

## Introduction

Our project aims to address the complex challenge of developing an effective recommender system based on graph embeddings. In the era of information overload, users often struggle to find relevant content tailored to their preferences. The goal of our research is to leverage advanced graph neural network techniques, natural language processing (NLP), and machine learning to improve recommendation accuracy and provide users with more personalized and context-aware suggestions. This involves navigating the challenges of sparse and dynamic graph structures, as well as incorporating additional node attributes for a richer understanding of the relationships within the graph. In parallel, proactive risk management aligned with the principles of PMBOK 6th edition will be integrated throughout the project lifecycle, ensuring a comprehensive and anticipatory approach.

## Data Comprehension

Data comprehension is a critical aspect of our project, involving a multifaceted approach to extract, process, and understand information from the PMBOK 6th edition document. Our initial step includes the conversion of the PDF document into text and image data. The text data is extracted from each page, capturing the rich content of the PMBOK guide. Simultaneously, images are saved to preserve visual information such as diagrams, graphs, and figures embedded in the document.

### Data Preprocessing

Data cleaning    Data transformation

### Text data processing:

- Text cleaning
- Stop words
- Text sectioning using regex expressions: Regular expressions are employed to address any inconsistencies, such as instances of repeated characters or four-digit numbers that may not contribute to the meaningful information.



Transformers



### Image processing:

- Images are processed to identify and extract Regions of Interest (ROIs) containing relevant visual content. We employ edge detection techniques using the Canny algorithm, followed by contour analysis to identify and filter out areas of interest. These ROIs are then saved as separate images, each associated with a unique UF.



## Feature extraction :

### 1. Synonyms and Definitions:

- Leveraging **spaCy's en\_core\_web\_md** model for precise NLP.
- Identifying synonyms and definitions to capture nuanced meanings.

### 2. Concepts, Objects, and Predicates:

- Custom parsing algorithms for entity extraction.
- Fuzzy string matching using **fuzzywuzzy** for robust relationship recognition.

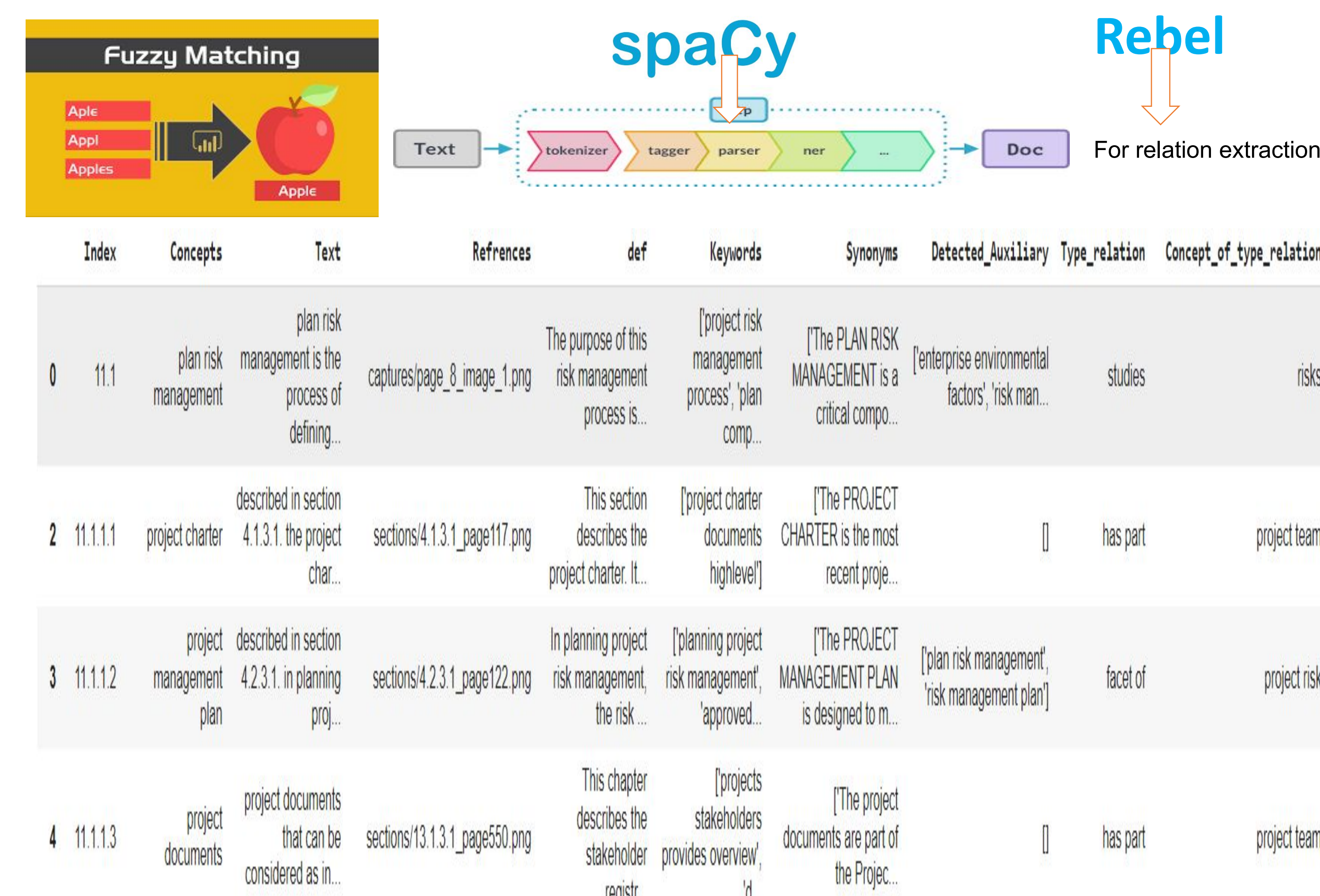
### 3. References:

- Integration of associated URLs for visual context.
- Enriching the dataset with both textual and visual references.

### Transformer-based Models:

- Utilizing transformers library for advanced pattern recognition.
- Enhancing system capabilities for nuanced content analysis.

Our feature extraction process combines these elements, providing a comprehensive understanding of the PMBOK guide's content and laying the foundation for a powerful and context-aware recommendation system.



## Conclusion and perspective

In conclusion, our project successfully tackled the challenge of developing a robust recommender system, utilizing advanced NLP techniques and GNN models. The feature extraction process, covering synonyms, definitions, and essential entities, lays a strong foundation for nuanced content understanding within the PMBOK guide.

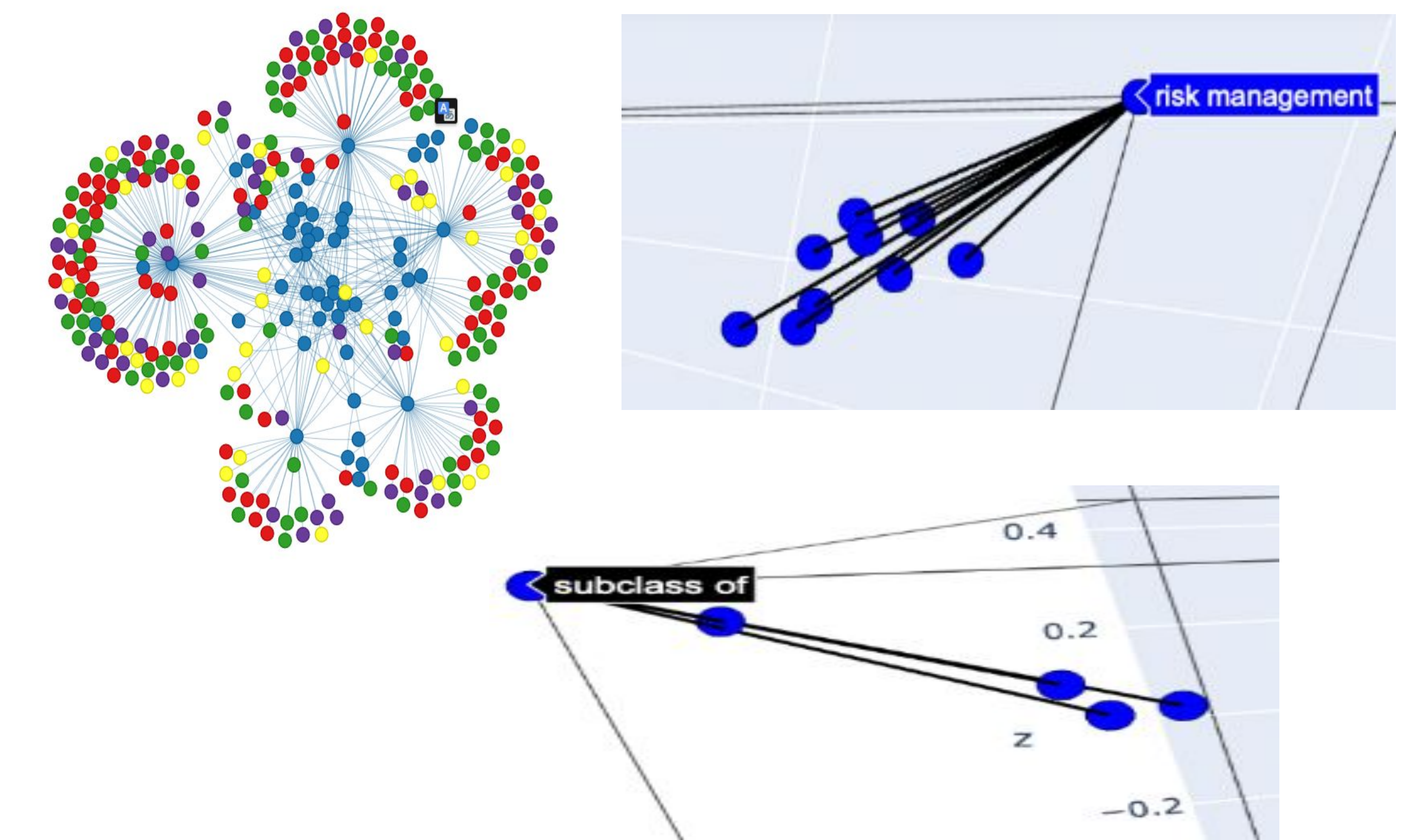
Looking ahead, our perspectives involve the integration of transformer-based models for refined content analysis. We aim to explore dynamic graph structures and evolve the GNN architecture to adapt to different PMBOK editions, ensuring continuous optimization. Envisioning a chatbot interface aligned with PMBOK 6th edition principles adds a real-time, interactive dimension to our system, transforming it into an integral part of project management workflows. These perspectives reflect our commitment to advancing the system in line with evolving project management practices.

## Detailed Approach

### Node and edge embeddings:

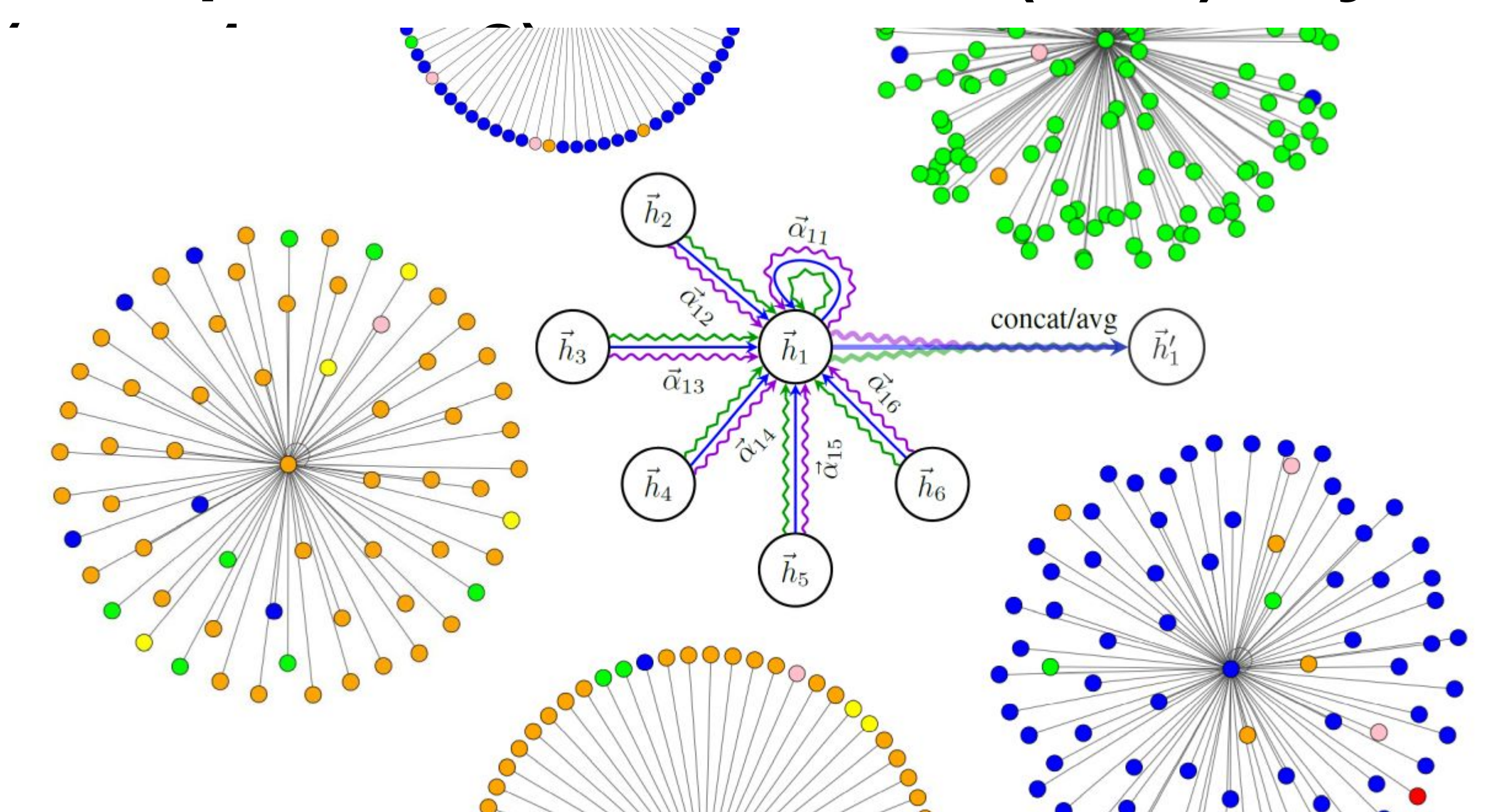
- Node2Vec
- Bert embeddings (bert uncased)

### Graph construction:



## GNN architecture :

### 2 Graph Attention Network (GAT) Layers



## Evaluation

Used RMSE criteria for 20% of data test

Mean Squared Error on Test Set: 2.2422869205474854