il RAG **Retrieval-Augmented Generation**, (in italiano *Generazione Aumentata tramite Recupero*) è un sistema che permette di migliorare l'output di un LLM estendendo la sua conoscenza con nuove informazioni, al di fuori dai suoi dati di addestramento. Allo scopo di:

- ottenere risposte personalizzate provenienti da librerie e codice custom;
- migliorare il codice generato rendendolo più specifico al dominio riducendo le allucinazioni;
- facilitare l'assistenza da parte del modello nella fase di debugging migliorando la sua comprensione di sistemi complessi;
- supportare la creazione di documentazione aggiornata;
- permettere all'interno di un Team di migliorare la coerenza del codice scritto da diversi programmatori;
- realizzare naturalmente senza forzature, codice più moderno proponendo librerie e standard comuni.
- evitare risposte imprecise a causa della confusione terminologica, in cui diverse fonti utilizzano la stessa terminologia per parlare di cose diverse.

3.2 Funzionamento

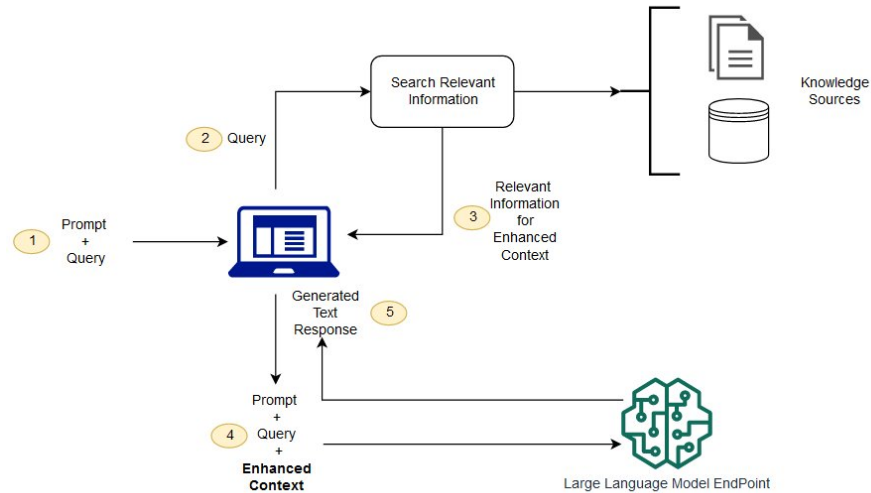


Figure 3.1: Flusso di una request ad un LLM integrato con un RAG

RAG è un sistema che integra il processo di generazione del linguaggio con un meccanismo di recupero delle informazioni. Il funzionamento si articola in diverse fasi, come illustrato in fig. 3.1.

3.2.1 Creazione degli Embedding

Il sistema RAG utilizza dati esterni al training set originale del LLM, provenienti da diverse fonti come:

- API e database interni
- Archivi documentali
- File di testo e codice

Questi dati vengono convertiti in rappresentazioni numeriche (embedding) e archiviati in un database vettoriale, creando una knowledge base accessibile dal RAG.

3.2.2 Fase 1: Function Calling

Il sistema RAG inizia con una chiamata di funzione per ricercare nei dati di embedding:

- La query dell'utente attiva una chiamata di funzione
- Il sistema cerca nei dati di embedding le informazioni pertinenti
- Se trovate, queste informazioni vengono aggiunte al prompt

3.2.3 Fase 2: Recupero delle Informazioni

Quando l'utente sottopone una query:

- La domanda viene convertita in un vettore
- Il sistema cerca nel database vettoriale le informazioni più pertinenti
- Viene calcolata la rilevanza attraverso calcoli matematici vettoriali

3.2.4 Fase 3: Aumento del Prompt

Il sistema RAG arricchisce il prompt dell'utente:

- Aggiunge le informazioni recuperate al contesto
- Utilizza tecniche di prompt engineering per ottimizzare la comunicazione con il LLM
- Fornisce al modello un contesto arricchito per generare risposte più accurate

3.3 Gestione dell'Aggiornamento dei Dati

Per mantenere l'efficacia del sistema nel tempo:

- Ricalcolo degli embedding per i nuovi dati
- Possibilità di aggiornamenti in tempo reale o batch

Questo approccio permette di superare le limitazioni dei LLM, fornendo risposte più accurate e contestualizzate grazie all'integrazione di conoscenze esterne aggiornate.

Chapter 4

Caso Studio: Implementazione di un Sistema RAG per lo Sviluppo del Software

4.1 Obiettivi del Caso Studio

Questo caso studio si propone di:

- Dimostrare l'applicazione pratica dei concetti teorici presentati nei capitoli precedenti
- **Scenario:** Creare un RAG di supporto allo sviluppo software in una software-house
- **Problematica:** Necessità di generare codice coerente con gli standard e le pratiche aziendali
- **Soluzione Proposta:** Sistema RAG integrato con LLM per la generazione di codice contestualizzato

4.2 Architettura del Sistema

(da completare) [Dati Aziendali] → [Chunking] → [Embedding] → [Vector DB] → [Retrieval] → [LLM] → [Risposta contestuale]

4.3 Software Utilizzati

4.3.1 Ollama

Ollama [Oll24] è un software che permette di utilizzare in locale LLM, facilitando la creazione, l'addestramento e l'integrazione di modelli di machine learning nelle applicazioni senza dover dipendere da servizi cloud esterni.

4.3.2 LLM

Per rendere più completa la ricerca sono stati utilizzati due modelli di linguaggio per il caso studio:

Llama 3.2

Llama 3.2 3B [AI24], un modello di linguaggio open source sviluppato da Meta AI. Il modello, con 3 miliardi di parametri, è ottimizzato per compiti di dialogo multilingue e si distingue per le sue capacità di recupero e sintesi delle informazioni. La scelta è ricaduta su questa versione per il suo equilibrio tra prestazioni e requisiti computazionali che permettono il suo utilizzo senza hardware troppo potente.

qwen2.5-coder:3b

qwen2.5-coder [Tea24] è stato sviluppato da Qwen AI ed è anch'esso open source, specializzato nella generazione di codice e documentazione tecnica. Con 3 miliardi di parametri, il modello è stato addestrato su un ampio dataset di codice sorgente e documentazione tecnica, permettendo di generare codice coerente e ben strutturato. La scelta di questo modello è stata dettata, a differenza di llama3.2, dalla sua specializzazione nella programmazione e dalla sua capacità di generare codice di alta qualità.

4.3.3 BGE-M3

BGE-M3 [BAA24] è un database vettoriale open source per la gestione di dati strutturati e non strutturati multilingue, sviluppato da BigGraph Engine.

4.4 Dataset

Il dataset utilizzato per il caso studio è stato creato a partire da documenti tecnici e codice sorgente di progetti interni all'azienda.

4.5 Implementazione

4.5.1 Modello embedding

Utilizzando BGE-M3 il contenuto del dataset viene suddiviso in chunk e tramite il modello embedding è convertito in rappresentazioni vettoriali numeriche. esempio

```
1 model = BGEM3FlagModel('BAAI/bge-m3', use_fp16=True)
2 frase = "<cis-ui:button skin='red' text='Bottone rosso' />"
3 embeddings = model.encode(frase)['dense_vecs']
4 print (embeddings)
```

```
[-0.03712651  0.02731708 -0.03263818 ...  0.02965933 -0.01724784  0.00510089]
```

l'output è un vettore numerico con 1024 elementi.

4.5.2 Creazione Rag

4.5.3 AI Agent

E' un componente che può chiamare delle funzioni (function calling) nel mio caso search il LLM capisce se deve chiamare o meno il RAG.

function calling

Il function calling è un meccanismo che permette di chiamare delle funzioni definite in un file esterno, in questo caso il RAG.

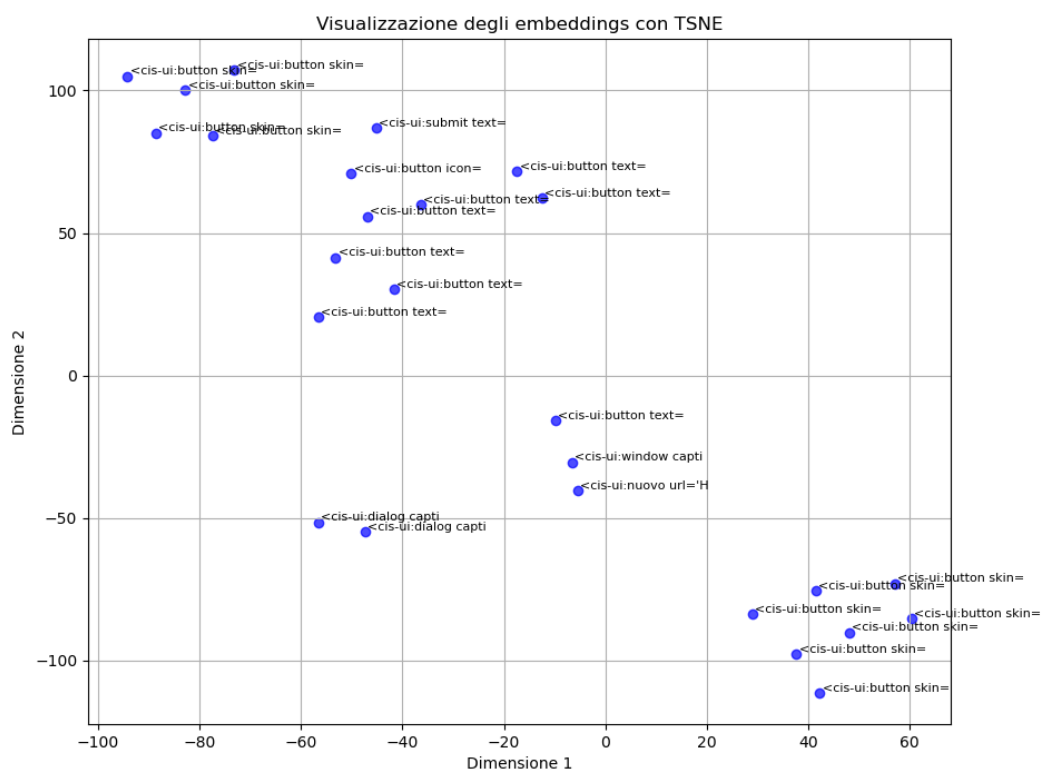


Figure 4.1: Visualizzazione degli embeddings con TSNE

Chapter 5

Conclusioni

5.1 Risultati Ottenuti

5.2 Impatto sullo Sviluppo Software

L'integrazione di strumenti basati su AI nel processo di sviluppo software sta rivoluzionando il settore. Durante il periodo di sviluppo di questa tesi (Ottobre 2024 - Gennaio 2025), abbiamo osservato:

- Rapida evoluzione degli strumenti di AI per lo sviluppo software
- Crescente disponibilità di soluzioni open source
- Miglioramento continuo nelle capacità di generazione e comprensione del codice

5.3 Sfide e Prospettive Future

Bibliography

- [AI24] Meta AI. Llama-3.2-3b: Open foundation and fine-tuned chat models, 2024.
- [BAA24] BAAI. Bge-m3: A multi-modal model understanding images and text, 2024. HuggingFace model repository for BGE-M3, a multi-modal model for image and text understanding.
- [Doc24] Huggingface Docs. Lora, dec 2024.
- [Doc25] GitHub Docs. Asking github copilot questions in your ide, jan 2025.
- [FGT⁺20] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155*, 2020.
- [Git24] GitHub. Github copilot is more than a tool, it’s an ally, dec 2024.
- [Met24] Meta. Llama-3.3-70b-instruct, dec 2024.
- [Oll24] Ollama. Ollama documentation, 2024. GitHub repository.
- [SBO23] Ahmed R. Sadik, Sebastian Brulin, and Markus Olhofer. Coding by design: Gpt-4 empowers agile model driven development, 2023.
- [Tea24] Qwen Team. Qwen2.5-coder-3b: A code-specialized language model, 2024.

- [WDS⁺20] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45. Association for Computational Linguistics, 2020.