# Code Compass

Midterm Review

**ie**
UNIVERSITY

# 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**
Keti Sulamanidze

**Data Engineer + Data Scientist**
Maud Hovland

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

**ie**
UNIVERSITY

**Agenda**

**Code Compass**

*Navigating Developers to Personalized GitHub exploration and Insights*

❖ Product Goals

❖ Demonstrate Progress

❖ Development Methodology

❖ Future Roadmap

❖ Collect Feedback

**ie**
**UNIVERSITY**

# Product Goals

UNIVERSITY

# Vision

*Enhance the GitHub experience for developers and learners by providing personalized guidance and resources. This system aims to simplify the process of finding relevant projects and learning materials, making the path to knowledge and collaboration more accessible and efficient.*

# Solution

*Addressing these challenges, Code Compass proposes a novel solution: a combination of a recommendation system and a chatbot interface. This dual approach aims to:*

- ***Curate Personalized Recommendations***: *Utilizing user-specific data such as project involvement, coding languages, and collaboration networks, Code Compass will recommend relevant GitHub projects, technologies, and connections.*
- ***Interactive Learning and Exploration***: *The integrated chatbot will serve as an interactive guide, helping users to delve deeper into recommendations, understand technology applications, and navigate GitHub more effectively.*

ie
UNIVERSITY

**Backlog**

## Todo 7
This item hasn't been started

**⊙ Github-Recommendation-System #7**
Develop chatbot framework and integrate NLP for user interaction.
`Phase 2` `ML Engineer`

**⊙ Github-Recommendation-System #8**
Implement the recommendation logic in the chatbot system.
`Phase 2` `ML Engineer`

**⊙ Github-Recommendation-System #10**
Develop a deployment strategy for the recommendation system and chatbot. (Documentation)
`Phase 3` `ML Engineer`

**⊙ Github-Recommendation-System #11**
Set up a CI/CD pipeline for automated deployment.
`Phase 3` `MLOps`

**⊙ Github-Recommendation-System #13**
Implement MLOps practices for monitoring and maintaining the system.
`Phase 3` `MLOps`

**⊙ Github-Recommendation-System #14**
Prepare final project presentation and documentation.
`Phase 4` `Project Manager (PM)`

**⊙ Github-Recommendation-System #15**
Evaluate the project against the criteria and make necessary adjustments.
`Phase 4` `Project Manager (PM)`

## In Progress 3
This is actively being worked on

**⊙ Github-Recommendation-System #4**
Build a data pipeline for continuous data ingestion and preprocessing.
`Phase 1` `Data Engineer`

**⊙ Github-Recommendation-System #9**
Design and develop user front end
`Phase 2` `ML Engineer`

**⊙ Github-Recommendation-System #12**
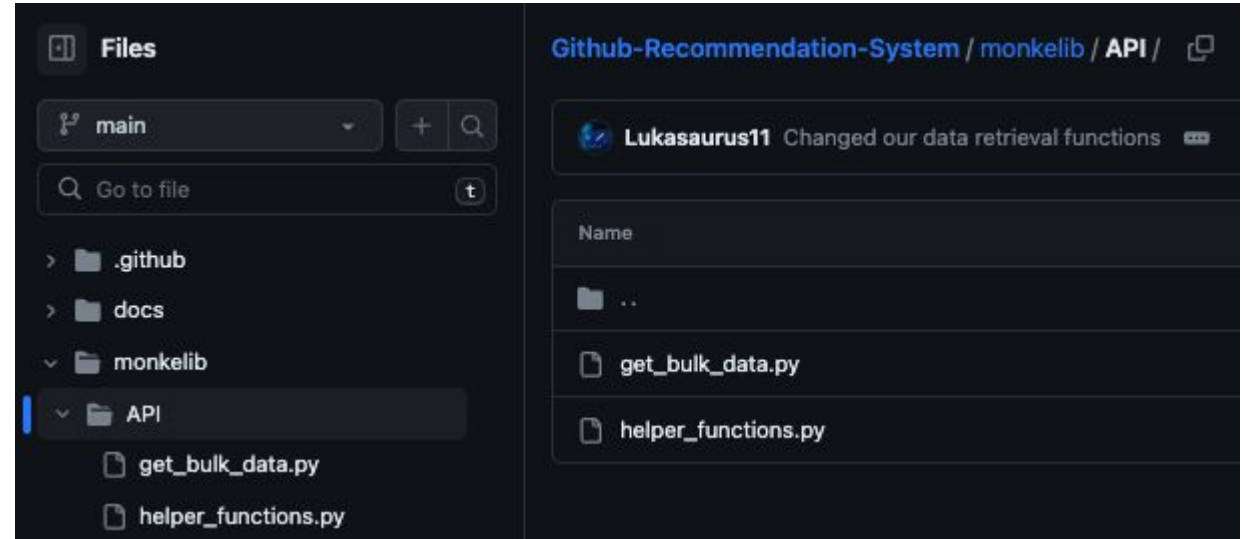Implement testing frameworks (unit tests, integration tests).
`Phase 3` `MLOps`

## Done 8
This has been completed

**⊘ Github-Recommendation-System #1**
Define Project scope and Objectives (Documentation)
`Phase 1` `Project Manager (PM)`

**⊘ Github-Recommendation-System #2**
Set up GitHub repository with branch rules, PR guidelines, and issue templates.
`Phase 1` `Project Manager (PM)`

**⊘ Github-Recommendation-System #6**
Design recommendation algorithms
`Phase 2` `Data Scientist`

**⊘ Github-Recommendation-System #27**
Automated tests for API code

**⊘ Github-Recommendation-System #25**
Path compatability and secrets

**⊘ Github-Recommendation-System #3**
Access and Retrieve data from the GitHub API.
`Phase 1` `Data Engineer`

**Github-Recommendation-System #34**
Data Modelling for Repository Recommendation System

**⊙ Github-Recommendation-System #5**
Perform Exploratory Data Analysis (EDA) to understand the data.
`Phase 1` `Data Scientist`
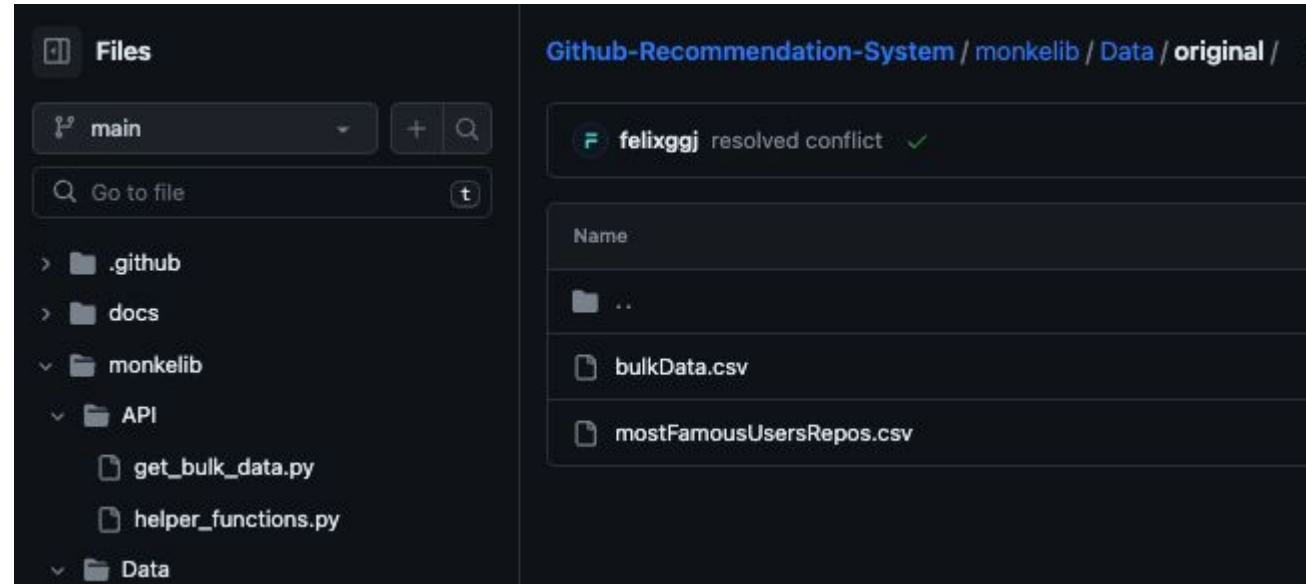
ie UNIVERSITY

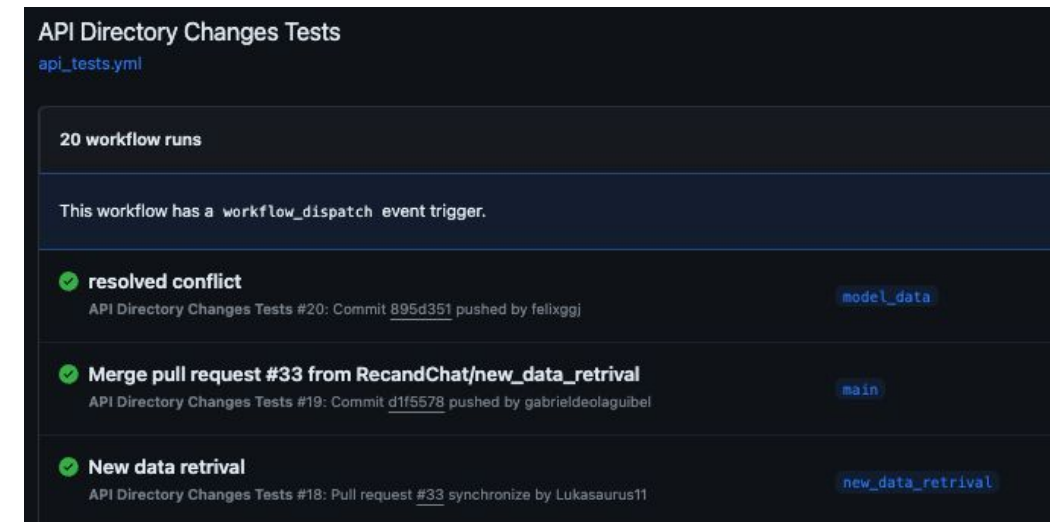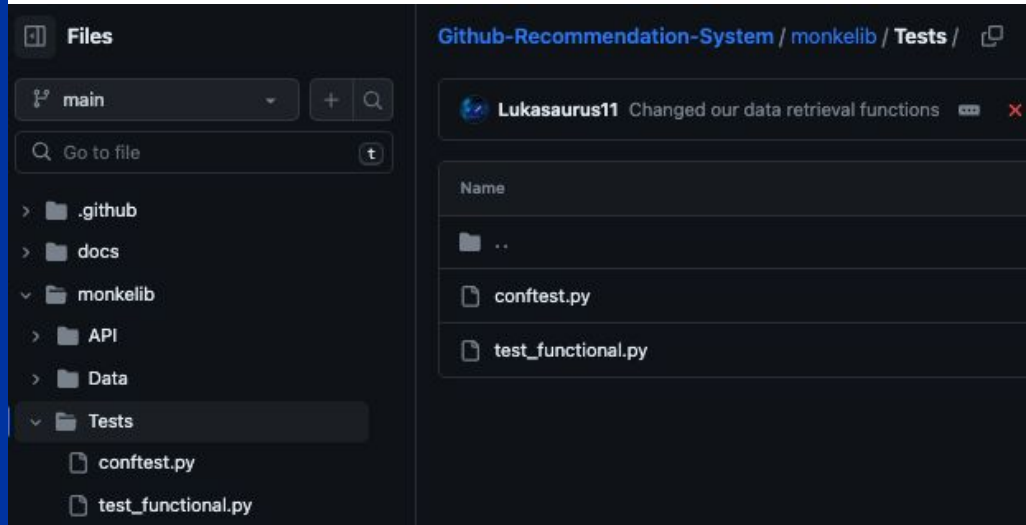# Progress

# Data Engineering / Prep



- *Utilizing Github API to retrieve repository and profile information*
- *All code files contained within monkelib/API folder*
  - *get_bulk_data.py:*
    - *Loads Github API token*
    - *Extracts Repository data*
    - *Saves data to CSV*
  - *Helper_functions.py: Contains utility functiontions for tasks like loading the authentication token, saving data, and structuring repository information.*

# Data Engineering / Prep



*Outputting results of data extraction into /Data/original folder for the Data Scientists to use for modeling*

# Data Engineering Tests

**Tests (monkelib/Tests)**: *Ensures the reliability and functionality of data retrieval from the GitHub API.*
- *conftest.py*
  - *defines fixtures for standardized test data.*
  - *Provides sample users and configurations for other test scripts.*
- *Test_functional.py*
  - *Contains functional tests for data retrieval functions.*
  - *Validates user, follower, and repository data fetching.*
  - *Tests both success and failure scenarios to ensure robustness.*

**Workflow Triggers  (.github/workflows/api_tests.yml)**: *Automated testing in response to changes within the project, specifically for code in the monkelib/API directory.*
- *On push, PR, Manual*

# EDA / Feature Engineering



```python
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
import numpy as np

# Handling missing values in text columns
df['description'].fillna('', inplace=True)
df['name'].fillna('', inplace=True)
df['language'].fillna('', inplace=True)

# Concatenating the text columns for vectorization
text_data = df['name'] + " " + df['description'] + " " + df['language']

# Vectorizing the text data
tfidf_vectorizer = TfidfVectorizer()
tfidf_matrix = tfidf_vectorizer.fit_transform(text_data)

# Calculating cosine similarity
cosine_sim = cosine_similarity(tfidf_matrix)

# Since the goal is to give a score to each repo, we can average the cosine similarities for each repo
# This gives a single score representing how similar each repo's text data is to all other repos
similarity_scores = np.mean(cosine_sim, axis=1)

# Adding the new column to the dataset
df['cosine_similarity_score'] = similarity_scores

# Displaying the updated dataset with the new column
df[['name', 'description', 'language', 'cosine_similarity_score']].head(100)
```

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.

# Recommendation Models



## Models

*cosine_modelling_repos.ipynb:*
- *Uses cosine similarity to recommend repositories based on language preferences.*
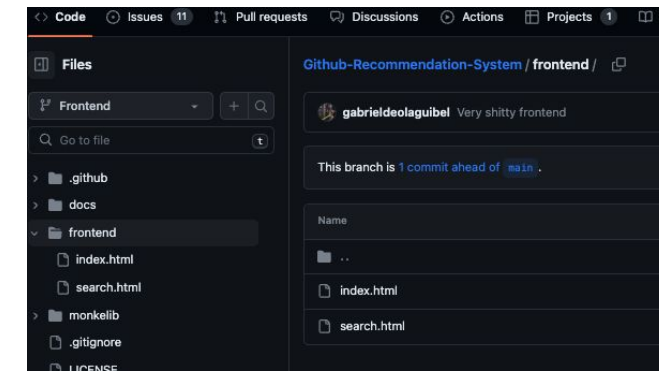- *Prepares text data for similarity scoring and vectorization.*

*modelling_repos.ipynb:*
- *Explores combining repository data with user interactions for recommendations.*
- *Implements feature engineering, including encoding and vectorization, for model input preparation.*

*modelling_differentiating_repos.ipynb*
- *Utilized domain specific embedding method to vectorize users.*
- *Finds the most differentiating project among the most similar users of the target user.*

# Simple Front end

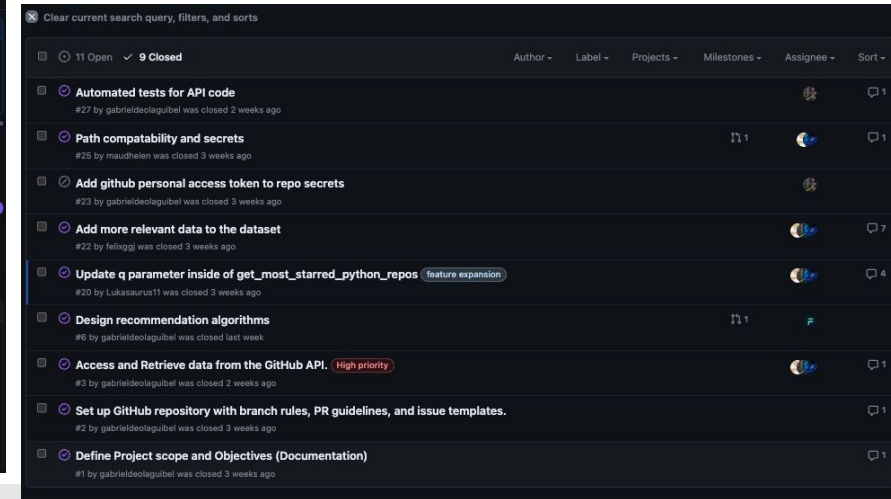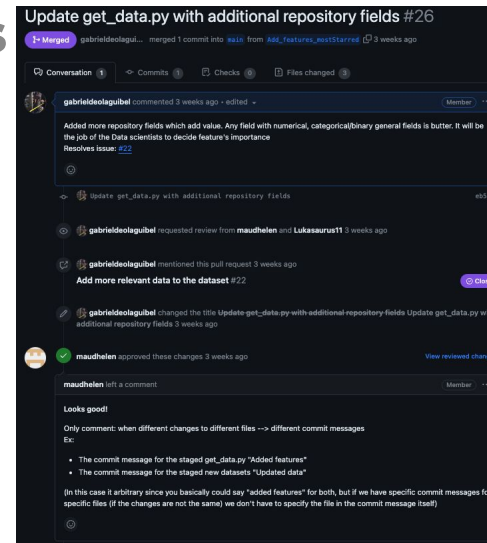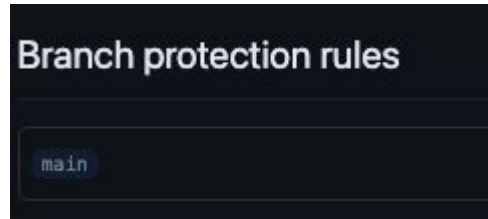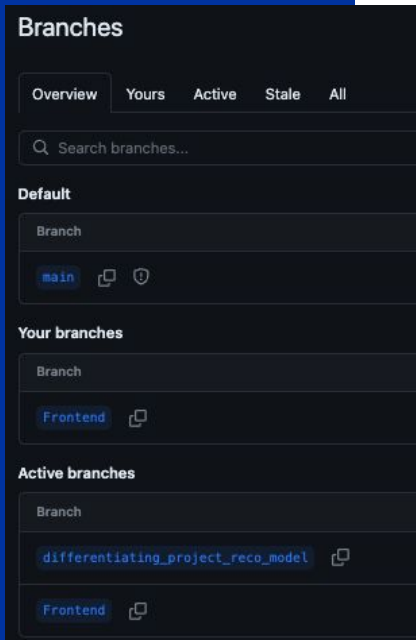**Under development within 'frontend' Branch**
- *Users will provide their github Username and receive a list of the most recommended repositories they may be interested in based on their profile*
- *Requires flask connection to the model*

# Development Methodology

UNIVERSITY

# 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*
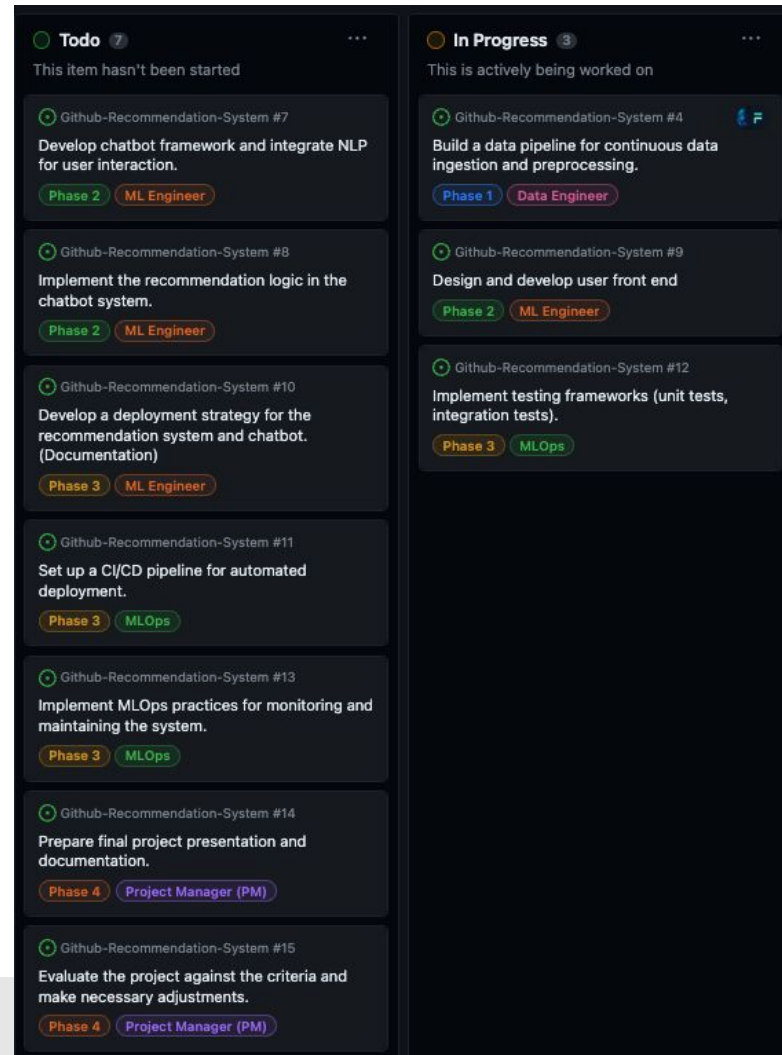- *PRs directly address issues*

*Libraries and Repository design:*
- *Notebooks and scripts call on monkelib library or external libraries*
- *Code is neatly organized in its respective folders following "Project template" design by Miguel.*

IE UNIVERSITY

# Future Roadmap

UNIVERSITY

# Roadmap



*Future progress focus:*
*Deployment to frontend*
*MLOps*
*Chatbot*

# Feedback / Questions

UNIVERSITY