

Building a natural language code search tool with vector embeddings and LLM

Dima Yanovsky

Context

Initial Frustrations

- ✗ Can't follow all the functions in a chain
- ✗ Can't read codebase like a book
- ✗ Limited access to the code authors

You need to understand

- ✓ How processes are implemented
- ✓ Workflows and business logic
- ✓ Codebase on a conceptual level

How to learn fast?

- ! Ask senior SWE a lot of questions
- ✗ Seniors have limited time
- 📉 Learning velocity decreases

Solving the problem

- 🤖 Build code knowledge tool
- ? Receives a question in natural language
- 👍 Answers the question and provides relevant code from the codebase

Background

 **Code elements**

Functions, interfaces, structs and other parts of the code that serve a purpose (aka chunks)

 **Code parsing and splitting**

A technique to split the codebase into code elements.

 **Vector Embeddings**

Representation of some object in a multi-dimensional space. Neural nets capture the semantic meaning of the object and express it as a vector. Objects with similar meanings are close to each other

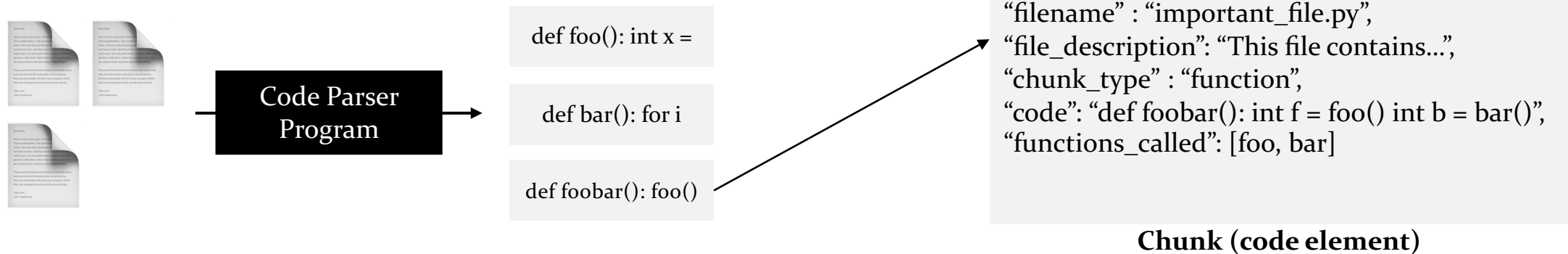
 **Vector Storage Database**

Database that stores vector embeddings and is used to find closest vectors to an input vector.

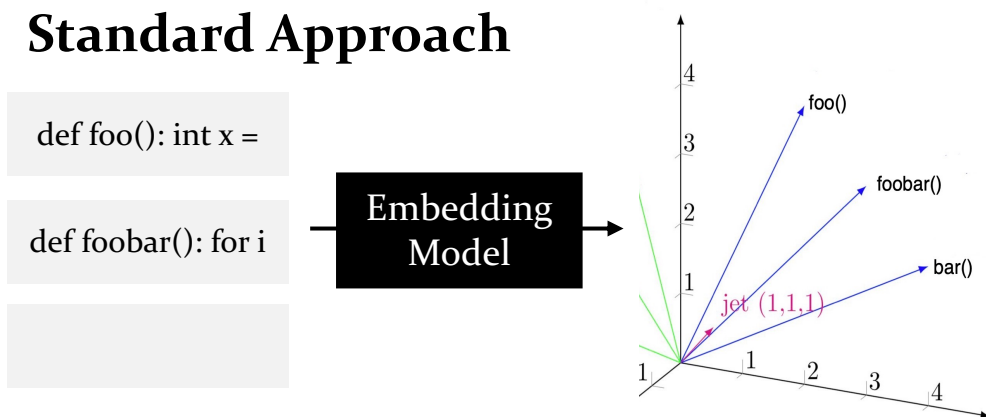
 **Large Language Model**

Model that can understand and generate human-like text based on the vast amount of data it's been trained on

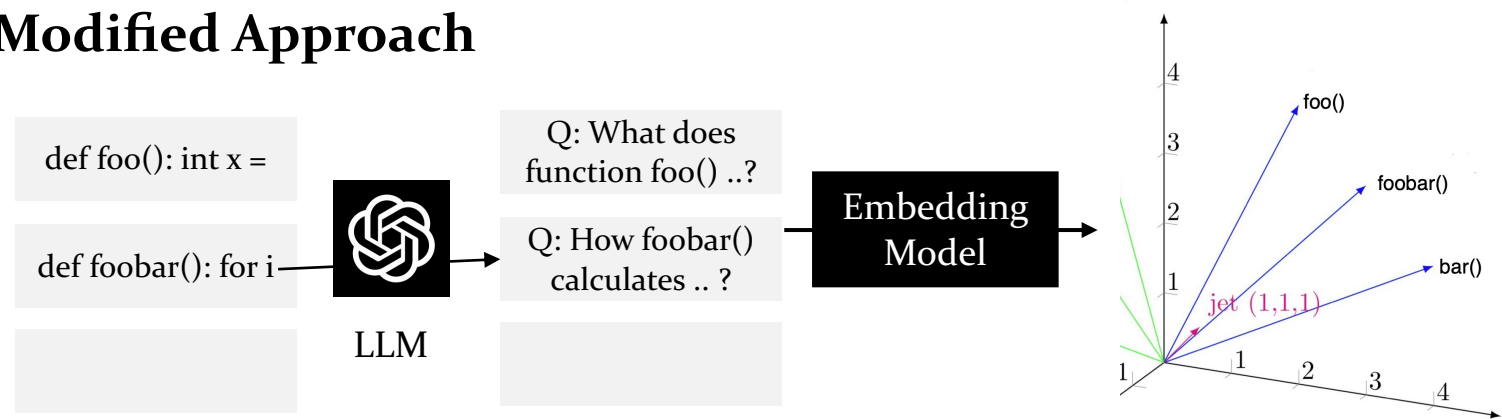
Parsing and Splitting



Standard Approach



Modified Approach



User: "Explain foo() function"

