

# COMP47250: Final Report

## Analysing Character Arcs and Story Development in Fiction Writing using Natural Language Processing

**Student:** Declan Atkins 14388146

**Supervisor:** Derek Greene

**Abstract.** When reading a work of fiction, character arcs and story development are often the most fundamental and defining aspects. While the details may fade from memory over time, readers tend to recall the overall plot and the evolution of the characters throughout the narrative. The success and failure of the characters we identify with become significant to us and we become emotionally invested in their journeys. As such, the capacity to craft coherent character arcs is a critical component of the fiction writing craft, and one of the most vital skills that authors must master. It also forms the backbone of a critique of a novel or other piece of writing. In this dissertation, we investigate the challenge of automatically extracting the character and plot arcs from novels using Natural Language Processing techniques. Through the application of these techniques, our aim is to uncover patterns and extract meaningful insights from literary texts, shedding light on the structural and emotional journeys of characters within narratives. This research not only contributes to the field of computational literature analysis, but also provides insights into the art of storytelling and character development, bridging the gap between literary analysis and computational approaches.

## 1 Introduction

The art of storytelling has been an integral part of human culture since time immemorial. Whether through ancient myths, epic poems, or contemporary novels, narratives have the power to captivate our imaginations and transport us to different worlds. Within the realm of fiction, the development of characters and the progression of their stories are central elements that engage readers on a profound level. While the specific details of a story may fade from memory, the overall plot and the evolution of characters tend to leave a lasting impression.

Understanding the character arcs within a novel is vital not only to the author but also to us as readers and to critics reviewing the quality of the story. By unravelling the complexities of these narrative components, we gain insight into the art of storytelling and the underlying themes and messages conveyed through literature. However, the analysis of character arcs and story development in the literature has been predominantly based on subjective interpretations and manual examinations of texts [21, 23].

We can define a *character arc* as a transformative journey undertaken by a character throughout the course of a novel [9]. The character begins the novel as one kind of person and subsequently changes based on the events that occur over the course of the novel. As discussed in [5], the main protagonist of a novel is the most likely to have a well-defined character arc, although side characters may also experience similar journeys throughout the plot. Some notable examples of character arcs include the downfall of Dorian Gray in *The Picture of Dorian Gray* by Oscar Wilde and the redemption of Ebenezer Scrooge in *A Christmas Carol* by Charles Dickens. Although diametrically opposed to each other in terms of the direction of their arcs, both of these characters experience a drastic change in their personalities over the course of their respective novels.

In recent years, advanced Natural Language Processing (NLP) techniques, combining machine learning and computational analysis, have revolutionised literary research. These techniques improve traditional subjective analysis by literary critics, offering a more quantitative approach to understanding fiction. This computational approach also enables the exploration of complex analyses that were once hindered by labour-intensive manual methods. One prominent method in this realm is “distant reading,” which focuses on extracting patterns and implicit rules in literature through objective data analysis across multiple works. Unlike traditional “close reading,” which delves into the individual meaning of a single work, distant reading compiles and analyses data from numerous pieces. In character arc analysis, distant reading has been used to examine structural changes that define narrative progressions. For instance, [8] investigated plot development through shifts in character social network graphs. In this dissertation, we will present our research on understanding character arcs in fiction writing through the use of NLP. We define the goal of this project as answering three key research questions, namely:

1. How can we extract characters from fictional texts?
2. How can we quantitatively model a character arc?
3. How can we compare quantitative models of different character arcs?

We will undertake this by examining four novels of 19th century literature, chosen because they have been widely studied from a literary criticism perspective and have well-defined character arcs involving the central protagonist. They are the following:

1. The Picture of Dorian Gray, Oscar Wilde (1890)
2. A Christmas Carol, Charles Dickens (1843)
3. Silas Marner, George Elliot (1861)
4. Pride and Prejudice, Jane Austen (1813)

The remainder of this report will be structured as follows. In Section 2, we will discuss existing state-of-the-art research and other work related to our project. In Section 3 we will describe the development work undertaken to answer these research questions. In Section 4 we will analyse the results we have obtained and compare them with the established literary critique of the novels we have studied. Finally, in Section 5 we discuss our conclusions and some avenues for future work.

## 2 Related Work

Several studies have explored the analysis of character arcs and story development in fiction writing using Natural Language Processing (NLP) techniques. These investigations have provided valuable information on narrative structures and character dynamics within literary works.

Character extraction is one of the key tools used in the computational analysis of fiction. This, as the name suggests, is the process of extracting unique characters from the text and involves a combination of Named Entity Recognition (NER) and co-reference resolution. Named Entity Recognition (NER) is a crucial task in NLP that involves identifying and classifying named entities within text [17]. Named entities refer to specific types of entities, such as people, organisations, locations, and any other class of entities that we wish to extract from text content. As discussed in [3], traditional NER pipelines can struggle with literary texts. These systems are often trained in news articles, blog posts, and social media content, which contain entities that are very different from those of fiction novels. It was also shown in [1] that the use of specific datasets, drawn from fiction novels, for training NER models, can aid in character extraction. The work presented in [4] provides us with BookNLP<sup>1</sup>, a library designed to process fiction novels. This library includes a comprehensive set of built-in character extraction techniques.

To study novels from a structural point of view, we can apply social network analysis (SNA). This involves the study of a graph which represents associations between different characters within a text. This approach has been shown to be a powerful tool to shed light on the interconnectedness of characters and their influence within various social contexts. In [8], the authors used the change in the social network of a novel over time to study plot development. This could also be used to study character arcs by adapting this technique to work with the changing subgraph of a chosen character.

We could also focus on more content-based changes over the course of a novel to study character arcs. Sentiment analysis and topic modelling provide two alternative methods by which this could be done. In [18], sentiment analysis applied to the dialogue between characters was used to track the progression of their relationships. It was also applied in [2] to study the works of H.P. Lovecraft. Topic modelling has been applied similarly to study literary works, such as in [20], where it was applied to uncover themes in the works of Charles Dickens.

The fusion of literary and computational research has ushered in a new era of exploration, marrying the rich heritage of literary analysis with the potency of computational methodologies. By merging the nuanced comprehension of texts honed in literary studies with the analytical prowess of computational tools, scholars can unlock profound insights into the intricacies of literature. However, this confluence of fields has caused some apprehension. The virtues and shortcomings of distant and close reading are extensively deliberated in the existing literature. Close reading, involving meticulous analysis of textual passages within

---

<sup>1</sup> <https://github.com/booknlp/booknlp>

a novel to reveal a profound understanding of content, represents the conventional approach to literary study. On the contrary, distant reading, as introduced by [16], employs graphs and tools to create a comprehensive overview of a single novel or a collection of novels. [19] contrasts these methodologies and presents a framework for their fusion: employing distant reading to identify text segments that merit detailed examination through close reading. Our project similarly employs a hybrid approach. To capture a character arc, both detailed, close-reading understanding of characters and interactions, as well as the high-level, distant-reading comprehension of plot, are crucial.

### 3 Research and Implementation

As discussed in Section 1, the work carried out in this dissertation had the aim of answering three key research questions. We will now discuss the methods by which these questions were tackled and the success and failures we encountered.

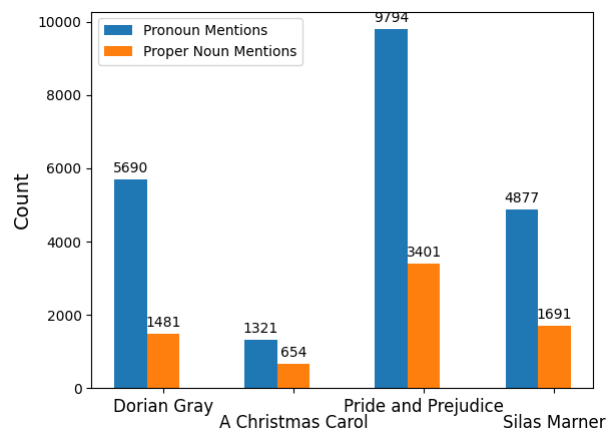
#### 3.1 RQ1: How can we extract characters from fictional texts?

Extracting characters from fiction novels requires a combination of tools such as Named Entity Recognition (NER) and co-reference resolution [10]. First, entities relating to people are extracted from the text, before co-reference resolution techniques are applied to link the entities that refer to the same characters. In the initial phase of this project, we investigated two methods for character extraction with the aim of building a pipeline that could be reused on multiple texts. The first candidate was a pipeline built using SpaCy<sup>2</sup> and a set of heuristic rules, while the second was based on BookNLP.

Our first attempt at character extraction was made using SpaCy NER. We initially extracted the named entities from the target novel before filtering for those that refer to people. Although we then had a list of people mentioned within the novel, this is still quite far from a complete character extraction pipeline. We must post-process this list to determine the unique characters in the novel and find where they are mentioned. This raises several difficult challenges. First, this pipeline faces specific complexities regarding how characters are referred to. Characters in stories often have various forms of identification, such as full names (“Dorian Gray”), shorter versions (“Dorian”), honorifics (“Mr. Gray”), or even nicknames. Further complicating this is the fact that novels often deal with several members of the same family. As such there may be several syntactically-correct solutions to which character a named entity refers too, and the true solution may only be deduced from the context of the entity mention. For example, in *Pride and Prejudice*, the mention “Miss Bennet” could refer to any one of the five Bennet sisters. Without the contextual knowledge of which characters are currently present in the scene, we have no way of determining to which character this mention refers.

---

<sup>2</sup> <https://spacy.io>



**Fig. 1.** Counts of pronoun vs proper noun usage in character mentions within the four studied novels.

A second challenge facing this pipeline is that of pronouns. For a character extraction pipeline, NER alone is not enough. Since the use of pronouns in novels is generally far more common than full names, such character mentions will not be picked up by NER. This disparity is shown in Figure 1. In *A Christmas Carol*, where the counts of pronoun usage to proper noun usage are closest, the ratio is still 2:1. Clearly, by not taking into account these pronouns mentions, our data will be greatly limited. In [22], the issue of pronouns in co-reference resolution is extensively discussed. Even with deep learning transformer-based approaches, such as those proposed in that paper, this is a very difficult challenge to overcome. Specifically, the authors found that only 30% of all examples were solved by all implemented methods and 15% were not solved by any.

Due to the limitations of the SpaCy pipeline and the weaknesses around its implementation for the current task, we decided to investigate alternative methods of character extraction. BookNLP, first presented in [4], was the obvious choice. This Python library provides an extensive set of tools for processing long texts, particularly fiction novels. For our project, the relevant functionality included entity extraction, dialogue speaker identification, and coreference resolution pipelines. With the results of almost a decade’s worth of research available in this library, we found that it far surpassed our attempts at character extraction, so we decided to adopt it as the basis for the rest of our work. However, even this is not a perfect solution. The character clustering algorithm used is not completely accurate. As such, drawing on the lessons learned from building the SpaCy pipeline, we added a post-processing step to further combine characters based on name structures. For example, in *A Christmas Carol*, we were able to merge “Mrs. Fezziwig” and “Old Fezziwig”, “Dick Wilkins” and “Poor Dick”, and “Belinda Cratchit” with “Miss Belinda”. Each of which had originally been parsed as separate characters. Similar merges were done for the other studied novels, providing us with an acceptably correct list of characters

for each novel. As such, our first research question was successfully answered, so we could proceed with our next research question.

### 3.2 RQ2: How can we quantitatively model a character arc?

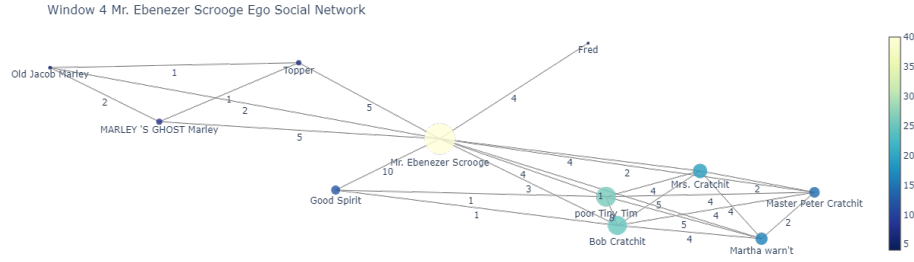
Defining a quantitative model of a character arc requires referring back to our original definition from Section 1. In [9], we are told that a character arc represents a transformative journey undertaken by a character over the course of a novel. However, this is a fairly subjective definition and is open to different interpretations. In [6], character arcs were defined and tracked according to the events that occurred within the novel. Although this is certainly a valid approach and was shown in the paper to provide interesting insight into the plot, we decided to go a different route and focus on more abstract concepts to which we could apply quantitative measures. In particular, we considered three different perspectives: 1) character social networks, 2) sentiment of character mentions and dialogue, and 3) shifting topics around character mentions.

**Character Social Network Analysis.** Social network analysis (SNA) in fiction involves a systematic exploration of the relationships and interactions between characters within the narrative. This analytical approach adapts principles from network science to understand the intricate web of connections, affiliations, and dynamics that shape characters' roles and impact the storyline. By visualising and quantifying these relationships, SNA provides insights into character centrality, influence, and the overall structure of the fictional world. Here our focus was primarily on the evolving *ego network* of the protagonist as the narrative unfolded. The ego network in this case is defined as the subgraph of characters that are connected to the protagonist, together with the ego itself.

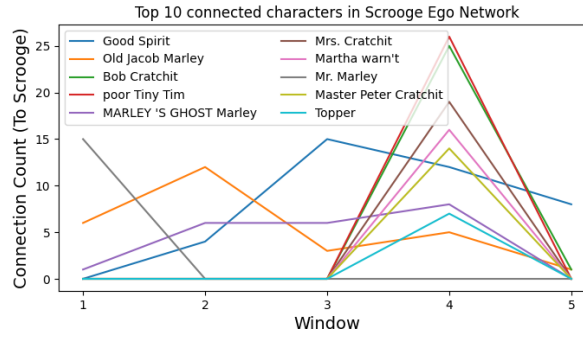
To show how such a structure changes over time, we divided the books into windows consisting of multiple paragraphs, before plotting the network for each of these windows. The number of windows is a hyperparameter that had to be selected based on the length of the novel and the level of detail desired. We experiment with several values for this for each novel, settling on different values for each. For *A Christmas Carol*, we chose to split the novel into five windows, roughly corresponding to one window per chapter<sup>3</sup>. In Figure 2 we can see the Ego Social-Network for window #4.

While the ego networks provide a snapshot of the current state of a character arc, they do not show us how a character has developed with respect to the previous windows. As such, a better way to visualise this evolution is to calculate metrics from each graph and plot how they change over time. In particular, we looked at the characters most connected to the protagonist and how this connectedness changed over the course of the novel. In Figure 3 we can see this graph for *A Christmas Carol*. This allows us to identify some significant moments

<sup>3</sup> It is simply coincidental that one window per chapter worked here. In *A Christmas Carol* the chapters are each roughly equal length. For the other novels, multiple chapters were contained within single windows.



**Fig. 2.** Ego Social-Network for the character of Scrooge for window #4 of 5 for the novel *A Christmas Carol*.



**Fig. 3.** Top-10 characters most closely connected to Scrooge in *A Christmas Carol*, plotted by connection count over 5 windows.

within the novel. In Stave Four, Scrooge visits the Cratchit family. We can spot this easily based on the sudden rise in their connections to him.

As we have shown, the application of SNA to fiction novels can provide us with insights into the progression of character arcs. However, it is somewhat limited in what it shows. While how the protagonist is connected to other characters is an important aspect of their character arc, it is not the only perspective that we should investigate. Of equal importance is how these characters interact. For instance, while a rise in co-occurrences could indicate a growing friendship between two characters, it could also be indicative of the exact opposite. As such, further analysis is required to build our quantitative model.

**Sentiment Analysis.** To provide a more nuanced look at how a character's personality and interactions evolve over the course of a novel, we considered the use of sentiment analysis. As discussed in [12], such an approach can be used to uncover the personalities of characters within novels. Although this does not correspond directly with the aim of our research, it does indicate that we can use sentiment analysis as a tool to model how those personalities change. For our work, we used the VADER sentiment model [11] as provided by NLTK<sup>4</sup>.

<sup>4</sup> <https://www.nltk.org/>

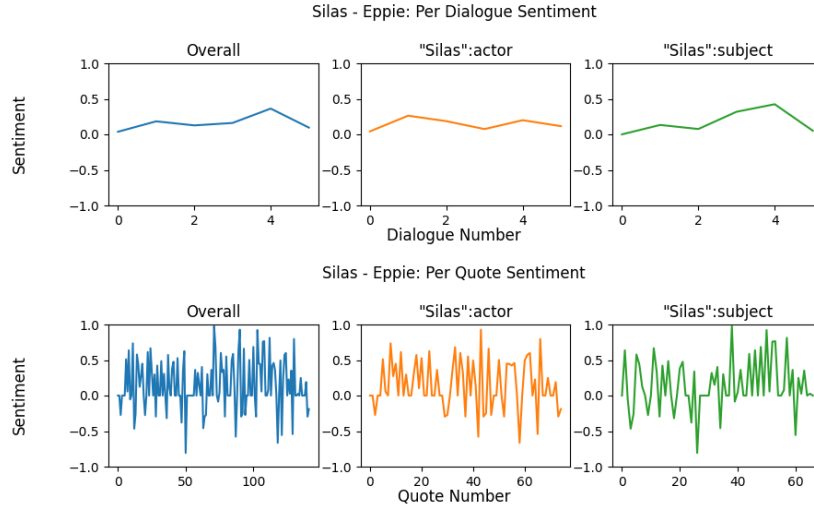
Although this model is not specifically trained on literary data, we found that its results were adequate for our needs. However, as will be discussed in Section 4, we encountered certain limitations in terms of accuracy. In this dissertation, we approach sentiment analysis in two ways. First, we looked at how the sentiment of dialogue involving the protagonist changed throughout the novel, before looking at how it changed on a character mention level throughout the text.

BookNLP provides a pipeline for identifying the speaker in an instance of quoted dialogue. We decided to analyse this in two ways: 1) by looking at how the dialogue between the protagonist and a specific second character evolved over the course of the novel; 2) by looking at how the sentiment of the protagonist’s dialogue with all characters changed. We grouped the instances of quoted dialogue into individual conversations by defining a threshold count of tokens between one instance of quoted dialogue ending and the next starting. We then filtered this list to contain only those conversations where the protagonist was an active participant. This then provided us with a list of conversations suitable for further analysis.

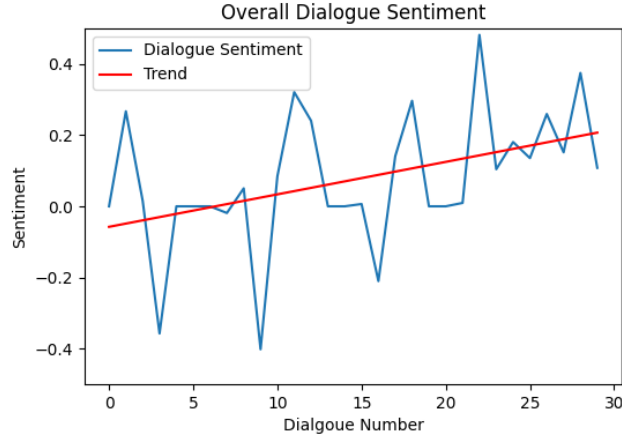
For the first part of our dialogue analysis, we focus on how the dialogue between the protagonist and another character shifts over time. We plotted this from an overall sentiment perspective and with respect to each unique pair of characters (i.e., the sentiment of their quotes). In Figure 4 we can see this for the pair of characters Silas and Eppie in *Silas Marner*, both on a per-dialogue and a per-quote level. While we observe that this has the potential to show interesting insight into how the dialogue between two characters changes, which may be relevant to certain character arcs, we observed that the data for character pairs were too sparse to provide a reliable representation, even for central characters. As mentioned previously, in [18] this method was used to study characters in drama. However, the density of dialogue in plays is much higher, and moreover, it is much easier to detect. The second aspect of our dialogue analysis provided far more useful results. Here, we simply provide a sentiment score for each individual conversation involving the protagonist and plot how this changes over time. In Figure 5, we can see this plotted for *Silas Marner* and we can see an upward trend as the novel progresses.

We also investigated an alternative approach to sentiment analysis, focusing on character mentions. To do this, we took the 50 tokens surrounding each character mention and plotted the rolling average of the sentiment in this text over the entire novel. During this process, we focused primarily on the novel protagonists. For example, in Figure 6 we can see a visualisation of this approach for mentions of Silas in *Silas Marner*. This is a very useful plot, as it contains information that is not captured within the dialogue. Often, while the actual quoted text of the dialogue may be neutral, the words describing the speech that occur outside of the quote can tell us more what the true sentiment of the quote is. For example, in chapter 27 of *Pride and Prejudice*, the text “La! my dear,” said Maria, quite shocked at the mistake appears. Passing the quoted dialogue to Vader gives a compound sentiment of 0.44, while including the text following the quote gives a compound sentiment of -0.4. Clearly, the context of





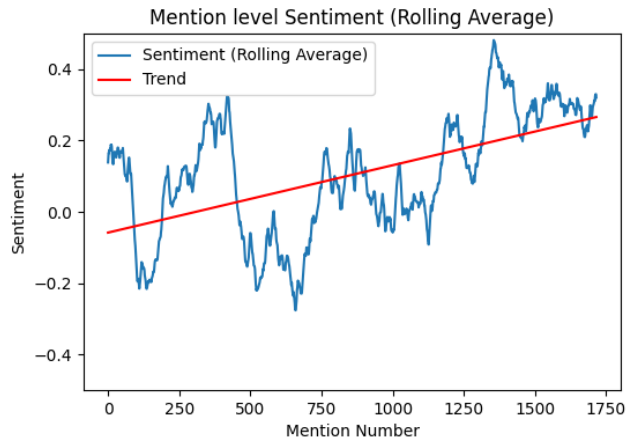
**Fig. 4.** Dialogue-based sentiment analysis between Silas and Eppie in *Silas Marner*.



**Fig. 5.** Overall dialogue sentiment and dialogue sentiment trend in *Silas Marner* of conversations involving Silas.

the quote is as important as the quote itself. Furthermore, for certain novels, the amount of quoted dialogue is limited, and the plot would be subject to outlier weighting. However, this is not the case for *Silas Marner*, as although the dialogue plot includes much more variability than the one at the reference level, the trend and the peaks are well matched.

As we have seen, sentiment analysis can provide us with a means of looking beyond the structural aspects of the novel that might be uncovered through network analysis. Specifically, it allows us to examine more qualitative changes in how characters are being discussed. For certain types of character arc, such



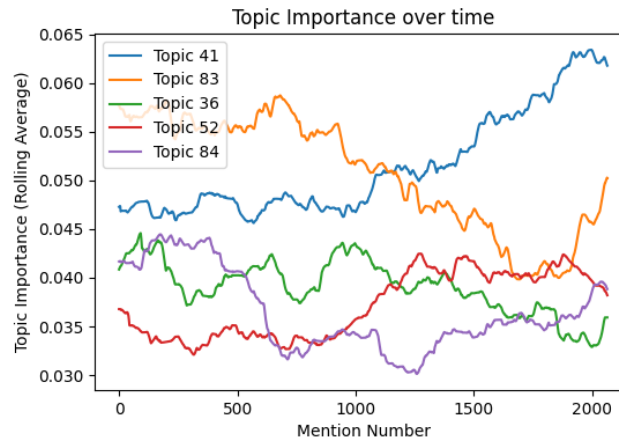
**Fig. 6.** Rolling Average and Trend for Sentiment of the text surrounding mentions of the character of Silas in *Silas Marner*.

as those appearing in novels like *The Picture of Dorian Gray*, *Silas Marner*, and *A Christmas Carol*, this might help us to map the trajectory of a character who has a dramatic change in personality (e.g. from innocence to evil or from a greedy miser to a compassionate person). As will be discussed in Section 4, this is not as effective for novels such as *Pride and Prejudice*, where the arcs are more neutral and subtle.

**Topic Modelling** Topic modelling is a common unsupervised text mining approach that provides an automated procedure to encode the content of a corpus of texts into a set of substantively meaningful categories called “topics”. By monitoring the changes in these topics over the course of the novel, we can uncover thematic shifts in the content and, as such, track concepts that are not as visible through metrics such as sentiment. In our work, we focus on the Latent Dirichlet Allocation (LDA) algorithm [7], as implemented in the Python Gensim library<sup>5</sup>. This algorithm takes a probabilistic view of representing each document in a corpus as a mixture of different topics which are uncovered through the analysis of a large corpus.

For our research, we built a single LDA model for the entire corpus of four novels with 100 topics. This model was trained on documents of each paragraph from the novels. To create the documents, we fed each paragraph through a pre-processing pipeline, involving stopwords and proper noun removal, to avoid capturing irrelevant information and topics that simply referred to characters. We had previously attempted to train the model with shorter documents, each being one sentence within a novel; however, these proved to be too short to impart useful information.

<sup>5</sup> <https://radimrehurek.com/gensim/>



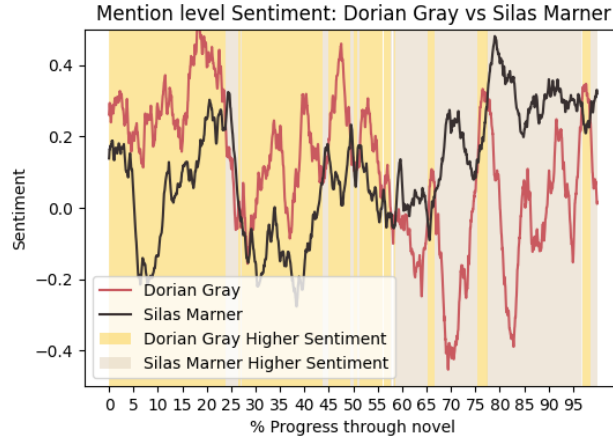
**Topic Top-10 Words:**

- Topic 41: time, feeling, evening, sort, minute, master, month, story, society, mere
- Topic 83: life, sense, live, reason, lose, form, beauty, age, create, beautiful
- Topic 36: feel, gold, hair, trust, lip, handsome, blue, youth, worship, scarlet
- Topic 52: speak, eye, care, heart, intend, guess, curious, concerned, heal, demean
- Topic 84: marry, wife, meet, marriage, laugh, party, forget, charm, dine, absurd

**Fig. 7.** Topic importance over time in *The Picture of Dorian Gray* for a selection of key topics, together with the top-10 words for these topics.

The metric that we create from applying topic modelling is the shift in “topic importance”. For each document that is used to query our LDA model, we retrieve a set of topics that are relevant to the document and a score for how relevant they are. This score can be thought of as the importance of the given topic to the document. We created further documents from the novels, defined as fifty tokens before and after a mention of the protagonist in a novel, and used these to query our LDA model. By tracking the rolling average importance of each topic over time, we should be able to see how these topics relate to the character arc.

Figure 3.2 shows the rolling average importance of the top 5 topics for *The Picture of Dorian Gray*. The decrease in importance of topics 83 and 36 is worthy of note, as the novel describes the downfall and descent into evil of the protagonist, it is interesting to see a drop in topics where words like “beauty” and “youth” are important. However, this also highlights one of the downsides of topic modelling as a visual metric: It is quite difficult to understand at a glance why certain topics exhibit a rise in importance, making it easy to fall into the trap of confirmation bias. Retrieving the next 5 words for Topic 41 reveals words such as “sorrow” and “terror”. however, we have to draw a line at some point to ensure that we are not simply adding words until our hypothesis is satisfied.

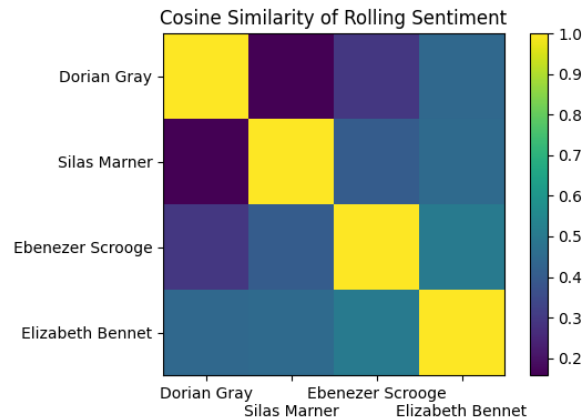


**Fig. 8.** Comparison of rolling-average sentiment across the novels *The Picture of Dorian Gray* and *Silas Marner*.

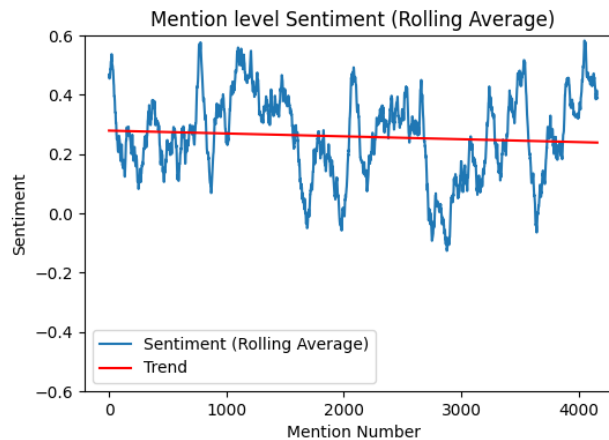
### 3.3 RQ3: How can we compare quantitative models of different character arcs?

Our final research question was substantially answered by the second. In order to analyse the character arcs across these novels, we can utilise the measurements established in the previous section. By representing these metrics on a shared graph, we can visually highlight the similarities and variations between them. A good example of this is shown in figure 8. In this chart, we see the rolling average sentiment of the mentions of Dorian Gray plotted against the rolling average sentiment of the mentions of Silas Marner. We also highlight the areas where each is higher than the other. What we can see from the plot is that, in general, these arcs are the opposite of each other, with higher sentiment for Dorian at the start and lower at the end, and vice versa for Silas Marner.

We can extend this approach further by applying second-level metric computations to the results. One such metric is cosine similarity. By applying this metric to the rolling average sentiment of the four protagonists, we can give a quantitative value for their similarity, as determined by the primary metric of rolling average of mention level sentiment. This graph is shown in Figure 9. In this graph, we can see a low similarity between the character of Dorian Gray and the characters of Scrooge and Silas Marner. As we will discuss in more detail in Section 4, this is what we intuitively expected, as the first is a downfall arc, while the other two are stories of redemption. However, we can also see a drawback to taking only one of the primary metrics to compute similarity. The character of Elizabeth Bennet is determined to be reasonably similar to all three other protagonists, despite the fact that we know that there are stark differences between the other three. This is because her character arc is much less defined by a change in sentiment, as shown in Figure 10. In this chart, we see that the sentiment varies little between the start and end of the novel.

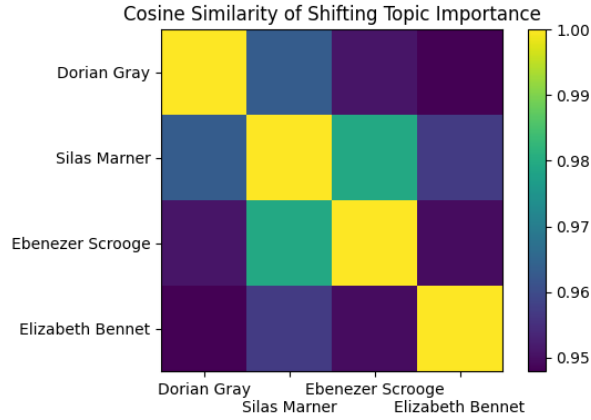


**Fig. 9.** Heatmap illustrating the cosine similarity for rolling sentiment of the protagonists in each of the four novels.



**Fig. 10.** Rolling average and trend for sentiment of the text surrounding mentions of the character of Elizabeth in *Pride and Prejudice*.

In Figure 11 we see a graph showing the cosine similarity of the change in topic importance between the novels. In this chart we see a very high similarity between Scrooge and Silas Marner. Also worthy of note is the fact that there is a much higher similarity between Dorian, Silas, and Scrooge than between any and Elizabeth. This further speaks to the difference in content between *Pride and Prejudice* and the other novels. While the initial and concluding themes in these novels differ significantly, their overall subject matter exhibits greater similarity to each other than to the thematic content and its changing significance in *Pride and Prejudice*.



**Fig. 11.** Cosine similarity of shifting topics for the protagonists in our four novels.

As discussed in Section 3.2, we have considered multiple quantitative ways to define a character arc. Through an amalgamation of these metrics, we can compare the character arcs from different novels and determine the true similarity between them, thus answering our third and final research question.

## 4 Analysis and Discussion of Results

To assess the validity of the results obtained in the previous section, we explore them in the context of previous work in literary criticism. The novels we have chosen are well studied and, as such, many relevant critical examinations exist for them. In this section, we consider a number of existing relevant studies and show how our results conform to the expert analysis of our chosen texts.

### 4.1 Defining the Expected Arcs

To evaluate the accuracy of our methods for extracting and quantifying character arcs, we must compare our results with well-established literary analyses of the selected texts. In this subsection, we will define the expected arcs for the protagonists of our chosen novels based on the literature.

**Dorian Gray, *The Picture of Dorian Gray*.** As discussed in [14], Dorian finds himself trapped within the influence of two characters, Basil Hallward and Henry Wotton. These characters and the conflict between them represent “two mutually exclusive interpretations of human experience: one, optimistic, religious and emotional; the other, pessimistic, cynical and intellectual.” Basil is eager for Dorian to follow a righteous path, promoting Dorian’s attempts at philanthropy and, after a time, his engagement to Sibyl Vane. This is in turn ridiculed by

Henry, who seeks to influence and dominate Dorian's imagination and encourage a life devoted to self-gratification and refusal of any moral standard.

The death of Sibyl Vane shows the victory of Henry and begins the downfall of Dorian. As discussed in [13], despite his acceptance and willingness to adopt a life of sin, Dorian still understands the importance of morality in his culture. Toward the end of the novel, he desperately tries to reform, yet he has done too much evil to erase what is shown on the canvas. In the end, to cover up his immoral life, he slashes the canvas, killing himself in the process. Dorian's character arc serves as a cautionary tale about the dangers of unchecked hedonism and the corrupting influence of a life devoid of moral principles. It explores themes of vanity, morality, and the consequences of one's choices.

**Silas Marner, *Silas Marner*.** The character arc of Silas is a study of the power of society to bring an individual back from a place of despair and meaninglessness to the enjoyment of life [24]. After fleeing his former community due to false accusations, he settles in the village of Raveloe, where he lives a reclusive life centred around his work and hoarded gold. However, when his gold is stolen, he spirals into despair until a young girl named Eppie enters his life. Silas raises her as his daughter, leading to a transformation in his outlook and values. Through Eppie, he learns to love, engage with the community, and let go of his materialistic obsessions. This evolution culminates in his complete integration into the village and a newfound sense of purpose and belonging. The arc of Silas Marner underscores themes of redemption, the power of human connections, and the potential for personal growth.

**Ebenezer Scrooge, *A Christmas Carol*.** Ebenezer Scrooge embarks on a similar journey to Silas Marner over the course of *A Christmas Carol*. He transforms from a miserly and cold-hearted old man to a generous and compassionate individual. At first, Scrooge is known for his cruelty and lack of empathy, valuing money above all else. Through visits from three spirits, he witnesses scenes from his past, present, and future. These experiences help him realise the impact of his actions on others and the emptiness of his current life. Overcome with regret, he vows to change. On Christmas morning, Scrooge awakens with a renewed spirit, embracing kindness and generosity. He becomes a benefactor to those in need and forms meaningful relationships. His journey reflects themes of redemption, self-discovery, and the potential for change.

**Elizabeth Bennet, *Pride and Prejudice*.** Elizabeth Bennet's journey is much less dramatic than those of the other three protagonists. Instead of a transformation from good to evil or vice versa, her character arc is a much more nuanced journey of self-reflection and growth. As discussed in [15], initially, Elizabeth stands out for her intelligence and wit, which challenge the traditional roles expected of women in her time. Her feminism is evident in her refusal to conform to the social pressures of marrying solely for financial security. However, as she

navigates complex social interactions and encounters with Darcy, her perspective evolves. While maintaining her independence and self-respect, Elizabeth also comes to appreciate the value of love and companionship. Her feminist beliefs remain integral, as she seeks a partnership built on mutual respect and understanding rather than societal expectations. Elizabeth’s character arc underscores her journey towards a balanced view of feminism, intertwining personal agency with a recognition of the complexities of relationships and societal constraints. In summary, her character arc is multidimensional. True, it involves a personal transformation in terms of looking beyond her original prejudices of people, but it also involves influencing a change in others. This makes her arc much more complex than those of the other three discussed protagonists.

## 4.2 Evaluation of Results

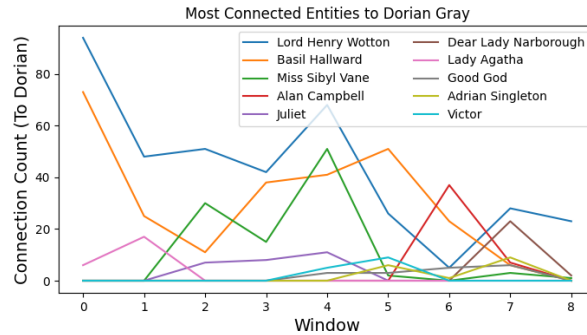
We believe that our quantitative results have accurately captured the arcs of Dorian Gray, Silas Marner, and Ebenezer Scrooge. Scrooge and Marner have very similar character arcs, whereas Gray’s is almost an inversion of the other two. This is quite clearly shown through each of the quantitative metrics defined in Section 3.2. On the other hand, our metrics do not seem to accurately capture the character arc of Elizabeth Bennet. This more nuanced and subtle character arc would require a much deeper investigation to quantify. Moreover, it is so different in concept from the other three that visualising this difference would likely require the study of other novels with characters similar to Elizabeth.

**The case of Silas Marner, Dorian Gray, and Ebenezer Scrooge.** In relation to these three protagonists, we can demonstrate the ability of our work to capture their character arcs. As mentioned in Section 4.1, we would expect the arcs of Silas Marner and Ebenezer Scrooge to be similar, while Dorian Gray’s should be almost the opposite. As shown in Figure 8, in terms of sentiment, our hypothesis seems to hold. In Figure 9, we can also see how the cosine similarity of the rolling average sentiment of mentions of Silas Marner and Ebenezer Scrooge is much higher than the similarity between these characters and Dorian Gray. This once again supports what was discussed in the expert analysis.

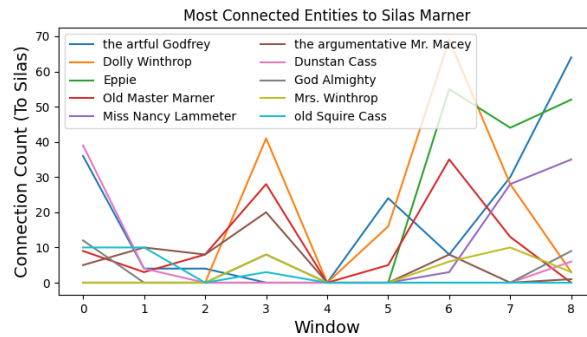
Notable are the journeys of Dorian Gray to become isolated and the embrace of society by Silas Marner. In Figure 13, we see that, as the story progresses, there is an increase in the connections between Silas and other characters in the novel. On the other hand, Figure 12 shows that Dorian becomes less connected to other characters. The death of certain characters, such as Sibyl Vane, likely contributed to this decline. However, the appearance of Alan Campbell in the top ten entities, who arrives late in the story and also dies, suggests that, as these characters either die or fade from Dorian’s social circle, no others come to replace them. Figure 3 presents an equivalent graph for Scrooge. In this case, there is not a dramatic change in connection to other characters<sup>6</sup>. This is likely

<sup>6</sup> The exception is window #4, where Scrooge visits the Cratchit house. However, the drop in the following chapter suggests that this is an outlier rather than a trend.





**Fig. 12.** Top-10 characters most closely connected to Dorian in *The Picture of Dorian Gray*, as plotted by connection count per window.

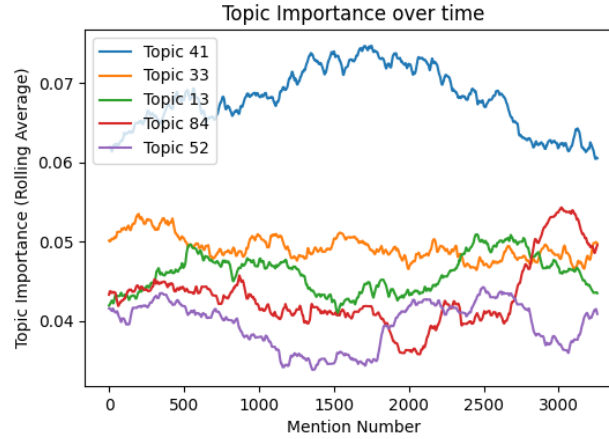


**Fig. 13.** Top-10 characters most closely connected to Silas in *Silas Marner*, as plotted by connection count per window.

due to the fact that the writing of *A Christmas Carol* is more focused on Scrooge than on his interactions with others. Scrooge's journey is more internal than that of Silas, whose arc relates more to how he becomes part of his community.

**The case of Elizabeth Bennet.** As observed in Section 4.1, the character arc of Elizabeth Bennet is much more complex and nuanced than that of our other three protagonists. We cannot simply define her arc in terms of sentiment and her connection to other characters, as the values of these metrics remain mostly static throughout the novel. Therefore, the last metric we have to look at is the shifting of topics.

Figure 4.2 shows the trends for a selection of key topics identified in *Pride and Prejudice* using LDA. We note the prevalence of Topic 41, which contains words such as society and words relating to time. This topic peaks around the middle of the novel and begins to decline after this. While it is difficult to definitely identify the reason for this change, it is likely caused by Elizabeth's discussions around societal norms with which she disagrees. The words for time could also be linked to pressure on female characters to marry as soon as possible to provide security



**Topic Top-10 Words:**

- Topic 41: time, feeling, evening, sort, minute, master, month, story, society, mere
- Topic 33: mind, hope, tone, doubt, cold, account, pity, drawing, sad, alter, treat, reproach, proper, shop, inform
- Topic 13: hear, moment, truth, glance, ready, ease, remark, spare, indifferent, finish, brush, hesitate, mix, 18th, 15th
- Topic 84: marry, wife, meet, marriage, laugh, party, forget, charm, dine, absurd, hospitality, scruple, honoured, 15th, heal
- Topic 52: speak, eye, care, heart, intend, guess, curious, heal, concerned, demean, disagreement, disrespectful, bounty, ladyship, occasional

**Fig. 14.** Topic importance over time in *Pride and Prejudice* for a selection of key topics, together with the top-10 words for these topics.

in their lives. We can also see a rise in Topic 84, which deals with positive words around marriage, likely linked to the marriage between Elizabeth and Darcy at the end of the novel. However, it is much more difficult to draw conclusions from the other topics, which certainly raises questions about the validity of looking at and drawing conclusions from topic trends within a novel. A more robust topic-modelling pipeline with a deeper level of analysis would give us a better understanding of other nuanced character arcs.

## 5 Conclusions and Future Work

In this dissertation, we have presented NLP methods for extracting, quantifying, and comparing character arcs from fictional texts. Our quantitative models for character arcs were defined by changes in character network connectivity, sentiment (both in dialogue and on a character mention level), and the importance of topics identified via topic modelling. For the novels *The Picture of Dorian Gray*, *Silas Marner*, and *A Christmas Carol*, our quantitative models derived from

character connectedness and sentiment analysis proved effective in tracking the transformations of the protagonists. We were also successful in comparing these arcs and our results aligned well with expert consensus drawn from the traditional literary study of these texts. However, our techniques proved less effective in tracking the arc of Elizabeth Bennet in *Pride and Prejudice*. In this novel, the characters do not undertake dramatic transformations (e.g. from good to evil, from isolated to embedded in society, or vice versa), meaning that techniques such as sentiment analysis and character connectedness are less effective. Our investigation of topic modelling showed some ability to reveal more complex arcs, although this approach suffers from issues related to interpretability.

Overall, this research has provided promising results. With further work, the understanding of any character arc from any novel through NLP might be achievable. As part of the completion of this dissertation, we have identified certain key areas of focus for future expansion of our research. First, although our topic modelling approach did not fully provide the solution to understanding more complex character arcs, with further research and development we believe that it could prove to be an essential part of understanding more nuanced character arcs. With more refined topics drawn from a wider corpus of novels, the topic trend charts presented within this dissertation could show a more understandable measure of how the concepts within a novel change as the plot progresses.

Using a large corpus might allow us to compare character arcs across different genres. The difficulties faced with extracting the arc of Elizabeth Bennet were certainly in part due to the lack of a similar character arc within our corpus. The other three arcs studied draw upon similar concepts, even if Dorian Gray's is travelling in the opposite direction. This made comparing them much simpler. This is further demonstrated in Figure 11, where we can see that these three arcs are much more similar to each other than to the arc of Elizabeth Bennet.

Finally, it would be interesting to apply these techniques to the character arc of a secondary character within a novel. Although not all characters in a novel have well-defined arcs, many texts, including *Pride and Prejudice*, contain a host of well-developed secondary characters, which might be considered in a further study. In general, we believe that we have been successful in answering the research questions we proposed. We have developed methods for quantitatively defining and comparing character arcs from fiction novels, we have validated our conclusions with comparisons to traditional literary studies, and we have presented future avenues of research for the expansion of our work.

**Acknowledgments.** I want to express my sincere thanks to my supervisor, Dr. Derek Greene, whose guidance, insights, and support were crucial in shaping this thesis. I am grateful for his expertise and encouragement throughout this research journey.

## References

1. Amalvy, A., Labatut, V., Dufour, R.: Bert meets d'artagnan: Data augmentation for robust character detection in novels. In: Workshop on Computational Methods

- in the Humanities 2022 (COMHUM 2022) (2022)
2. Arroyo-Barrigüete, J.L.: Sentiment analysis of lovecraft’s fiction writings. *Heliyon* **9**(1) (2023)
  3. Bamman, D., Popat, S., Shen, S.: An annotated dataset of literary entities. In: Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics. pp. 2138–2144 (2019)
  4. Bamman, D., Underwood, T., Smith, N.A.: A bayesian mixed effects model of literary character. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. pp. 370–379 (2014)
  5. Bell, J.S.: *Write great fiction-plot & structure*. Penguin (2004)
  6. Bhyravajjula, S., Narayan, U., Shrivastava, M.: Marcus: An event-centric nlp pipeline that generates character arcs from narratives. In: Text2Story@ ECIR. pp. 67–74 (2022)
  7. Blei, D.M.: Probabilistic topic models. *Commun. ACM* **55**(4), 77–84 (Apr 2012)
  8. Fischer, F., Göbel, M., Kampkaspar, D., Kittel, C., Trilcke, P.: Network dynamics, plot analysis: Approaching the progressive structuration of literary texts. In: DH (2017)
  9. Gerke, J.: Plot versus character: A balanced approach to writing great fiction (2010)
  10. Givon, S.: *Extracting information from fiction* (2006)
  11. Hutto, C., Gilbert, E.: Vader: A parsimonious rule-based model for sentiment analysis of social media text. In: Proc. ICWSM’14. vol. 8, pp. 216–225 (2014)
  12. Jacobs, A.M.: Sentiment analysis for words and fiction characters from the perspective of computational (neuro-) poetics. *Frontiers in Robotics and AI* **6**, 53 (2019)
  13. Kidd, C.E.: The uselessness of art: Critique and contradiction in the picture of dorian gray. *Interdisciplinary Journal of Undergraduate Research: Vol* **6**, 16 (2017)
  14. Liebman, S.W.: Character design in” the picture of dorian gray”. *Studies in the Novel* **31**(3), 296–316 (1999)
  15. Magtulis Cano, D.L.: A heroine of change and consolidation: Elizabeth benet: A harbinger of change in jane austen’s pride and prejudice (2022)
  16. Moretti, F.: *Graphs, maps, trees: abstract models for a literary history*. Verso (2005)
  17. Nadeau, D., Sekine, S.: A survey of named entity recognition and classification. *Linguisticae Investigationes* **30**(1), 3–26 (2007)
  18. Nalnick, E.T., Baird, H.S.: Character-to-character sentiment analysis in shakespeare’s plays. In: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). pp. 479–483 (2013)
  19. Stefan, J., Franzini, G., Cheema, M.F., Gerik, S., et al.: On close and distant reading in digital humanities: A survey and future challenges. In: Proc. EuroVis’15 (2015)
  20. Tabata, T.: Mapping dickens’s novels in a network of words, topics, and texts: Topic modelling a corpus of classic fiction. *JADH 2017* p. 73 (2017)
  21. Varner, R.M.: *Lady of the joust: Defining and classifying flat character arcs* (2023)
  22. Webster, K., Recasens, M., Axelrod, V., Baldrige, J.: Mind the gap: A balanced corpus of gendered ambiguous pronouns. *Transactions of the Association for Computational Linguistics* **6**, 605–617 (2018)
  23. Weiland, K.: *Creating character arcs: The masterful author’s guide to uniting story structure, plot, and character development*. (No Title) (2016)
  24. Williamson, L.A.: *A critical study of character-analysis in the novels of George Eliot*. Ph.D. thesis, University of British Columbia (1926)