



# Code Compass

Final Review

# Agile Team



**Product Manager + Data Engineer**  
Gabriel de Olaguibel



**Data Scientist + ML Ops**  
Miranda Drummond



**Data Engineer + Data Scientist + ML Engineer**  
Luca Cuneo



**Data Scientist + MLOps**  
Ketí Sulamanidze



**Data Engineer + Data Scientist**  
Maud Hovland



**Data Scientist + ML Engineer**  
Félix Gómez-Guillamón

## Code Compass

*Navigating Developers to  
Personalized GitHub exploration  
and Insights*

## Agenda

- ❖ Product Goals
- ❖ Demo
- ❖ Data Engineering
- ❖ Recommender EDA
- ❖ Recommender Model
- ❖ Chatbot
- ❖ Tests/MLOps



# Product

# Vision

*Enhance the GitHub experience for developers and learners by providing personalized experience.*

## Solution

*A combination of a **recommendation system** and a **chatbot**:*

- **Personalized Recommendations:** Utilizing user-specific, Code Compass will recommend relevant GitHub projects.
- **Interactive Learning and Exploration:** The chatbot will serve as an interactive guide, helping users to delve deeper into repositories, understand technology applications, and navigate more effectively.



# Demo

# Chatbot

## Github Repository Chatbot

Exit Chat



Tell me about the following repository: <https://github.com/recommenders-team/recommenders>



The repository "recommenders" contains a wide range of files and directories. Some key files and directories include:

- Main Files:
  1. README.md
  2. LICENSE
  3. CODE\_OF\_CONDUCT.md
  4. CONTRIBUTING.md
  5. SETUP.md
  6. SECURITY.md
  7. AUTHORS.md
  8. MANIFEST.in
  9. NEWS.md



# Recommender

## GitHub Repo Recommendation System

Enter the target user's username:

miguelgferro

Get Recommendations

### Recommendations

1 - Repo ID: 537923335.0, Owner: khizarsiddiqui, [Link](#)

Description: A Python program that generates ASCII art from graphical images using Pillow, a fork of the Python Imaging Library (PIL) and numpy.

2 - Repo ID: 403642804.0, Owner: shazron, [Link](#)

Description: Adobe I/O Photoshop API SDK

3 - Repo ID: 251512047.0, Owner: Ewenwan, [Link](#)

Description: MindSpore is a new open source deep learning training/inference framework that could be used for mobile, edge and cloud scenarios.

4 - Repo ID: 317381518.0, Owner: noahgift, [Link](#)

Description: Some recipes for doing with serverless technologies

5 - Repo ID: 604684135.0, Owner: giswqs, [Link](#)

Description: A Python package demo for interactive mapping

6 - Repo ID: 10714550.0, Owner: petrounias, [Link](#)

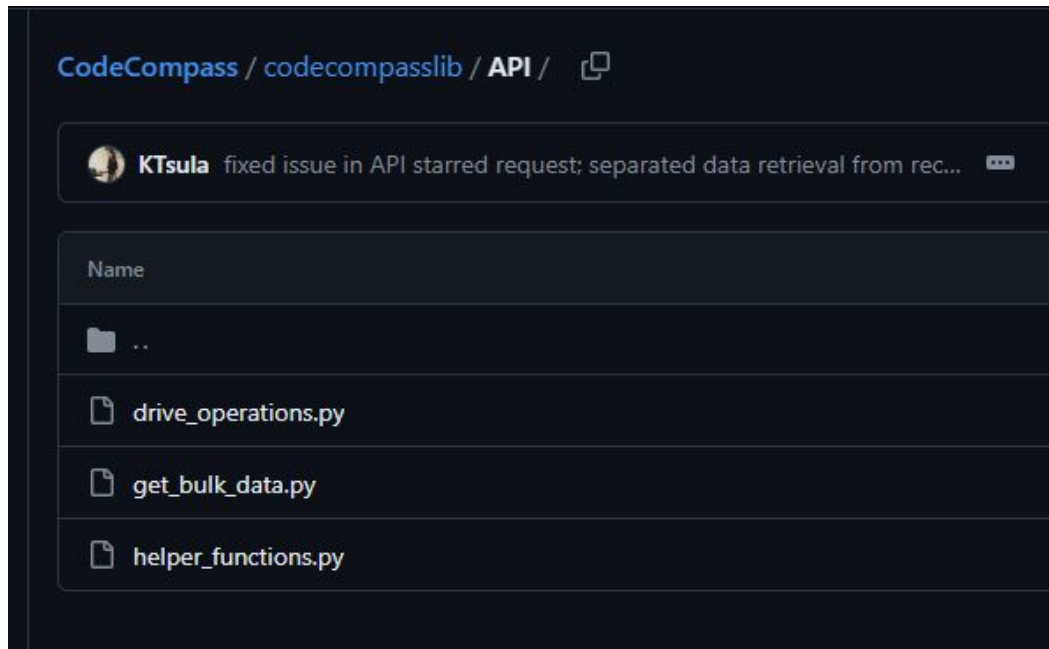
Description: Programmatic building of JSON schemas (document and field mappings) with validation.





# Data Engineering

# Data Engineering



The screenshot shows the CodeCompass file explorer for the 'data' directory. It displays a table of files with columns for Name, Last modified, and File size. The files listed are various embedded datasets in CSV format, including df\_with\_embeddings.csv, Embedded\_C\_dataset.csv, Embedded\_C#\_dataset.csv, Embedded\_C++\_dataset.csv, Embedded\_Java\_dataset.csv, Embedded\_JavaScript\_dataset.csv, Embedded\_Jupyter Notebook\_dataset.csv, Embedded\_PHP\_dataset.csv, Embedded\_Python\_dataset.csv, Embedded\_Ruby\_dataset.csv, Embedded\_Shell\_dataset.csv, test.csv, and uploaded\_dataset.csv.

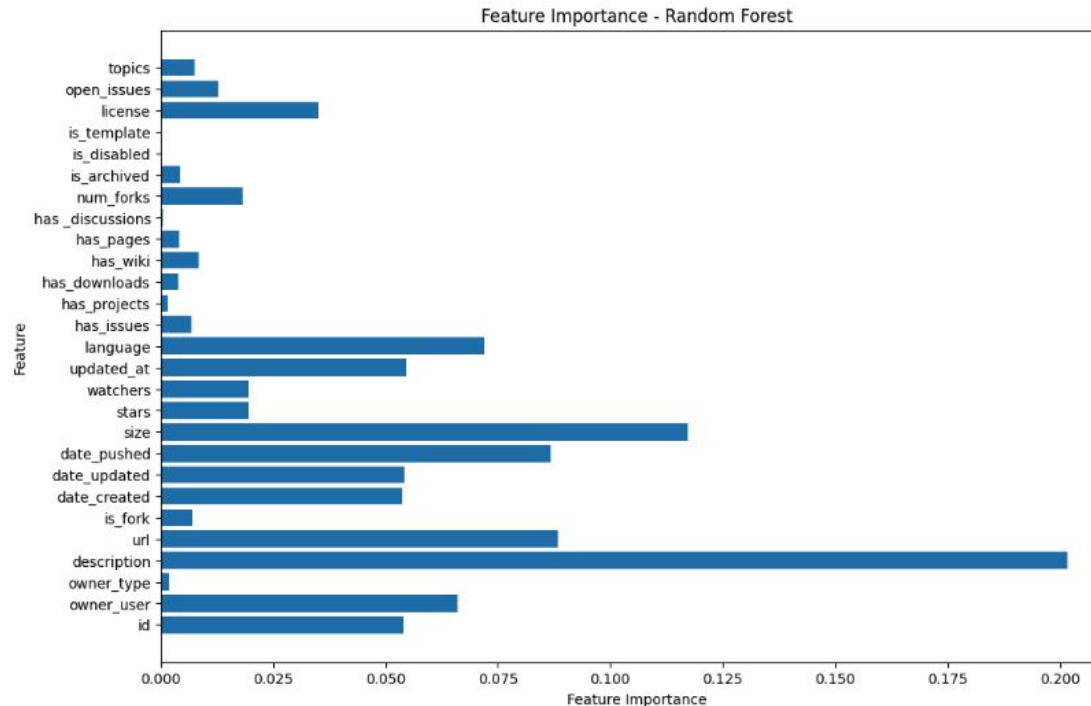
Name	Last modified	File size
df_with_embeddings.csv	27 Mar 2024 mhovland.ieu2...	190 MB
Embedded_C_dataset.csv	29 Mar 2024 mhovland.ieu2...	114.9 MB
Embedded_C#_dataset.csv	29 Mar 2024 mhovland.ieu2...	101.1 MB
Embedded_C++_dataset.csv	29 Mar 2024 mhovland.ieu2...	143.2 MB
Embedded_Java_dataset.csv	29 Mar 2024 mhovland.ieu2...	243.9 MB
Embedded_JavaScript_dataset.csv	29 Mar 2024 mhovland.ieu2...	653.4 MB
Embedded_Jupyter Notebook_dataset.csv	29 Mar 2024 mhovland.ieu2...	77.8 MB
Embedded_PHP_dataset.csv	29 Mar 2024 mhovland.ieu2...	109.6 MB
Embedded_Python_dataset.csv	29 Mar 2024 mhovland.ieu2...	397.6 MB
Embedded_Ruby_dataset.csv	29 Mar 2024 mhovland.ieu2...	134.3 MB
Embedded_Shell_dataset.csv	29 Mar 2024 mhovland.ieu2...	96.1 MB
test.csv	03:44 mhovland.ieu2021@s...	16 bytes
uploaded_dataset.csv	24 Mar 2024 mhovland.ieu2...	837.5 MB

- Utilizing Github API to retrieve repository and profile information
- All code files contained within codecompass/API folder
  - *get\_bulk\_data.py*:
    - Contains the logic for building our dataset (this has been improved since last review)
  - *helper\_functions.py*:
    - Contains functions to load secrets, or extract useful information from repositories
  - *drive\_operations.py*
    - Contains the functions to upload our data to our google drive.



# Recommender EDA

# EDA / Feature Engineering



## EDA and Feature Engineering:

- **Data Preparation:** Loads and cleans data, removing unnecessary columns and handling missing values.
- **Visualization and Correlation:** Utilizes plots to visualize data and performs correlation analysis to identify key relationships.
- **Encoding and Vectorization:** Applies one-hot encoding to categorical data and uses `TfidfVectorizer` for textual data to prepare it for machine learning.
- **Normalization:** Scales numerical features to a uniform range using techniques like `MinMaxScaler`.
- **Data Preparation for ML:** Converts processed data into tensors, making it suitable for training with PyTorch models.



# Recommender Model



# Model Embeddings

## OpenAI

```
def generate_openAI_embeddings(strings_to_embed, client):
    """
    Generates OpenAI embeddings for the given strings using the specified OpenAI client.
    Args:
        strings_to_embed (list): A list of strings to generate embeddings for.
        client: The OpenAI client object used to make API requests.
    Returns:
        dict: A dictionary containing the embeddings generated by OpenAI.
    Raises:
        OpenAIException: If there is an error while making the API request.
    """
    response = client.embeddings.create(
        input=strings_to_embed,
        model="text-embedding-3-large", # You can choose the model you prefer
        dimensions=256 # You can choose the number of dimensions you prefer
    )
    return response
```

## Word2Vec model trained in SWE domain

```
word_vect = KeyedVectors.load_word2vec_format("./codecompasslib/PretrainedModels/SO_vectors_200.bin", binary=True)
return word_vect
```

## Codellama locally

```
def generate_codellama_embeddings(text):
    """
    Generates embeddings for the given text using the OllamaEmbeddings model.
    Args:
        text (str): The input text for which embeddings need to be generated.
    Returns:
        query_result (list): A list of embeddings for the input text.
    """
    embeddings_model = OllamaEmbeddings(model='codellama:7b', device='gpu') #
    query_result = embeddings_model.embed_query(text)
    return query_result
```

## Sentence transformer | roberta-base

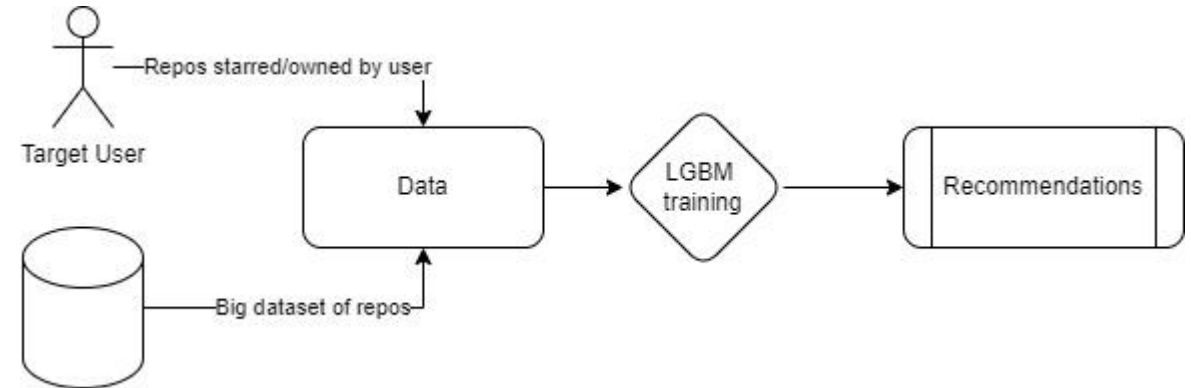
```
def generate_sentence_transformer_embeddings(text):
    """
    Generates Sentence Transformer embeddings for the given text.
    Parameters:
        text (str): The input text to generate embeddings for.
    Returns:
        embedding (numpy.ndarray): The Sentence Transformer embeddings for the given text.
    """
    # Load a pre-trained Sentence Transformer model
    model_name = 'stsb-roberta-base'
    model = SentenceTransformer(model_name)
    embedding = model.encode(text)
    return embedding
```

- Word2Vec - the fastest | not advanced enough.
- Roberta-base - Slightly more advanced sentence transformer | quite slow and space consuming
- Codellama - Much more advanced and applicable | very slow and space consuming
- OpenAI text embeddings 3-large - Advanced, applicable, adjustable is size | not done locally

# Final Model

## LightGBM- why?

- Efficiency and Scalability
- Gradient Boosting Strengths
- Categorical Feature Handling



## Models Considered

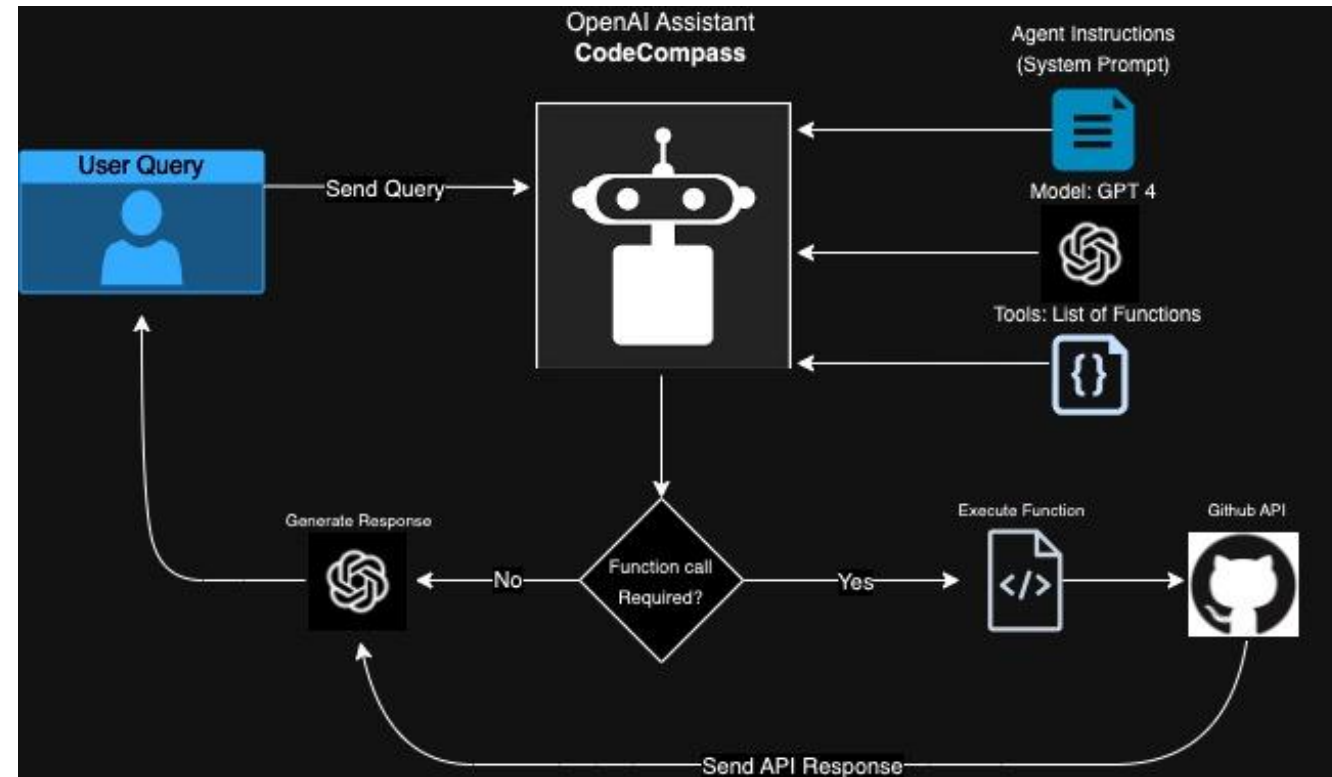
- Cosine similarity - assigns similarity scores to user language choices
- Differentiating repos (couldn't find related papers suggested by others)
- LSTUR (no timestamps for user interaction data) -> could be considered for future.
- RAG - could be used to fill in missing descriptions in repos





# Chatbot Architecture

# Chatbot Architecture



## Based on OpenAI Assistants (Beta)

- Performs Function calling to External API's when Repository Information is needed
- Function Calling Capabilities:
  - Fetch repo file tree
  - Fetch Contents of specific files
  - Fetch branches in the repo
  - Fetch commit history
  - Search by keyword
  - Search commits by keyword
  - Find repositories by keyword

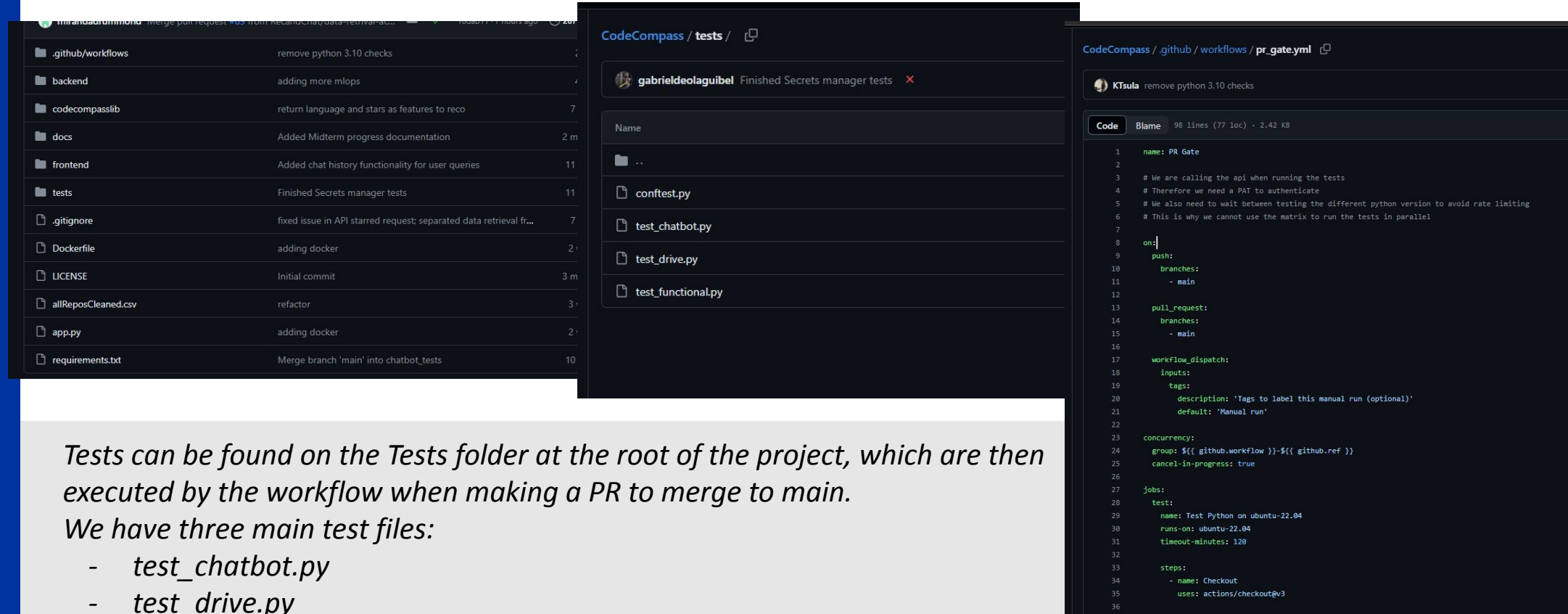
## Development Iterations

- Vanilla OpenAI API + "openplugin"
- OpenAI Chat Completion
- OpenAI Assistant (Final)



# MLOps/Tests

# Testing/Workflows



The image displays three screenshots from a GitHub repository, illustrating the project structure and testing workflow.

**Left Screenshot:** A file browser view of the repository root. The 'tests' folder is highlighted, showing its commit history. Other files and folders visible include .github/workflows, backend, codecompasslib, docs, frontend, .gitignore, Dockerfile, LICENSE, allReposCleaned.csv, app.py, and requirements.txt.

**Middle Screenshot:** A view of the 'tests' folder. It lists the following files: confest.py, test\_chatbot.py, test\_drive.py, and test\_functional.py. The folder is named 'gabrieldeolaguibel' and shows 'Finished Secrets manager tests'.

**Right Screenshot:** A view of the 'pr\_gate.yml' workflow file. The file is 98 lines (77 loc) and 2.42 KB. It defines a 'PR Gate' workflow that runs on the 'main' branch. The workflow includes a 'push' trigger, a 'pull\_request' trigger, and a 'workflow\_dispatch' trigger. It uses a 'concurrency' group to ensure only one run is active at a time. The workflow includes a 'test' job that runs on 'ubuntu-22.04' with a 120-minute timeout. The job steps include a 'Checkout' step using 'actions/checkout@v3'.

*Tests can be found on the Tests folder at the root of the project, which are then executed by the workflow when making a PR to merge to main.*

*We have three main test files:*

- *test\_chatbot.py*
- *test\_drive.py*
- *test\_functional.py*

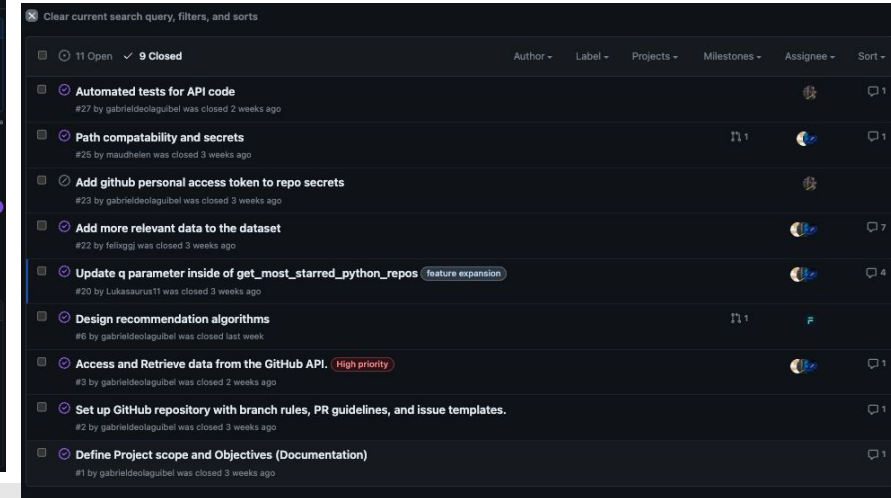
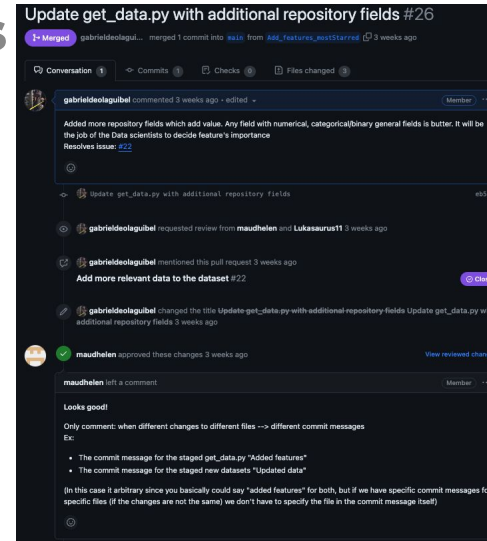
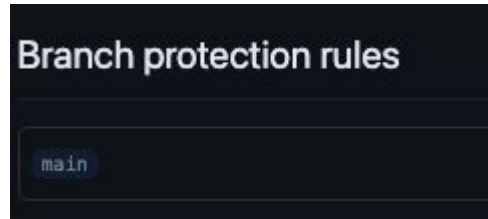
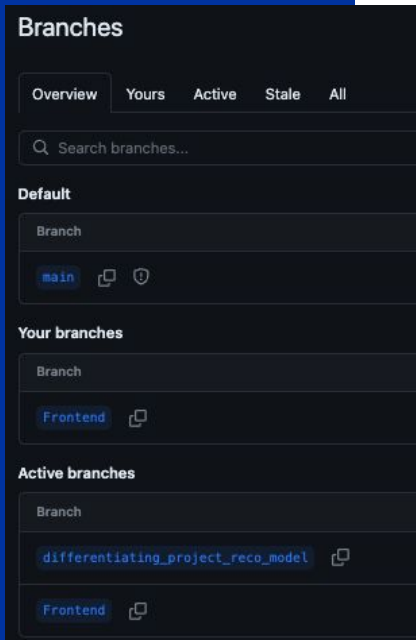


# Development Methodology



# Development Methodology

## Good Coding Practices



*Required development in short-lived branches:*

- *No Member is allowed to push to main*
- *Branches are deleted after merges to main*

*Add code via PRs:*

- *Requires another member to review and comment.*
- *Requires automated workflow tests to pass*

*Issues:*

- *Issues are part of the github project backlog*
- *Members submit issue with what need to be fixed/done and get feedback from others on the ideas*
- *PRs directly address issues*

*Libraries and Repository design:*

- *Notebooks and scripts call on codecompasslib library or external libraries*
- *Code is neatly organized in its respective folders following “Project template” design by Miguel.*

# Further Improvements

## To be Implemented by Sunday April 7

- **Full Code Documentation:** Explanation/Math behind the algos, Demo videos/gifs
- **Finalize Library (codecompasslib):** Perform release and move a few files around/cleanup
- **Further Tests:** Improve code coverage for (Reco Model functions)

## Longer term Improvements

- **Implement different Reco Models:**
  - Collaborative filtering based models;
  - Time considerate models (ex. LSTUR to capture short term and long term representations)
- **Improve chabot functionality:** Implement more repository data retrieval
- **Cloud Deployment:** Host fontends on Cloud
- **Improve Datastore:** Try Vector Databases to store embeddings for faster Data Retrieval
- **Should work for any user:** Implement function that embeds and adds new user's repos to dataset automatically





# Feedback / Questions