

Unveiling Narratives: Analyzing Character Roles and Themes in Novels

Jeeva Saravana Bhavanandam¹, Subramaniya Siva T S¹, Om Sai Krishna Madhav Lella¹, Thejesh Mallidi¹

¹2nd year MS in Data Science, College of Natural Science, College of Engineering, Michigan State University
{sarava25, tsubrama, lellaom, mallidit}@msu.edu

Abstract

This project develops a natural language processing (NLP) pipeline to analyze novels, focusing on named entities, their roles, and dominant narratives. The pipeline identifies and classifies named entities as protagonists, antagonists, or other characters. It further classifies narratives within the novel using embedding-based topic modeling and generates narrative explanations through a Hybrid Retrieval-Augmented Generation (RAG) approach. A Knowledge Graph, enriched with relationships between entities and literary themes, is integrated with a VectorStore to produce coherent narrative explanations, offering deeper insights into plot progression and thematic elements within the novels.

The code for the project is available at: [GitHub](#)
The results data can be accessed at: [Results](#)

1 Introduction

Understanding the structure of narratives and the roles of characters is a fundamental aspect of literary analysis. This project aims to develop a natural language processing (NLP) pipeline that can analyze novels by identifying key entities, such as characters, and classifying their roles within the story. In addition to identifying and categorizing these entities, the pipeline will also focus on extracting dominant narratives and providing explanations for how these narratives evolve and interact within the text. The system will offer insights into character relationships, plot progression, and overarching themes. By focusing on both the extraction and explanation of narratives, this project seeks to shed light on the deeper connections between characters, events, and literary themes. The approach aims to provide a tool that enhances the understanding of complex narrative structures and how they shape the overall storytelling experience in novels, contributing new perspectives to the field of literary analysis.

2 Datasets

We downloaded a collection of 3,036 English books written by 142 authors from the Gutenberg Dataset for our analysis.

3 Literature Survey

Recent advancements in NLP and narrative analysis have explored various aspects of entity extraction, narrative classification, and role identification, but a comprehensive system that combines these elements for novel analysis is still lacking.

3.1 Entity Role Identification

While much of the early work on narrative classification focused on analyzing persuasion and narrative framing techniques across different domains, including novels, these efforts often failed to fully integrate character role identification. Studies like (Coan et al., 2021) and (Piskorski et al., 2022) explored narrative structures but did not establish connections between entities like protagonists and antagonists with their roles in driving the narrative. In addition, works such as (Kotseva et al., 2023) advanced fine-grained narrative classification, offering methodologies relevant to literature but did not extend to linking roles with the overarching narrative framework. Our project addresses this gap by combining character role identification with narrative detection in a unified system for novels.

3.2 Narrative Theme Identification

BERTopic employs transformer embeddings to facilitate effective thematic extraction in novels, allowing for a nuanced analysis of narrative structures and cultural contexts (Grootendorst, 2020). Recent studies have harnessed this approach to uncover common themes across genres and examine reader responses, thereby deepening our understanding of literary works (Khan et al., 2021) (Lee and Kim, 2022). In our analysis, we utilized topic

modeling to identify the major themes within the narratives of the novels.

3.3 Narrative Explanation Generation

Narrative explanation generation is the task of identifying meaningful narratives from datasets and providing contextual explanations that align with their themes and relationships. This process is critical in areas like literary analysis, historical research, and thematic exploration, where understanding both the content and context of narratives is essential. By combining semantic analysis with structured relational models, narrative explanation systems aim to deliver enriched, accurate, and contextually relevant insights, addressing the limitations of traditional keyword-based methods.

4 Methodology

This section outlines the key stages of the NLP pipeline, focusing on four essential components: Named Entity Recognition (NER), entity role classification, narrative classification, and narrative explanation generation. Each component is integral to constructing a cohesive system for novel analysis.

4.1 Named Entity Recognition (NER)

The Named Entity Recognition (NER) module identifies key entities in the text, such as characters, locations, and objects, to support narrative analysis. We initially used SpaCy's pre-trained models but found a high rate of false positives. To improve precision, we fine-tuned a RoBERTa model on the CoNLL-2003 dataset, which provides reliable annotations for person, organization, and location entities. Our final approach combines the RoBERTa model with rule-based enhancements to capture the nuances of literary text, ensuring accurate extraction of entities needed for subsequent analysis. Performance was evaluated using Precision, Recall, and F1-score metrics to confirm its reliability across varied texts.

4.2 Entity Role Classification

The Entity Role Classification module identifies characters' roles within a novel, such as Protagonist, Antagonist, or Supporting Character, and assigns nuanced traits that align with their behavior and influence in the narrative. The approach uses a hierarchical, zero-shot classification method, combining contextual analysis and trait extraction to capture character dynamics comprehensively. The main steps are outlined below:

4.2.1 Preprocessing and Text Splitting

Each novel is split into individual sentences using a regular expression-based tokenizer. This segmentation allows the classifier to accurately capture character interactions and context-dependent information throughout the narrative.

4.2.2 Character Context Extraction

To analyze characters effectively, their mentions in the text are identified and grouped using contextual information derived from prior Named Entity Recognition (NER) results.

The process involves:

- **Character Mention Consolidation:** Mentions of a character, including name variations and aliases, are identified and grouped. Surrounding sentences (two before and two after each mention) are included for added context.
- **Contextual Sentence Aggregation:** Sentences with a character's name are combined into a single document, capturing their interactions, traits, and role within the story.

This approach ensures that the character's context is comprehensive and representative of their function within the story, enabling accurate and consistent classification.

4.2.3 Hierarchical Role Classification

Characters are categorized into main and sub-roles using a hierarchical, zero-shot classification approach implemented with the facebook/bart-large-mnli model. The process involves two primary steps:

1. **Main Role Assignment:** Each character's aggregated context is analyzed to assign a primary role from predefined categories such as Protagonist, Antagonist, or Supporting Character.
2. **Sub-Role Classification:** Based on the assigned main role, a secondary classification is performed to determine more specific sub-roles, providing a nuanced understanding of each character's function within the narrative.

This two-tiered classification enables a detailed portrayal of character dynamics, capturing both their overarching roles and specific contributions to the story.

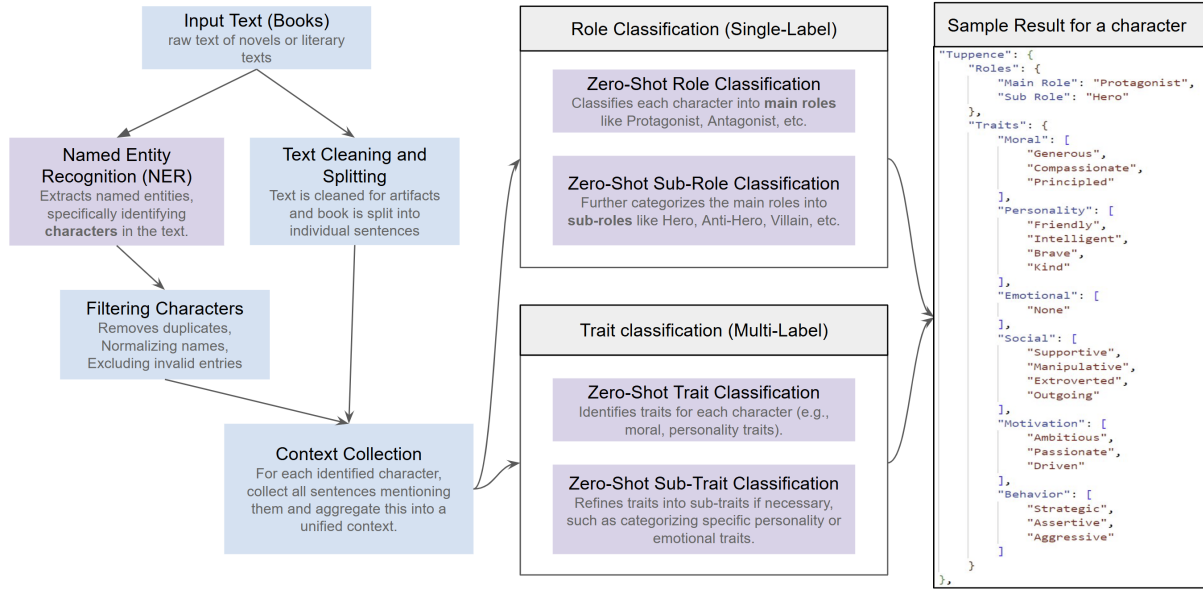


Figure 1: Entity Role and Trait Classification Pipeline.

4.2.4 Trait Classification

To enrich character profiles, traits are identified across multiple domains using a multi-label, zero-shot classification method with the same facebook/bart-large-mnli model. The classification process includes:

1. **Trait Domains:** Traits are organized into categories such as Moral, Personality, Emotional, Social, Motivational, and Behavioral.
2. **Trait Extraction:** For each character’s aggregated context, relevant traits are assigned within each domain. A confidence threshold of 0.6 ensures that only significant traits are retained, allowing for a comprehensive yet precise characterization.

4.2.5 Output Format

The final output includes each character’s main role, sub-role, and relevant traits, facilitating subsequent narrative analysis and enabling a structured examination of character dynamics and thematic influence.

4.3 Narrative Theme Identification

To uncover the various themes, we split each novel into multiple smaller documents, assuming each document corresponds to a single theme. We then transformed these documents into vector embeddings using sentence-transformer models. Next, we applied UMAP and HDBSCAN to reduce the dimensionality of the embeddings and cluster the

documents into common themes. We labeled each cluster as a theme by prompting the LLaMA model with the TF-IDF keywords and the most representative documents. Finally, based on the novel-to-document mapping, we identified multiple themes in each novel. The quality of the theme clusters will be assessed using the Cluster Coherence metrics. Below is a description of the steps used to extract narrative themes from all the novels.

4.3.1 Embedding

We split each novel into chunks of 250 words and converted each chunk into a 384-dimensional vector using the all-MiniLM-L6-v2 model. Considering the maximum number of tokens processed by the all-MiniLM-L6-v2 model, we split each novel into 250-word chunks to ensure no information was lost in the analysis.

4.3.2 Dimensionality Reduction

We reduced the dimensionality of the vectors from 384 to 5 using UMAP with n neighbors = 15. While reducing the dimensionality of the vectors, we assumed a non-linear relationship for each chunk with its 15 neighbors. We avoided the curse of dimensionality by reducing the dimensions from 384 to 5.

4.4 Clustering

We employed a hierarchical density-based cluster model (HDBSCAN) to identify latent themes in the novels. HDBSCAN identifies latent clusters by

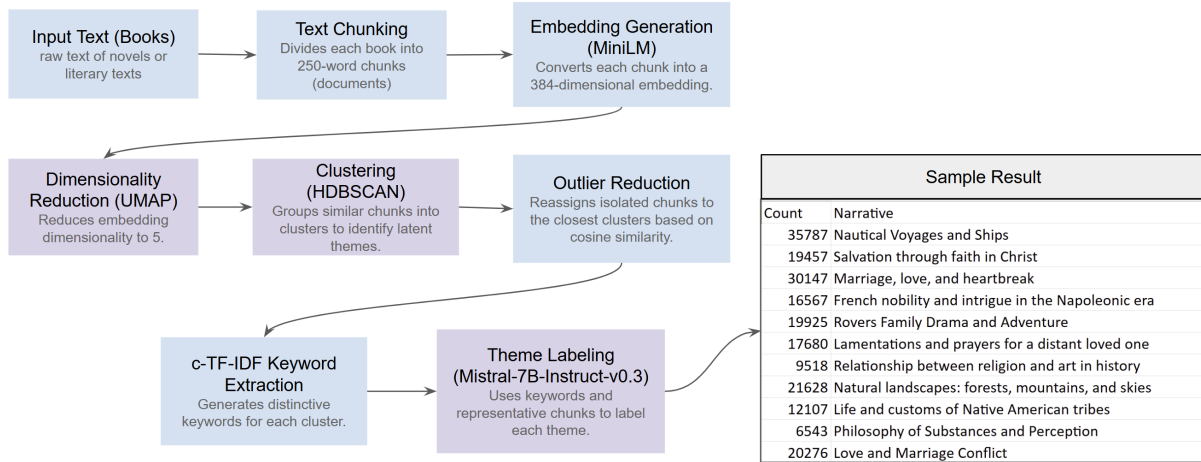


Figure 2: Narrative Theme Identification Pipeline.

analyzing the density of data points and constructing a hierarchy of clusters based on their connectivity and varying densities, allowing it to uncover clusters of different shapes and sizes.

4.4.1 Outlier Reduction

Following HDBSCAN clustering, several chunks were identified as outliers. We reduced the number of outliers by reassigning them to the closest clusters based on cosine similarity.

4.4.2 Theme Labeling

Next, we applied a c-TF-IDF strategy to extract keywords for each cluster. c-TF-IDF is an adjusted TF-IDF representation that considers what makes the documents in one cluster different from those in another cluster. After that, we utilized the KeyBERTInspired() model to fine-tune the keywords. Finally, we used the text-generative large language model, Mistral-7B-Instruct-v0.3, to generate theme labels. At the corpus level, the model was prompted with keywords and the top three most representative chunks for each cluster. At the book level, it was prompted with keywords and ten chunks to ensure more specific labeling. Theme labels were generated separately at the book level because the corpus-level labels, derived from diverse genres, were less specific to individual books.

4.4.3 Coherence Scores

We conducted several iterations by varying the min cluster size from 50 to 500, selecting the optimal model based on the coherence scores of the top 25 clusters. To evaluate the iterations, we utilized the 'c_v' coherence measure, which analyzes word co-occurrence patterns and their distribution within

the documents. This approach offered valuable insights into the interpretability and quality of the identified clusters. After evaluating the scores, we chose iteration 7, with min cluster size = 350, as our best iteration.

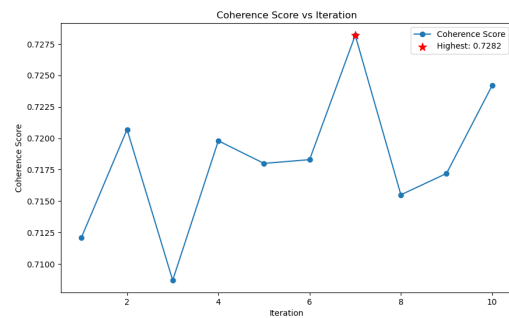


Figure 3: Coherence Scores vs. Iterations.

4.5 Narrative Explanation Generation

Narrative explanation generation extracts meaningful narratives from datasets and provides contextual explanations, critical for tasks like literary analysis and thematic exploration. It combines semantic understanding with structured relationships.

This work introduces a hybrid retrieval system using semantic embeddings and knowledge graphs to enhance narrative retrieval. The objectives are:

1. Develop a **vector-based** semantic search system.
2. Build a **knowledge graph** to encode narrative relationships.
3. Combine both methods into a **hybrid retrieval** system with reranking.

This approach addresses challenges such as capturing thematic connections, structuring relational data, and aligning semantic and structural methods.

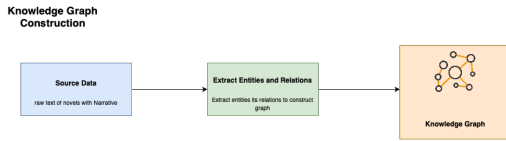


Figure 4: Knowledge Graph Construction pipeline.

4.5.1 Knowledge Graph Construction

The knowledge graph encodes relationships between narratives, explanations, authors, and books for structured and context-rich retrieval.

4.5.2 Graph Schema Design

The schema defines:

- **Entities:** Author, Book, Narrative, Explanation.
- **Relationships:**
 - **Author** → **WROTE** → **Book**
 - **Book** → **HAS** → **Narrative**
 - **Narrative** → **HAS** → **Explanation**

4.5.3 Implementation

Data was ingested into Neo4j, creating nodes and relationships using Cypher queries. This supports efficient traversal and advanced queries, such as retrieving narratives by author or theme.

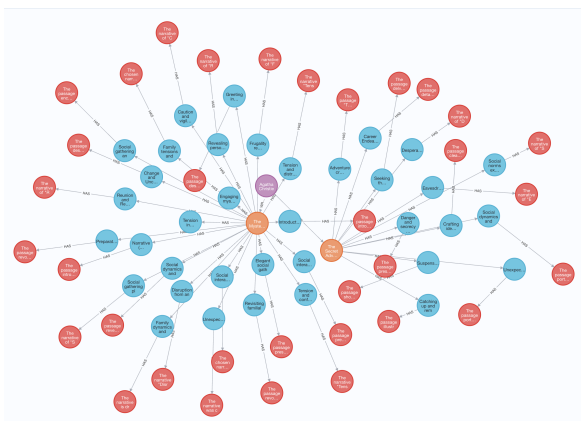


Figure 5: Knowledge Graph in Neo4j.

4.5.4 Semantic Integration

Using principles of semantic and knowledge graphs (Lystopadskyi et al., 2023), the structure supports context-enriched retrieval and reasoning through labeled, real-world connections.

4.5.5 Vector-Based Retrieval

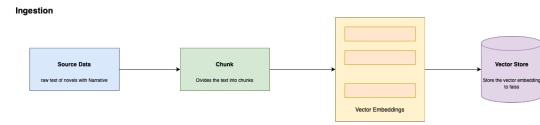


Figure 6: Data Ingestion in Vectorstore.

Vector-based retrieval and reranking were used to enhance semantic search capabilities (Lewis et al., 2021). Narratives were encoded into high-dimensional vectors using the pre-trained all-MiniLM-L6-v2 model from SentenceTransformers, which maps texts into a semantic space. These embeddings were indexed using FAISS with an IndexFlatL2 configuration optimized for L2 distance similarity. Queries were encoded into embeddings and matched with the top- k closest narratives from the FAISS index.

4.5.6 Reranking

Reranking refined the results retrieved via FAISS. A cross-encoder model (cross-encoder/ms-marco-MiniLM-L-6-v2) evaluated query-narrative pairs, capturing fine-grained relationships and generating relevance scores (Dong et al., 2024). Unlike FAISS, which relies solely on vector proximity, the cross-encoder ensures precise ranking by directly modeling query-narrative interactions.

This two-stage system combines the efficiency of FAISS with the precision of the cross-encoder, delivering contextually relevant results across diverse queries and large datasets.

4.5.7 Hybrid Retrieval with Graph and Vector Store

The hybrid retrieval system combines the strengths of graph-based and vector-based methods for enhanced performance. Graphs, using Neo4j, excel at structured queries involving entity relationships, such as retrieving narratives by a specific author or explanations linked to keywords (Sarmah et al., 2024). However, they lack the ability to capture semantic similarity, which is addressed by vector-based methods like FAISS, enabling thematic and nuanced similarity searches in high-dimensional spaces.

4.5.8 Hybrid Integration

The system integrates graph querying for relational accuracy and vector querying for semantic similar-

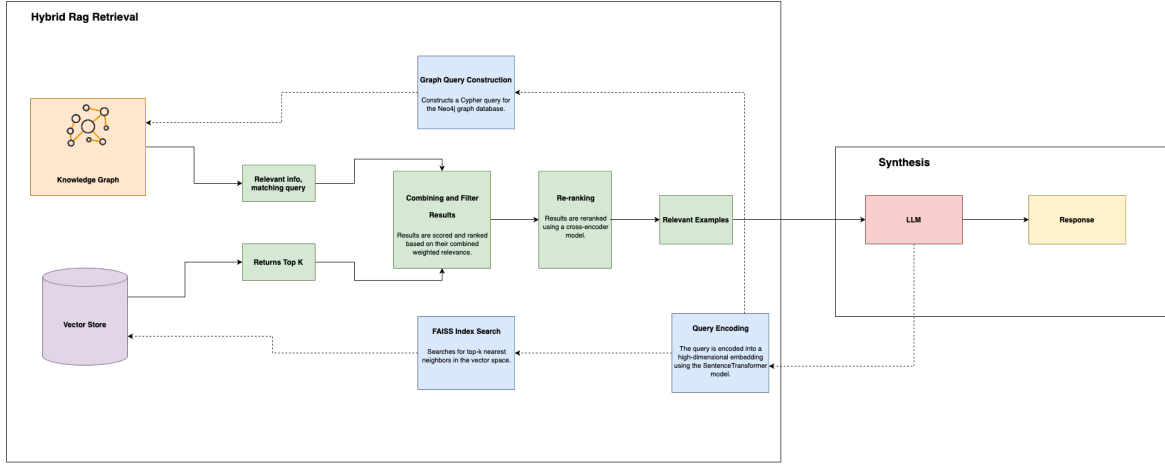


Figure 7: Narrative Explanation Pipeline.

ity. Results are combined using a weighted scoring formula:

$$\text{Score} = (\text{Graph Weight} \times \text{Graph Score}) + (\text{Vector Weight} \times \text{Vector Score})$$

Graph weights (e.g., 0.6) prioritize relational accuracy, while vector weights (e.g., 0.4) capture semantic richness. Duplicate results are filtered, and optional reranking with a cross-encoder ensures contextual alignment with the query.

This hybrid approach leverages the precision of graph queries and the adaptability of vector searches, providing robust, contextually, and semantically relevant retrievals.

5 Experimental Results

This section outlines the outcomes achieved in various stages of our analysis.

5.1 Named Entity Recognition (NER)

Fine-tuning the NER model on the CoNLL-2003 dataset improved accuracy, with the BERT-base model achieving an F1 score of 0.89 and the RoBERTa-base model reaching 0.91, demonstrating enhanced precision for literary text.

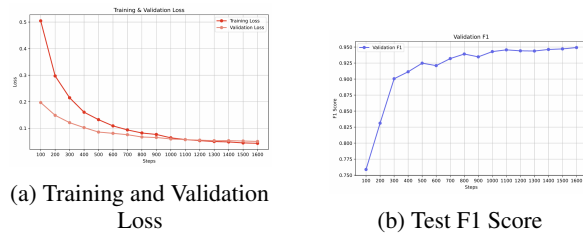


Figure 8: Training/Validation Loss Curves and Test F1 Score for RoBERTa-based NER Model

Below are sample NER outputs, which will be utilized in the subsequent Entity Role Classification task.

```
{
  "entity": "Person",
  "score": 0.9938389658927917,
  "word": "Tuppence",
  "start": 765,
  "end": 773
},
```

(a) NER Output for Tuppence

```
{
  "entity": "Person",
  "score": 0.9912531971931458,
  "word": "Whittington",
  "start": 20955,
  "end": 20966
},
```

(b) NER Output for Whittington

Figure 9: Sample NER Outputs from RoBERTa Model on Novel Excerpts

5.2 Entity Role Classification

The Entity Role Classification module was applied to analyze character roles and associated traits within Agatha Christie's *The Secret Adversary*. The model successfully identified primary and secondary roles for key characters, categorizing them as protagonists, antagonists, or other functional roles within the narrative. Additionally, it extracted various character traits across dimensions like moral alignment, personality, emotional tendencies, social behavior, motivation, and behavioral patterns.

The model effectively classified key characters, capturing their roles and traits with strong alignment to the narrative. Tuppence is identified as a **Protagonist** with a sub-role of **Hero**, characterized by traits like bravery, intelligence, and strategic thinking. Her tactical use of manipulation for morally justified actions, such as impersonating "Jane Finn," reflects her resourcefulness and determination. Whittington, in contrast, is classified



Figure 10: Entity Role Classification Results for Characters in *The Secret Adversary*

as an **Antagonist** with a sub-role of **Oppressor**, embodying traits such as arrogance, hostility, and recklessness. His manipulative and self-serving actions drive much of the conflict but ultimately lead to his downfall.

The analysis highlights the contrast between Tuppence's moral courage and Whittington's corruption, with shared traits like manipulation used in starkly different ways. The hierarchical classification successfully captured main roles and nuanced sub-roles, providing depth to the character analysis. However, limitations include challenges with evolving roles and restricted context length, which can affect accuracy for complex characters. Overall, the classifications enrich our understanding of character dynamics and their contributions to the plot.

5.3 Narrative Theme Identification

We generated a total of 830,642 chunks from 3,036 books by splitting each book into 250-word segments. Among all the chunks, we identified 194 latent themes. Below are two plots: the first displays the top 10 themes and their corresponding chunk counts for the entire corpus, while the second shows the same for Agatha Christie's *The Secret Adversary*.

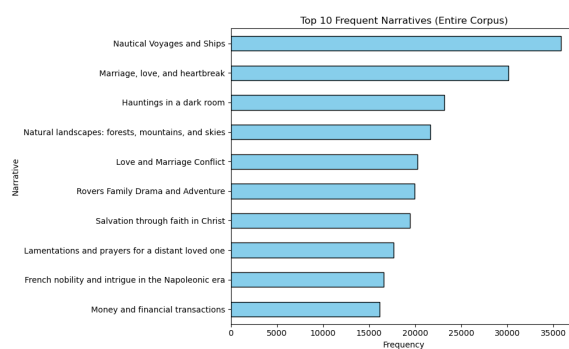


Figure 11: Top 10 Narrative Themes in the Entire Corpus.

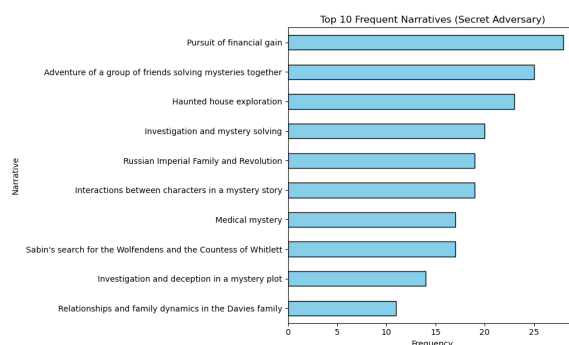


Figure 12: Top 10 Narrative Themes in *The Secret Adversary* Novel.

The analysis of narrative themes in Agatha Christie's *The Secret Adversary* highlights the central role of **"Pursuit of financial gain"**, which drives much of the characters' actions and motivations. Themes such as **"Adventure of a group of friends solving mysteries together"** and **"Investigation and mystery solving"** emphasize the camaraderie and suspense that unfold as the characters work to uncover hidden truths. The inclusion of **"Russian Imperial Family and Revolution"** provides historical depth, while **"Haunted house exploration"** adds an eerie layer to the intrigue.

Personal stakes are explored through themes like **"Interactions between characters in a mystery story"**, **"Medical mystery"**, and **"Sabin's search for the Wolfendens and the Countess of Whitlett"**, which highlight the characters' vulnerabilities and quests. Meanwhile, **"Investigation and deception in a mystery plot"** reinforces the story's focus on uncovering hidden motives. Lastly, **"Relationships and family dynamics in the Davies family"** ties the plot together, illustrating how personal connections

and family ties influence the unfolding mystery. Together, these themes shape a complex narrative filled with intrigue, suspense, and emotional depth.

Overall, we successfully extracted the main themes from all the novels using a computationally efficient approach, without the need to explicitly prompt each chunk in the corpus.

5.4 Narrative Explanation Generation Evaluation

Evaluating language models for explanation generation requires assessing readability, contextual relevance, and diversity. Key metrics include:

Readability: The Flesch Reading Ease measures sentence simplicity, while the Gunning Fog Index evaluates complexity, balancing accessibility and informativeness (Lee and Lee, 2024).

Linguistic Diversity: Entropy captures textual variety, promoting engaging and coherent explanations (Wei et al., 2024).

Contextual Relevance: Semantic Similarity ensures alignment with source contexts, avoiding hallucinations (Mikolov et al., 2013). Topic Coherence, using LDA, evaluates consistency with dominant topics in the input (Röder et al., 2015).

These metrics provide an objective framework for evaluating explanations across readability, complexity, diversity, and alignment, ensuring accessibility and coherence across domains.

5.5 Model Performance

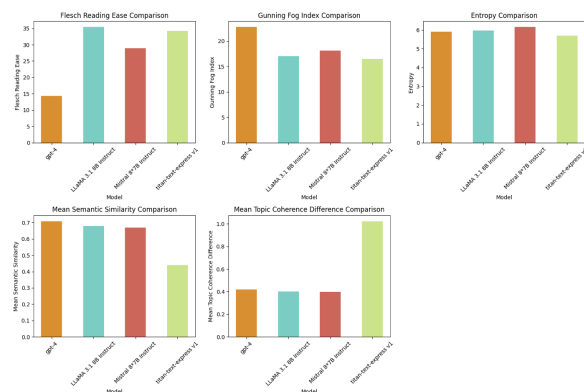


Figure 13: LLM's performance comparison on evaluation metrics Entropy, Flesch Reading Ease, Gunning Fog Index, Semantic Comparison.

The evaluation highlights strengths and weaknesses across metrics:

- **GPT-4:** Achieves highest Semantic Similarity (0.70) and lowest Topic Coherence Difference

Model	Semantic Sim.	Topic Diff.	Flesch	Fog	Entropy
GPT-4	0.70	0.41	14.33	22.75	5.91
LLaMA 3.1 8B	0.678	0.40	35.43	17.01	5.97
Mistral 87B	0.670	0.39	28.94	18.13	6.15
Titan-text v1	0.43	1.02	34.23	16.51	5.69

Table 1: Evaluation metrics for various models

(0.41) but produces complex outputs (Flesch: 14.33, Fog: 22.75).

- **LLaMA 3.1 8B Instruct:** Balances readability (Flesch: 35.43, Fog: 17.01) with moderate Semantic Similarity (0.678).
- **Mistral 8×7B:** Offers moderate coherence (0.39) and readability (Flesch: 28.94, Fog: 18.13) but lower Semantic Similarity (0.670).
- **Titan-text-express v1:** Simplest (Entropy: 5.69) but lags in Semantic Similarity (0.43) and Topic Coherence Difference (1.02).

While readability and complexity metrics assess surface-level quality, embedding-based similarity and topic coherence depend on the robustness of underlying models, highlighting the importance of combining quantitative and qualitative evaluations.

The comprehensive evaluation framework balances readability, diversity, and contextual alignment. Despite limitations in individual metrics, it supports objective insights and fosters trust across diverse applications.

6 Contributions

- **Thejesh Mallidi:** Named Entity Recognition (NER)
- **Subramaniya Siva T S:** Entity Role Classification
- **Om Sai Krishna Madhav Lella:** Narrative Theme Identification
- **Jeeva Saravana Bhavanandam:** Narrative Explanation Generation

7 Conclusion

This study presents a comprehensive NLP pipeline designed to analyze novels by identifying key entities, classifying character roles, extracting dominant narrative themes, and generating contextually rich narrative explanations. By leveraging models like fine-tuned RoBERTa for Named Entity Recognition, hierarchical zero-shot classification for role identification, and hybrid retrieval systems that integrate vector-based semantic search with knowledge

graph structures, the pipeline provides detailed insights into character dynamics, plot progression, and thematic elements. The results demonstrate the pipeline's ability to uncover nuanced relationships between characters, narratives, and themes, contributing to both literary research and the broader field of computational humanities. This work highlights the potential of automated systems to offer deeper understanding and novel perspectives on complex storytelling structures.

8 Future work

Future work will focus on broadening the applicability and functionality of the pipeline. Expanding support for multilingual and genre-specific analysis will enable the study of a wider range of global narratives, addressing the linguistic and cultural diversity of literary traditions. Incorporating dynamic role tracking will allow for more accurate analysis of characters as they evolve throughout the narrative. Sentiment and emotion analysis can add a layer of psychological and emotional depth to character and plot analysis, enriching the system's interpretative capabilities.

References

- T. G. Coan, W. Ahn, and J. Lanning. 2021. [Narrative structure and role identification in storytelling](#). *Journal of Narrative Theory*, 41(2):145–172.
- Jialin Dong, Bahare Fatemi, Bryan Perozzi, Lin F. Yang, and Anton Tsitsulin. 2024. [Don't forget to connect! improving rag with graph-based reranking](#).
- Maarten Grootendorst. 2020. Bertopic: Leveraging transformers for topic modeling. <https://github.com/MaartenGr/BERTopic>. Accessed: YYYY-MM-DD.
- S. Khan, J. Weber, and M. Thakkar. 2021. [Understanding reader responses through thematic analysis of novels](#). *Digital Humanities Quarterly*, 15(3).
- M. Kotseva, K. Novak, and T. Ivanova. 2023. [Advancing fine-grained narrative classification for literature](#). *Computational Narratology Review*, 5(4):222–234.
- Bruce W. Lee and Jason Hyung-Jong Lee. 2024. [Traditional readability formulas compared for english](#).
- H. Lee and J. Kim. 2022. [Cultural contexts in fictional narratives: An embedding-based approach](#). *Literary Studies and Digital Methods*, 8(1):23–38.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#).
- Daniil Lystopadskyi, André Santos, and José Paulo Leal. 2023. [Narrative Extraction from Semantic Graphs](#). In *12th Symposium on Languages, Applications and Technologies (SLATE 2023)*, volume 113 of *Open Access Series in Informatics (OASIs)*, pages 9:1–9:8, Dagstuhl, Germany. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#).
- J. Piskorski, J. S. Bień, and K. Wieczorek. 2022. [Character role extraction in literary texts: A comprehensive framework](#). *Literary and Linguistic Computing*, 37(1):65–87.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. [Exploring the space of topic coherence measures](#). In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM '15*, page 399–408, New York, NY, USA. Association for Computing Machinery.
- Bhaskarjit Sarmah, Benika Hall, Rohan Rao, Sunil Patel, Stefano Pasquali, and Dhagash Mehta. 2024. [Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction](#).
- Lai Wei, Zhiqian Tan, Chenghai Li, Jindong Wang, and Weiran Huang. 2024. [Diff-erank: A novel rank-based metric for evaluating large language models](#).