

Informative and Discernible Visualisations: Empowering Literature Students through Visual Analyses of Classic Literature

Thomas J. Davidson
Corpus Christi College



**UNIVERSITY OF
CAMBRIDGE**

*A dissertation submitted to the University of Cambridge
in partial fulfilment of the requirements for the degree of
Master of Engineering in Computer Science*

University of Cambridge
Computer Laboratory
William Gates Building
15 JJ Thomson Avenue
Cambridge CB3 0FD
UNITED KINGDOM

Email: tjd45@cam.ac.uk

May 30, 2019

Declaration

I Thomas J. Davidson of Corpus Christi College, being a candidate for the M.Eng in Computer Science, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 11,999¹

Signed:

Date:

This dissertation is copyright ©2019 Thomas J. Davidson.

All trademarks used in this dissertation are hereby acknowledged.

¹Calculated using: `detex dissertation.tex | tr -cd "0-9A-Za-z \n" | wc -w`

Acknowledgements

I would like to acknowledge the contributions of my collaborating users: Ellie Jackson, Anna Bondarenko, Ellie Hennessey and Zoe Black for providing the driving motivation and direction behind the development of this project. I would also like to thank Tim Regan for his excellent supervision, Sean Rintel for his guidance on qualitative evaluation and Alan Blackwell for his support as my internal supervisor, as well as my Director of Studies, Martin Kleppmann. Additionally, Charlotte Burrows for her outstanding proofreading and general patience. For his help and tolerance with my understanding of statistics, I must also credit Jacob Bradley. Finally I would like to thank the peers, friends and strangers who took part in my experiment.

Abstract

This report describes novel techniques for the visualisation of English literature. By building on previous approaches to abstract text visualisation I have developed informative and visually appealing images from text. The tool and visualisations enable a user to perform exploratory visual analysis of text and to generate narrative analyses to convey their findings to others. The visualisations and tool were evaluated using three different evaluative techniques. The final visualisations are promising supplementary analytical tools for literature students as well as stand alone analyses. The potential for usage extends beyond the scope of academia into the realms of visual art and the visualisation of other media.

Total word count: 11,999¹

¹Calculated using: `detex dissertation.tex | tr -cd "0-9A-Za-z \n" | wc -w`

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Goals	2
1.3	Overview of Approach	3
2	Background and Related Work	5
2.1	Visualisations as Expressive Art Forms	5
2.2	Analytical Tools and Visualisations	6
3	Design and Implementation	9
3.1	Recruitment of Collaborative Users	9
3.2	Design	10
3.2.1	Timeline	10
3.2.2	Initial Interviews	11
3.3	Selection of Literature	14
3.4	Selection of Visualisation Features - Authorial vs. Editorial Decisions	14
3.4.1	Structural interrogation	16
3.4.2	Character interaction	16
3.4.3	Location interaction	16
3.5	Software Design	17
3.5.1	Programming Language Selection	17
3.5.2	Backend	18
3.5.3	Frontend	18
3.6	Implementation	20
3.6.1	Iteration 1	20
3.6.2	Iteration 2	27
3.6.3	Iteration 3 - Final Modifications	36
3.7	“Finished Tool”	39

4	Evaluation	43
4.1	Evaluation Plan	43
4.2	Experiment	44
4.2.1	Design and Method	44
4.2.2	Results	47
4.2.3	Think-Aloud Walkthrough	51
4.2.4	Conclusion	54
4.3	Evaluative Interviews	55
4.3.1	Results	55
4.3.2	Conclusion	63
4.4	Think-Aloud Study	64
4.4.1	Usability	66
4.4.2	Conclusion	68
5	Conclusions	69
5.1	Improvements	69
5.2	Successes	69
5.3	Future Work	71
A		73
A.1	Initial Interview Responses	73
A.2	End of Iteration 1 Interview Responses	75
A.3	End of Iteration 2 Interview Responses	77
A.4	Overview of Literature	78
A.4.1	Alice’s Adventures in Wonderland (1865) by Lewis Carroll	78
A.4.2	The Story of Doctor Dolittle (1920) by Hugh Lofting	79
A.4.3	Great Expectations (1861) by Charles Dickens	79
A.4.4	Strange Case of Dr. Jekyll and Mr. Hyde (1886) by Robert Louis Stevenson	79
A.4.5	Pride and Prejudice (1818) by Jane Austen	79
A.4.6	The War of the Worlds (1897) by H.G. Wells	80
A.5	Progression of Visualisation in Cognitive Walkthrough	81
A.6	Final Summative Interviews	85
A.6.1	Questions	85
A.6.2	Transcriptions	87

List of Figures

2.1	A visualisation of 1984 by George Orwell from Stefanie Posavec’s <i>Writing Without Words</i> project.	6
3.1	Gantt chart of planned progress through project.	10
3.2	A chord diagram of character co-occurrences in <i>The Hobbit</i> [9].	11
3.3	A linear timeline of sentiment in <i>The Hobbit</i> [9].	11
3.4	A word cloud generated from the Bible [21].	12
3.5	An example of how the bounding boxes surrounding the location names will be generated.	17
3.6	The initial plan for the rough design of the software.	18
3.7	An abstracted overview of the data ingestion process. Green indicates fully automated, orange – some manual input required.	19
3.8	The initial database schema.	19
3.9	An excerpt from Chapter 2 of <i>Strange Case of Dr. Jekyll and Mr. Hyde</i> which was correctly parsed as a single sentence	21
3.10	The initial design for an ingestion system. The orange highlighting signifies external classes or packages.	22
3.11	An excerpt of a long thin image showing sentence length, created from <i>The War of the Worlds</i>	23
3.12	The first 350 sentences of <i>Pride and Prejudice</i> with Elizabeth (green) and Mr. Darcy’s (red) occurrences highlighted.	24
3.13	An excerpt of a visualisation of <i>Alice’s Adventures in Wonderland</i> with the garden (purple) and pool of tears (yellow) location-boxes highlighted.	26
3.14	A visualisation of the whole text of <i>Strange Case of Dr. Jekyll and Mr. Hyde</i>	28
3.15	The whole text of <i>Strange Case of Dr. Jekyll and Mr. Hyde</i> with the characters Utterson (blue), Jekyll (red) and Hyde (green) highlighted.	28

3.16	The whole text of <i>Strange Case of Dr. Jekyll and Mr. Hyde</i> with margins being used and the characters Utterson (blue), Jekyll (red) and Hyde (green) highlighted.	29
3.17	Example using the word ‘abomination’ of the effect of the paucity of data before 1800 on the results from Google’s Ngram viewer.	31
3.18	The colour spectrum assigning each year a colour.	32
3.19	The left half of this visualisation of <i>The War of the Worlds</i> is generated using the granular dot approach and the right half uses the aggregated, per-sentence approach.	33
3.20	<i>The Story of Doctor Dolittle</i> with the chapters highlighted. . .	33
3.21	<i>Alice’s Adventures in Wonderland</i> with sentence 480, selected through the interface by clicking on it, highlighted in white. .	34
3.22	Chapter 6 of <i>Alice’s Adventures in Wonderland</i> with Alice (pink), the Queen (cyan) and the Hatter (magenta) highlighted in the margin. The granular version of the year information is highlighted in situ.	35
3.23	The updated colour spectrum assigning each year a colour. . .	38
3.24	Part of the adjacent colour scheme for character-highlighting. There are more than 3 characters possible, but this is the default.	39
3.25	Part of the triadic colour scheme for location-highlighting. There are more than 2 locations possible, but this is the default.	39
3.26	Example visualisation of <i>Great Expectations</i> with default parameters set and the following information - aggregated year-highlighting, chapter-highlighting, character-highlighting with lines (Joe (cyan), Estella (purple), Havisham (magenta)), location-highlighting (London (green), Blue Boar Inn (blue)).	40
3.27	The final structure of the system – arrows represent data flow. Green headings are python scripts, orange headings indicate external libraries. There are elements missing from this diagram, but it gives an overview.	41
4.1	Example question from section 3 of the questionnaire asking the user to identify the book from the image (the first 350 sentences of the book with the main characters highlighted) - the correct answer is <i>The War of the Worlds</i>	46
4.2	Example visualisation of <i>Great Expectations</i> showing the appearances of Joe (cyan) and Estella (magenta) throughout the text.	57

4.3	Example visualisation of <i>Alice's Adventures in Wonderland</i> showing the episodic interactions of Alice (cyan) with the Hatter (magenta) and the Queen (purple), as well as the recurrence of location through the <i>pool of tears</i> (blue) and the <i>garden</i> (green).	58
4.4	The image on the left is <i>Pride and Prejudice</i> , published in 1818, the image on the right is <i>Great Expectations</i> , published in 1861. Despite only slightly more than 50 years difference in publication date, there is a clear visual difference between the two images.	59
4.5	Example of physical annotations from user E's studies.	61
4.6	The reference sheet used in the online form.	62
4.7	Visualisation of <i>Great Expectation</i> generated by user E in the cognitive walkthrough with default parameters set and the following information: aggregated year-highlighting, chapter-highlighting, character-highlighting with lines (Pip (cyan), Herbert (purple), Estella (magenta), Joe (green), Jaggers (red), Magwitch (yellow), Orlick (orange)); location-highlighting (London (green), Satis House (blue), Blue Boar Inn (orange)).	67
A.1	First visualisation created by the user highlighting characters Pip (cyan), Estella (magenta) and Herbert (purple).	81
A.2	Second visualisation created by the user highlighting characters Pip (cyan), Estella (magenta) and Herbert (purple); locations London (green), Satis House (blue) and Blue Boar Inn (orange).	82
A.3	Third visualisation created by the user highlighting the same as above and turning on chapter highlighting.	83
A.4	Fourth visualisation created by the user highlighting the same as previously and attempting to highlight POS through personal pronouns (green) and comparative adjectives (orange).	84
A.5	Fifth visualisation, removing POS-highlighting and adding aggregated year-highlighting.	85

List of Tables

3.1	Table of codified responses from initial interview with users' views on chord diagrams. Bracketed letter represents the user who made the comment.	12
3.2	Table of codified responses from initial interview with users' views on novel idea of character highlighting.	13
3.3	Summarising table of key information about the selected books.	15
3.4	Table of codified responses from interview with users' views on character-highlighting approach.	26
3.5	Table of codified responses from interview with users' views on year-highlighting approach.	36
3.6	Description and default values for the core changeable parameters of the visualisations.	38
4.1	Table of percentage of correct answers for the different sections, both for selecting the book name from the image and for selecting the generated image from a given book. All results are given to two significant figures. Underlined, italicised results indicate significance and are listed with their p-values.	48
4.2	Table of percentage of correct answers for the different sections relative to the minimum familiarity with the book, both for selecting the book name from the image and for selecting the generated image from a given book. All results are given to two significant figures. Underlined, italicised results indicate significance and are listed with their p-values.	50
A.1	Table of codified responses from initial interview with users' views on linear timelines.	73
A.2	Table of codified responses from initial interview with users' views on word clouds.	74
A.3	Table of codified responses from initial interview with users' views on novel idea of structural interrogation.	74

A.4	Table of codified responses from initial interview with users' views on novel idea of location-highlighting.	74
A.5	Table of codified responses from interview with users' views on structural interrogation presentation.	75
A.6	Table of codified responses from interview with users' views on location-highlighting approach.	75
A.7	Table of codified responses from interview with users' views on the potential POS-highlighting approach.	76
A.8	Table of codified responses from interview with users' views on the potential year-highlighting approach.	76
A.9	Table of codified responses from interview with users' views on vertical vs. horizontal presentation.	77
A.10	Table of codified responses from interview with users' views on lines connecting character mentions together.	77
A.11	Table of codified responses from interview with users' views on the margin approach.	78

Chapter 1

Introduction

1.1 Motivation

Analysis and interpretation of large bodies of text can be time-consuming and laborious. Previous work has shown that students learn in different ways, such as visually and verbally [1]. The initial motivation for this project grew out of frustration at a lack of ability to explore and present literature in a more visual manner.

Through discussions with my peers I discovered these frustrations were shared by many and new issues were brought to my attention. These introductory interviews are covered in depth in **3.2.2** but can be summarised with these three statements:

1. “I want to be able to find interesting parts of the text without trawling through masses of text.”¹

This was the most prevalent opinion and reflected the difficulty of searching large quantities of text to find interesting passages to study. This is not a counter-argument to the benefits of interaction with physical media, such as those put forward in *The Myth of the Paperless Office* [2], but an observation of the difficulty of assimilating large amounts of information. Whilst tech-

¹From a Classics undergraduate.

nology makes this easier by allowing users to search documents, this does not alleviate the problem of moving back and forth in a document. This problem leads to the second statement.

2. “I want to be able to analyse the text as a whole.”²

Participants felt constricted by the nature of written literature and the inability to analyse the work as a whole. It requires a significant amount of effort and time to obtain a basic understanding of a full book.

3. “I want to create visual art from literature.”³

Some of those interviewed discussed their desire for creating visually stimulating work from written material.

From this preliminary investigation I set out to create a visualisation tool for literature that would enable the creation of stimulating, engaging and informative images from English literature.

1.2 Goals

The aim of this project was to develop a fully usable software tool for generating informative visualisations from text files of literature. These visualisations should allow for *exploratory analysis*⁴ and the generation of *narrative analyses*.⁵ I refined this into three main categories: *data ingestion*, *visualisations*, and *user interface*. *Data ingestion* would deal with converting raw text into a usable form. *Visualisations* would generate the images from this data, and *user interface* would allow the users to interact and generate images.

As I discuss in **3.6.1** developing interesting visualisations was more important than developing a user interface (UI). This is not to trivialise the work that goes into developing a good UI, but to emphasise that the purpose of this project was to generate useful visualisations. It is an aim to develop a UI

²From a student studying for a Masters in English Literature.

³From a student studying for a Masters in Creative Writing.

⁴Investigating the images to discover trends – see **2.2**.

⁵Presenting information from the data to other people – see **2.2**.

for the tool in the future, hence the software was designed to facilitate the addition of one.

1.3 Overview of Approach

Throughout the project I utilised a *collaborative design* approach [3]. This means I worked towards a common goal with *collaborating users* (CUs) who drove the development of the tool at every stage – giving feedback on prototypes that was incorporated into the next iteration.

After recruiting users I carried out preliminary interviews, formulated a design plan and selected literature to visualise. I then worked through three iterations of development, progressing the areas the CUs reacted positively to further.

I undertook extensive evaluation of the visualisations through evaluative interviews, an online experiment and a think-aloud study [4]. This combination of evaluative techniques allowed me to determine the success of the project in achieving its aim of generating informative novel text visualisations facilitating explorative analysis and the generation of narrative analyses.

Chapter 2

Background and Related Work

The field of research that this project covers is described as (semi) abstract text visualisation [5]. This field creates images by directly querying the text of published work. Previous work in this area can be split roughly into two categories: work which produces expressive art, and work that attempts to create useful analytical tools. Sometimes projects which set out to do the latter end up being primarily used for the former. Brad Paley's TextArc [6] is a screen-based application which takes text and displays it around the edge of a large ellipse using a word-cloud approach and a novel anchoring technique. This was intended to create an exploratory tool for academics analysing literature, however, one of the most notable outcomes was creating large-scale prints as visual art. The project I have proposed spans across these two categories and hopes to create an effective analytic tool using text visualisation whilst also producing visually appealing images.

2.1 Visualisations as Expressive Art Forms

Work such as Stefanie Posavec's *Writing Without Words* [7] is an example of a piece of physical art generated from text. Posavec does not utilise computer driven text-mining but her techniques could be automated. One set of images

is generated by drawing a line equal to the length of a sentence and turning 90 degrees at the end of every sentence (figure 2.1). The images that are created are visually appealing, but do not convey immediately accessible information about the books they are generated from. A further example of this kind of work is Boris Müller’s visualisations of poems for the annual *Poetry on the Road* literature festival [8]. The images he creates here are captivating, but so far abstracted from the original text that it is difficult to garner useful information. This is a unifying feature of many artistic visualisations.

NINETEEN EIGHTY-FOUR
George Orwell
1949

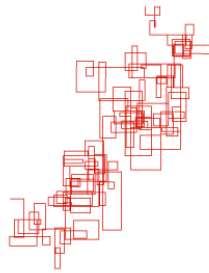


Figure 2.1: A visualisation of 1984 by George Orwell from Stefanie Posavec’s *Writing Without Words* project.

2.2 Analytical Tools and Visualisations

This work is often focussed on generating infographic-style visualisations – such as chord diagrams (figure 3.2) – which are useful for providing insight into the text but sever this information from any context. This can be seen in a 2016 project by Natalia Bilenko [9]. She uses chord diagrams to show the frequency of co-occurrence of characters in *The Hobbit*. This approach is

augmented with the addition of a timeline of sentiments (figure **3.3**) for each sentence throughout the book, relevantly highlighted when a character is selected, thus, adding context to this visualisation. The way it is portrayed, though, makes it difficult to understand how the frequency of co-occurrence fits into the context of the whole book.

A shortcoming of some of the work in this area is the development of tools to visualise specific literature. This is seen in the previously mentioned TextArc[6] where Paley focussed on the book *Alice's Adventures in Wonderland*. Using one book to focus on is a valuable approach, allowing the developer to have in-depth knowledge of that book and successfully design visualisations to highlight key features in it. However, this has the danger of the developer creating a narrative analysis tool when they set out to develop an exploratory tool. Narrative analysis refers to analysis that tells a story to the viewer [10]. In this context, the story being told is determined by the visualisation you choose to create as the developer. In contrast, an exploratory tool allows users to find a 'story' themselves [11]. TextArc does focus heavily on a single text, but it is not necessarily purely a narrative analysis tool. Visualisations such as Larkin's historic Dispensational Charts [12], however, are, as they have a narrow focus of selected literature, and a specific target with their narrative so cannot be taken and applied to other texts.

Ben Fry's visualisation of the evolution of Darwin's ideas [13] is a visualisation which maintains tight association with its context. It highlights the updates of revisions of Darwin's *On the Origin of Species* [14]. This is an example of a narrative analytical tool with exploratory scope – it could also be applied to other work where multiple editions of the same book are available.

Tim Regan's project visualising *His Dark Materials* by Philip Pullman [15] focussed on generating two forms of visualisation: *Whole Text*¹ and *Character Flowers*.² The *Whole Text* approach I found most appealing as it directly

¹Visualisations of the entire text of the book printed very small.

²Similar to chord diagrams, a visualisation highlighting frequently co-occurring words by arranging in a 'petal' arrangement at varying distances around the character name.

utilised the text of the book and superimposed character occurrences onto it. This project was focussed specifically towards this series but had the potential for application to further books and for exploratory analysis through altering the names of selected characters. The main downside of the visualisations created in this project were that they were designed to be printed in a large format. This makes for visually engaging and informative images, but makes it harder to justify the visualisations as a legitimate aid for academic analysis of text.

Martin Wattenberg's 2007 paper utilising 'chromograms' to visualise activity on Wikipedia visualises edit activity by assigning a colouring system to text sequences [16]. This is an example of an exploratory visualisation tool which is then used to narratively display various features which become apparent through his framework.

More recent work has looked at analysing and visualising the readability of paragraphs in *Harry Potter* [17]. This project does tie parts of its visualisation to the context of the book, but in a way that is sometimes visually confusing. In latter stages they sever the connection between the actual text and move towards aggregative analysis of the data. In addition to this, the measures they are attempting to visualise (readability and sentiment) are inherently subjective. In general this project, and others, tries to construct too much of a narrative analysis and obfuscates the authorial decisions behind layers of subjective analysis. This makes it hard for a user to analyse the text themselves using their tool. Another piece of more recent work has seen text visualisation approaches applied to social media activity [18]. The work I create here could be applied in this sphere too.

I aim to overcome the issues raised here and generate visually stimulating and informative visualisations through a tool which allows users to perform exploratory analysis of text. By building directly off the author's raw text I hope to tie the visualisations to the context of the physical book and remove layers of obfuscation. I will achieve this by building on previous research and art that is mentioned here, and by developing novel visualisation techniques.

Chapter 3

Design and Implementation

3.1 Recruitment of Collaborative Users

As the project was focussed on visualising English literature, I recruited students studying English as my main users. In total, five users contributed to the project at varying times. The users were recruited through my academic network and I sought to engage with a variety of people within the broader region of students for whom studying literature was key.

- User A: 22 year old English Literature Masters student at University of Bristol - BA in English from Cambridge University
- User B: 22 year old Creative Writing Masters student at University of Edinburgh - BA in English from Cambridge University
- User C: 21 year old Undergraduate Classics student at Cambridge University
- User D: 22 year old MPhil student in Classics at University of Cambridge - BA in Classics from Cambridge University
- User E: 21 year old Undergraduate English student at Cambridge University

User's A and B were the main CUs and were present throughout the whole process. User C was instrumental in the earlier stages of the development but reduced their commitment later in the project. User's D and E were contributors during the latter stages.

3.2 Design

The design and implementation of the visualisations was driven by feedback that I obtained from the CUs. In total, I conducted four rounds of interviews. The first of these consisted of pre-development interviews to shape the direction of development. The second two rounds of interviews were structured as iterative development processes [19] where I acted on the feedback from users and presented them with progress, before repeating. The final round of interviews identified potential avenues for further development of the tool.

3.2.1 Timeline

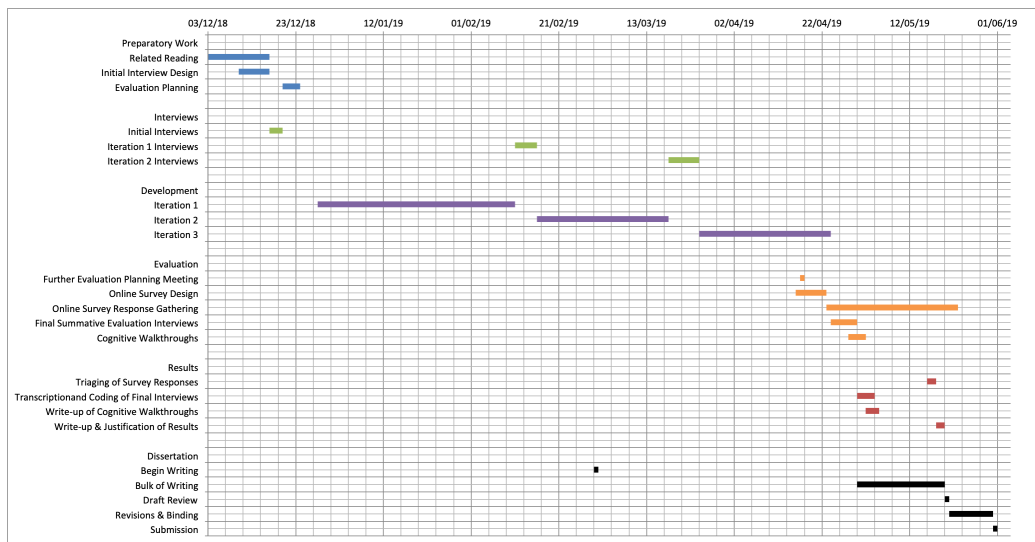


Figure 3.1: Gantt chart of planned progress through project.

3.2.2 Initial Interviews

The initial interviews were one-on-one, split into two parts. Firstly, exploring the interviewee’s opinions on previous approaches to visualising literature, before suggesting more novel approaches and asking for their potential ideas. In assessing the past approaches to visualisations I identified three main *traditional* techniques; chord diagrams (Figure 3.2), linear timelines (Figure 3.3) and word clouds (Figure 3.4). In order to analyse the users’ responses to these approaches I codified their responses [20]. Table 3.1 shows the results for chord diagrams and the other responses can be found in tables A.1 and A.2.

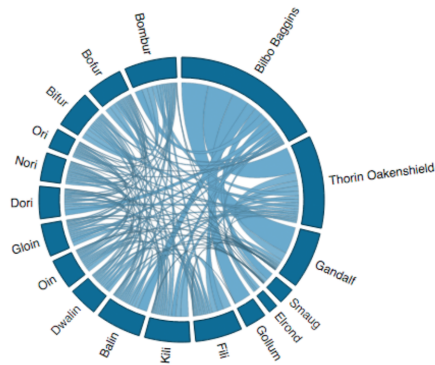


Figure 3.2: A chord diagram of character co-occurrences in *The Hobbit* [9].

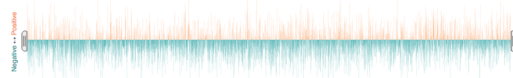


Figure 3.3: A linear timeline of sentiment in *The Hobbit* [9].

Positive	Negative
Good to be able to see characters that always occur together. (C)	Not easy to see how it fits into story. Importance of character co-occurrence to other aspects of the story. For example when Petrochilus and Achilles are not together in Iliad, something is wrong. (C)
Good for general view. (B)	Not clear. (A)
	Seems more artistic. (A)
	Not easy to use. (A)
	Raw data with no context. (A)
	Can't see kind of relationship. (B)
	Can't see where in story they interact. (B)

Table 3.1: Table of codified responses from initial interview with users' views on chord diagrams. Bracketed letter represents the user who made the comment.

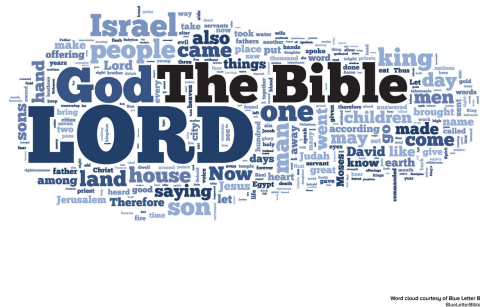


Figure 3.4: A word cloud generated from the Bible [21].

The users repeatedly reported that historic approaches to text visualisation did offer interesting insights into the text, however, the lack of context made them hard to interpret. Additionally they were frustrated by the lack of visual appeal of some of the images. With this in mind I developed initial ideas of more novel approaches to visualising literature. I built on previous work concerning character interaction [15] and more artistic work in this

Positive	Negative
Interesting to see how characters actually occur compared to your perception. (C)	Authorial voice might affect visibility of a character. (C)
Useful to be able to see in context (A)	Can't see if a character is being mentioned lots but their name isn't. (C)
Could use to see use of a character as a plot progression tool (C)	
Gives you context - where they occur and co-occur in the story. (B)	
Interesting different visual representation from the textual imagery in the book. (B)	

Table 3.2: Table of codified responses from initial interview with users' views on novel idea of character highlighting.

area [7]. Based on the users' comments I wanted to centre my visualisations around the physical shape of the text to allow the images to have direct context. The approaches initially considered were:

- Structural Interrogation:- The idea of directly displaying sentential length
- Character Interaction:- Highlighting occurrences of key characters names throughout the text
- Location Interaction:- Showing how different locations are mentioned in the text

I sought the CUs' opinions on the proposed techniques and codified these responses. The results for character interaction are shown in table **3.2** and the others in tables **A.3** and **A.4**. The nature of this interview section was slightly leading as the users were aware that these were my propositions and were not forthcoming with negative responses. The lack of negative should be viewed with this caveat. At later points in this dissertation I more directly address their concerns and issues with each of the visualisation features.

The feedback from the users towards the approaches was almost entirely

positive. There was a positive consensus towards the provision of context for the visualisations using the structural interrogation. The users also offered suggestions of investigating the age of words (**3.4**) and using NLP techniques to ascertain the parts of speech of words.

3.3 Selection of Literature

The criteria for selection of literature were dictated by necessity and my users' views. Firstly, the book should be out of copyright so that the book is accessible and there are no legal issues with redistribution of the text. Secondly, my users should be familiar with the book to enable me to perform thorough evaluation of the software. Additionally, books were selected that had been quite widely studied to increase the effectiveness of testing. I focussed on books that had been on exam syllabi in the UK. Finally, the users expressed interest in observing the effect of authorial voice¹ as well as a range of publication dates.

Eventually I settled on six books which were easily accessible from Project Gutenberg.² They are listed in table **3.3** with key details. An overview of each story is provided in **A.4**.

3.4 Selection of Visualisation Features - Authorial vs. Editorial Decisions

Another outcome of the initial interviews was the importance of focussing on authorial decisions over editorial decisions – by this, I mean decisions made by the writer, and not by the publisher or editor. A good example of where this can be difficult is in the structural interrogation visualisation feature.

¹First and third person.

²<https://www.gutenberg.org> – a volunteer-run website to digitise and archive cultural work.

Title	Author	Year Published	Key Characters	Key Locations	Authorial Voice
Alice's Adventures in Wonderland	Lewis Carroll	1865	Alice, Queen, Hatter	Garden, Pool of Tears	3rd Person
The Story of Doctor Dolittle	Hugh Lofting	1920	Doctor, Polynesia, Jip	Africa, Puddleby	3rd Person
Great Expectations	Charles Dickens	1861	Joe, Estella, Havisham	London, Boar Inn	1st Person
Strange Case of Dr. Jekyll and Mr. Hyde	Robert Louis Stevenson	1886	Utterson, Jekyll, Hyde	London, Soho	3rd Person
Pride and Prejudice	Jane Austen	1818	Elizabeth, Darcy, Bingley	Longbourn, Netherfield	3rd Person
The War of the Worlds	H.G. Wells	1897	Curate, Ogilvy, Henderson	Woking, Mars	1st Person

Table 3.3: Summarising table of key information about the selected books.

Instead of focussing on the line length – an editorial decision – I focussed on the sentence length, which – whilst potentially influenced by the editorial process – is a clearer authorial decision.

The initial three visualisation features I progressed with were the ones presented to the CUs in the initial interview: structural interrogation, character interaction and location interaction.

3.4.1 Structural interrogation

This took the form of generating a line for every sentence where the line's length was equal to the number of words in the sentence. This is a similar concept to the approach in Figure 2.1 and is a simple visualisation, though requires an effective way of splitting the text into individual sentences.

3.4.2 Character interaction

Combining with structural interrogation and building closely on Tim Regan's work[15] the aim of this approach was to highlight every occurrence of the main character(s) name, so that the section of the sentence where they appeared was highlighted. Different colours would distinguish different characters. Additionally, joining up each occurrence with a line could highlight the character's 'journey' through the book and address the desire of the users to see characters as plot progression tools.

3.4.3 Location interaction

This feature was a novel approach to visualising how the main locations of the book are mentioned by highlighting a bounding box around contiguous mentions of a specific place name without intervening location mentions. For a visual understanding of this see Figure 3.5.

Sentence 1: The man arrives at Leckhampton.
 Sentence 2: He wanders through the gardens there.
 Sentence 3: After a while, he leaves Leckhampton and cycles to Trumpington.
 Sentence 4: Half way along the road, he realises he has left his keys in the gardens at Leckhampton.
 Sentence 5: He quickly rushes back and collects them.
 Sentence 6: Now back on the road to Trumpington he is almost hit by a car.
 Sentence 7: In all the confusion he drops his wallet.
 Sentence 8: Several hours later he finally arrives in Trumpington.

Location	Start Sentence	Finish Sentence
Leckhampton	1	3
	4	4
Trumpington	3	3
	6	8

Figure 3.5: An example of how the bounding boxes surrounding the location names will be generated.

3.5 Software Design

The tool would consist of two distinct parts. A backend to ingest and store data from the text and a frontend to generate visualisations by querying stored data.

3.5.1 Programming Language Selection

I elected to use Java for the majority of this project. The primary motivation was my familiarity and proficiency meaning I could focus on genuine problems when implementing. I considered using a more specialist visualisation language such as D3.js³ but concluded that as the visualisations I would generate were reasonably simple in structure this would add unnecessary complexity. Moreover, the approach I took to development meant that it would be straightforward to use more complex visualisation tools in the future. Where appropriate throughout the project I used other languages – python for a web-scraping script, and SQL for database interaction.

³<https://d3js.org>

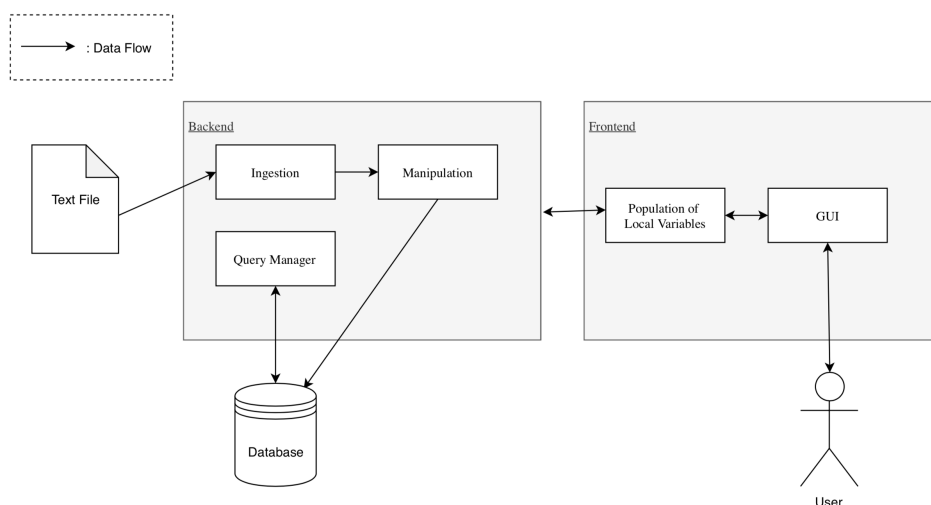


Figure 3.6: The initial plan for the rough design of the software.

3.5.2 Backend

The backend of the system is responsible for ingesting raw text of the book and storing this in a query-able form. This process was semi-automated, however some human intervention steps were initially essential. The workflow for this is detailed in figure 3.7. After downloading the plain-text from Gutenberg I manually stripped the surrounding text detailing distribution rules. Then NLP operations are performed on the text, splitting the text into its constituent sentences, as well as POS-tagging and stemming. This data is stored in an SQL database. Thus, with minimal human input the information in figure 3.8 for every word in each of the books could be stored.

3.5.3 Frontend

The frontend generates visualisations by querying stored data and editing local variables before displaying the image in a JFrame. Interactions with the frontend allow a user to dynamically change the image in real-time, gathering new data from the database.

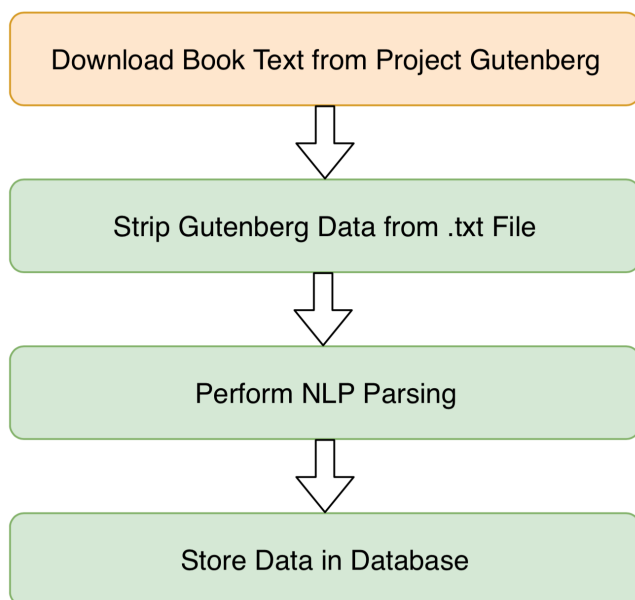


Figure 3.7: An abstracted overview of the data ingestion process. Green indicates fully automated, orange – some manual input required.

SENT_ID (INT)	WORD_ID (INT)	WORD (VARCHAR 25)	STEM (VARCHAR 25)	POS (VARCHAR 10)
What number sentence in the book the token comes from	The position of the token in the sentence	The token as it appears in the text	The stem of the token as returned by NLP techniques	The part of speech of the token as returned by the NLP techniques

Figure 3.8: The initial database schema.

3.6 Implementation

The main bulk of implementation was undertaken in three iterations with further interviews conducted at the end of the first two iterations. The discussions in these interviews formed the basis of alterations to the program and various additions to it.

3.6.1 Iteration 1

Development

Data Ingestion I began by focussing on developing a backend for ingesting raw text into the desired database form. I wrote custom functions utilising Java’s built-in string management functions to strip unnecessary information from the start and end of the .txt file. Then, using Stanford’s NLP package,⁴ I split the text into sentences using the *Document Preprocessor* class. The Stanford NLP toolkit is a frequently referenced toolkit and so I chose not to carry out significant testing into its sentence splitting [22]. As anecdotal evidence, figure 3.9 is a complex sentence successfully extracted from *Strange Case of Dr. Jekyll and Mr. Hyde*.

I then used Stanford’s *Max Ent Tagger* class to tag the POS and used a Porter Stemmer [23] to stem⁵ each of the words. This processed data was inputted into a MySQL database⁶ using a JDBC driver⁷ and a custom database management class. This ingestion system was encapsulated into a number of Java classes – both external, and custom-written – whose interactions can be seen in figure 3.10.

⁴<https://nlp.stanford.edu>

⁵The stem of a word refers to its most basic form without suffixes and other embellishments.

⁶<https://www.mysql.com>

⁷<https://www.mysql.com/products/connector/>

The will was holograph, for Mr. Utterson though he took charge of it now that it was made, had refused to lend the least assistance in the making of it; it provided not only that, in case of the decease of Henry Jekyll, M.D., D.C.L., L.L.D., F.R.S., etc., all his possessions were to pass into the hands of his “friend and benefactor Edward Hyde,” but that in case of Dr. Jekyll’s “disappearance or unexplained absence for any period exceeding three calendar months,” the said Edward Hyde should step into the said Henry Jekyll’s shoes without further delay and free from any burthen or obligation beyond the payment of a few small sums to the members of the doctor’s household.

Figure 3.9: An excerpt from Chapter 2 of *Strange Case of Dr. Jekyll and Mr. Hyde* which was correctly parsed as a single sentence

Structural Interrogation Whilst the backend benefitted from heavy use of external libraries, due to the novelty of the visualisations, almost all of the front-end modules were custom-written. The front-end consisted of a large class which made several calls to an instance of the database interaction class. Within the front-end class I developed methods to interpret data stored in the database and facilitate its display on screen. I used the concept of ‘database’ and ‘image’ co-ordinates for every word in the book. The database co-ordinates are the sentence- and word-id of each word. The image co-ordinates are the x and y co-ordinates of that word in the image. The key behind the interaction of the backend data with front-end data was developing methods which converted between image and database coordinates (listing 3.1).

I used array-lists when working with data-points due to their ability to dynamically add and remove items. In the initialisation methods of the visualisation class several lists are populated, most importantly the *activated* list, which stores the database co-ordinates of every word in the selected book. The *basePaint* method paints the basic structural interrogation aspect of the visualisation. This takes the *activated* list, converts them to image co-ordinates, and prints them in a Java Swing JPanel as rectangles

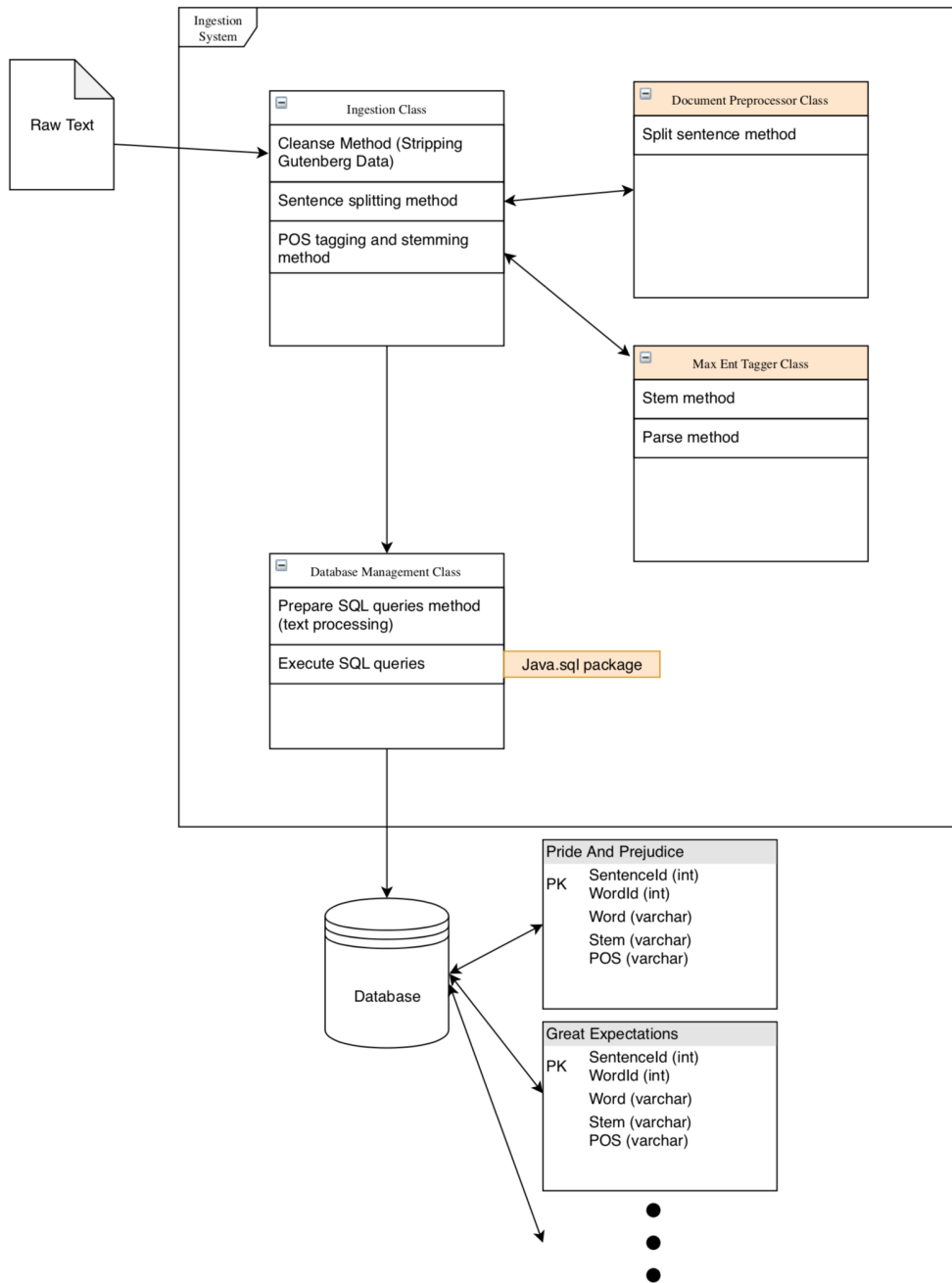


Figure 3.10: The initial design for an ingestion system. The orange highlighting signifies external classes or packages.

onto a background. This simple method initially allowed for the creation of long thin visualisations which showed the varying sentence length throughout a book.

```
589     private static Point2D convertToImageCoord(Point2D p){
590         int oldY = (int)p.getY();
591         int oldX = (int)p.getX();
592
593         int newY = oldY*pixelHeight;
594         int newX = oldX*pixelWidth;
595
596         return new Point2D.Double(newX, newY);
597     }
```

Listing 3.1: The method for converting from database co-ordinates to image co-ordinates.



Figure 3.11: An excerpt of a long thin image showing sentence length, created from *The War of the Worlds*.

Character Highlighting For the character-highlighting the *visualisation* class called various methods in the *database interaction* class to populate a new list of database co-ordinates representing where the character's name

appeared in the book. This is done using an SQL query to the database, such as listing 3.2. Then by selecting a different colour for each character it is possible to overlay this information on the base painted image and show every occurrence of the desired characters' names (figure 3.12).

```
SELECT sentid,wordid FROM Pride_And_Prejudice WHERE  
word = 'Elizabeth';
```

Listing 3.2: Example SQL query for returning database co-ordinates of mentions of the name Elizabeth from *Pride and Prejudice*.

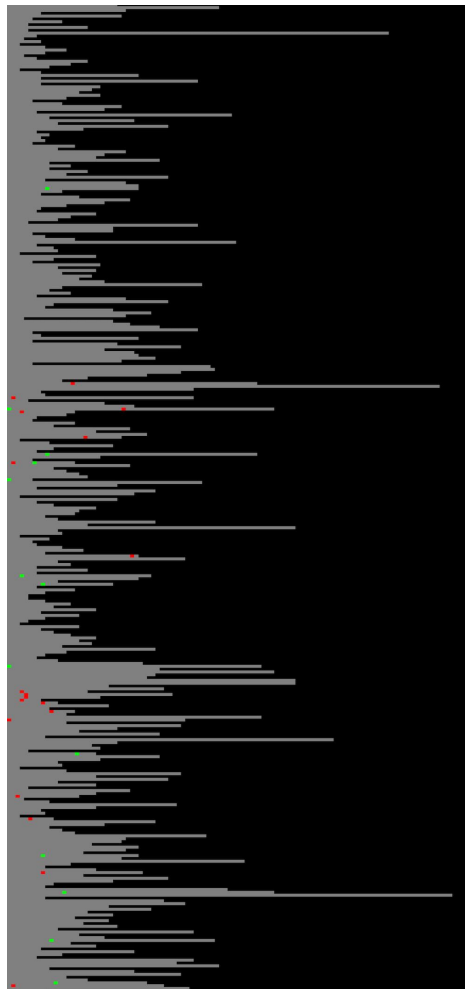


Figure 3.12: The first 350 sentences of *Pride and Prejudice* with Elizabeth (green) and Mr. Darcy's (red) occurrences highlighted.

Location Highlighting Developing location-highlighting required more initial steps. First, I added an additional database table to store place names referenced in the book⁸. I then wrote methods to populate a custom data structure I called a *location-box* to store information about the occurrence of each location. Each location box stored the first and last sentence where a location was mentioned with no other location mentioned in the sentences between.

On initialisation of a visualisation object every location in that book had a list of location-boxes generated. To then add this into the visualisation I simply looped over the selected location's list of location-boxes and painted chunks representing the range of sentences. The chunks were painted as an ellipse spanning from the first sentence number to the last. I invoked this method before the base paint method to then paint the sentences over the top to allow for the visual metaphor of locations in the visualisation being the backdrop for sentences.

Feedback

At this point I generated more copies of the initial visualisations for a couple of the books I was investigating and conducted a second round of interviews with my collaborating users. During the interviews I surveyed the users' opinions on the three initial visualisation forms. The codified results for character-highlighting are shown in table 3.4, the rest can be found in tables A.5, A.6, A.7 and A.8.

Their reaction to the initial images was positive, but they offered valuable suggestions for improvement. In particular, the lack of visibility of the character-highlighting approach was raised and potential ways of improving this suggested – such as connecting occurrences. The users agreed that the generated visualisations were cumbersome to interpret on a laptop screen due to their length, and proposed the idea of presenting in columns. The users were not forthcoming when discussing their desires in terms of user interface,

⁸Initially I added the location names manually.



Figure 3.13: An excerpt of a visualisation of *Alice's Adventures in Wonderland* with the garden (purple) and pool of tears (yellow) location-boxes highlighted.

Positive	Negative/Alterations
Use of dot increases visibility.	Could join up occurrences of characters.
<i>Good way of ascertaining key moments in the books when everyone comes together</i>	Annotation in margin could work well

Table 3.4: Table of codified responses from interview with users' views on character-highlighting approach.

though they did unanimously agree that the ability to query the image to ascertain the sentence associated with each line would be useful. In general, they were much more interested in developing informative and aesthetically pleasing visualisations. For this reason the majority of development focussed on this and left a fully-fledged user interface to future work.

During the secondary interview I discussed at more length the addition of investigating the age of the words used and highlighting specific parts of speech. The interviewees expressed the belief that these approaches offered lots of potential for exploratory analysis. In particular, being able to combine these approaches with the others to see how different language was used around different characters or locations. I consolidated the users' feedback and incorporated it during the second round of development.

3.6.2 Iteration 2

Development – Adaptations

Acting on feedback from the users after the first iteration, I implemented an approach to generate visualisations in multiple columns or rows. I adapted the *database to image co-ordinates* method to allocate each point to the correct column or row. By editing the central method for establishing image co-ordinates there was no need to adapt the implemented visualisation methods.

To increase the visibility of character-highlighting I painted circles larger than the original rectangles used to represent a character. I also introduced the ability to connect mentions of character instances, drawing a line between each dot (figure **3.15**). This raised the problem of visually distracting lines cutting across columns when there was space between mentions. To counter this I used linear interpolation to draw the line to the correct place at the bottom of the column and continue from the top of the next. In addition, I developed a basic interface to allow users to dynamically alter the size of the dots representing character mentions using the '+' and '-' keyboard keys.

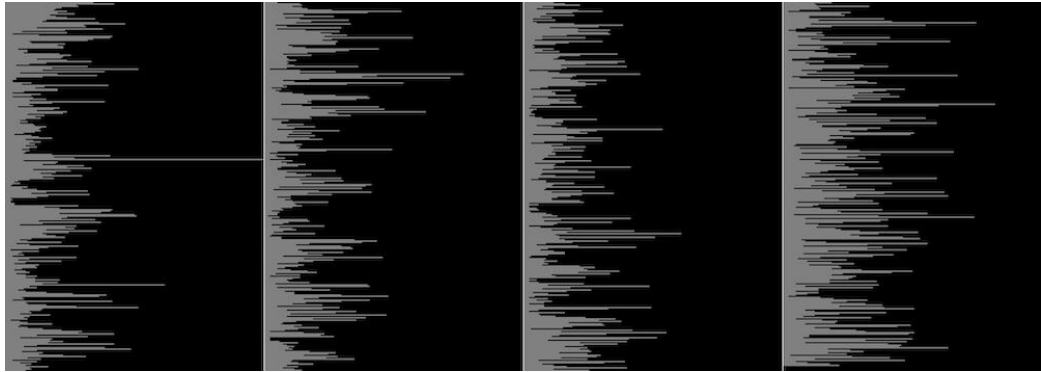


Figure 3.14: A visualisation of the whole text of *Strange Case of Dr. Jekyll and Mr. Hyde*.

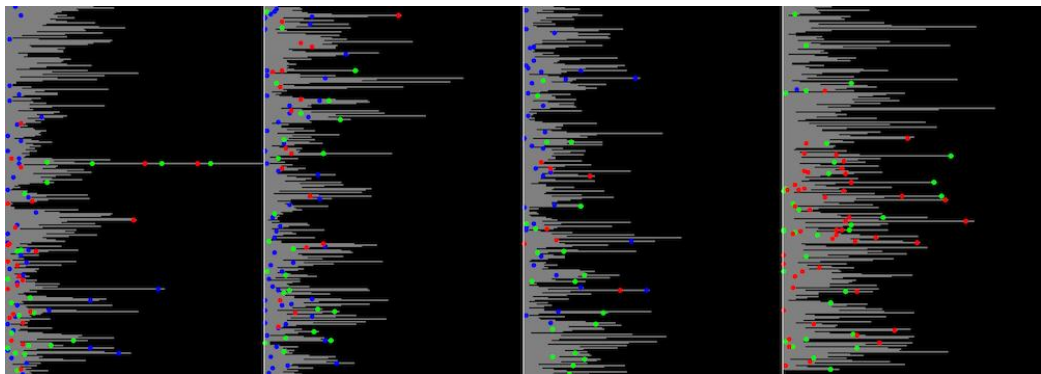


Figure 3.15: The whole text of *Strange Case of Dr. Jekyll and Mr. Hyde* with the characters Utterson (blue), Jekyll (red) and Hyde (green) highlighted.

I added the ability to highlight specific POS by creating a new method based on the character-highlighting function. This new function populated a list of points for a specified POS and then worked in the same way, placing coloured dots on their occurrences.

A final alteration was to introduce the idea of a margin. A criticism from one of the users was that the visualisations had begun to feel very ‘busy’ and that the addition of further information to the same visual ‘canvas’ could be overwhelming. I combatted this by introducing my margin approach, creating more ‘canvas’. The margins were small gaps between columns or rows, allowing the user to have the dots for character-highlighting placed in the margin instead of on the sentence – leaving this space free for other information, such as POS-highlighting.

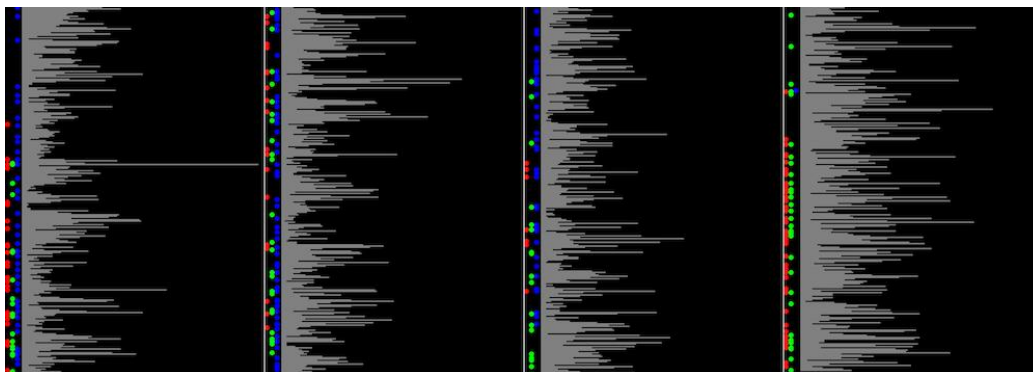


Figure 3.16: The whole text of *Strange Case of Dr. Jekyll and Mr. Hyde* with margins being used and the characters Utterson (blue), Jekyll (red) and Hyde (green) highlighted.

Additions – Data Ingestion/Handling

To aid the process of highlighting characters I added another database table per book to store the character names from each book. Furthermore, I introduced a new ingestion method allowing a user to populate the location and character tables by taking the proper nouns from the book and classifying each word as *character*, *location* or *ignore*. By having one user familiar with

the books do this and storing the data, the information could be re-used by a future user without needing re-classification.

During the secondary interview one of the users discussed the decision of where to delimit with chapters as a primarily authorial decision, and thus, I introduced this to the visualisations. I updated the underlying data by manually adding a delimiting character sequence at the start of every chapter. This sounds like a moderately intensive process, in reality it took under ten minutes for all the books. I updated the table for each book to include chapter-id as a field and populated this by incrementing a counter every time one of the delimiting character sequences was detected.

Additions – Age of Words

In order to find a way of visualising the age of the words used in each book, I first had to devise a method for calculating what year should be assigned to each word. I utilised Google’s Ngram project⁹ for this. The project works by having a large corpus of published work from each year and for a word that you enter, calculating the percentage usage of that word in the year. It is normalised for corpus size and a graph is generated to show the change in usage of the word. There are three main parameters when issuing a query: year range, corpus and smoothing. For corpus I elected to use their ‘English Fiction’ corpus as this most closely represented the literature I was utilising. I chose to use no smoothing to analyse the raw data. finally, limited the scope to between 1800 and 2000. The rationale for this was due to the relative paucity of data outside of this region in the Ngram corpora. This paucity of data – driven by, for example, limited publication before the 19th century – leads to spurious results for wider scopes than this (figure **3.17**). These parameter choices are easily modifiable at later dates.

There is no API for the Ngram project at time of writing but the webpage has a simple structure which returns the usage data as a JSON array. I wrote a python script to scrape this data and called this from my Java program.

⁹<https://books.google.com/ngrams>

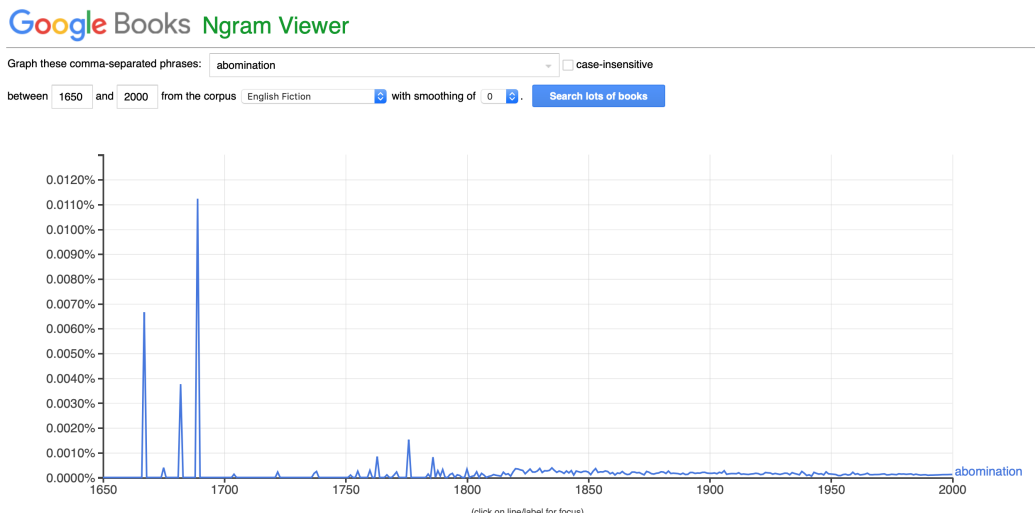


Figure 3.17: Example using the word ‘abomination’ of the effect of the paucity of data before 1800 on the results from Google’s Ngram viewer.

From the array I took the peak year of usage as the ‘age’ of that word. There is potential for further investigation on this front, such as words with bi-modal peaks, but I have left this as an extension for the future. I did add a threshold to ensure the peak year was significantly above the average usage over the range. Additionally, I tested to see if the word’s usage was overly common (greater than 0.1%). If the data failed either of these tests, the year was discarded and returned as null. I augmented the schema of the book tables to include year and called the script for every word. To minimise duplicated calls I created supplementary tables in the database to store words with peak years and those without. By checking these two tables first, I significantly reduced the number of times the python script had to run for each subsequent book.

Additions – Visual Features

I translated the age of words into visual features by assigning a spectrum of colour from 1800 - 2000 (3.18) using interpolation between the taking. I developed two ways of visualising this feature. Firstly, an adaptation of character- and POS-highlighting, placing appropriately coloured dots on the

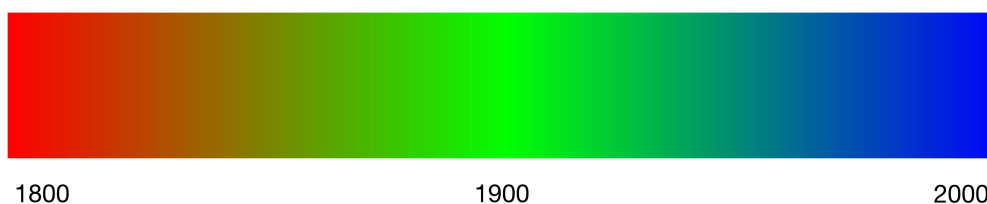


Figure 3.18: The colour spectrum assigning each year a colour.

locations of the words which had a year associated with them. Secondly, taking the average age of words in a sentence (ignoring nulls) and calculating the colour from this average age. Then, painting the sentence that colour. I am aware of the limitations a spectrum may face in terms of limiting accessibility of this project to colour-blind users. Addressing this issue is beyond the scope of this project and is left to future work.

I adapted the location painting method to paint the chapters of the book using an alternating grey and black pattern to show the end of each chapter. To improve visibility of this when using the margin approach the chapter-highlighting extends all the way into the margin whilst location does not.

Additions – Interface

Whilst my discussion with the users had led to a decision to focus on the visualisation aspect, they had all expressed a desire to interpret the image more directly. To deal with this I developed the tool to allow selection of a sentence using the mouse. The selection is highlighted and the relevant sentence printed adjacent to the visualisation frame (figure 3.22). Once a sentence is selected the arrow keys can be used to move up and down sentences as well as across column or row divides.

Finally, I implemented a zooming method to allow the user to visualise an individual chapter. I facilitated this by detecting a right mouse click and changing the *activated* list so that only records associated with that chapter were considered. I did this using methods on the database class to gather

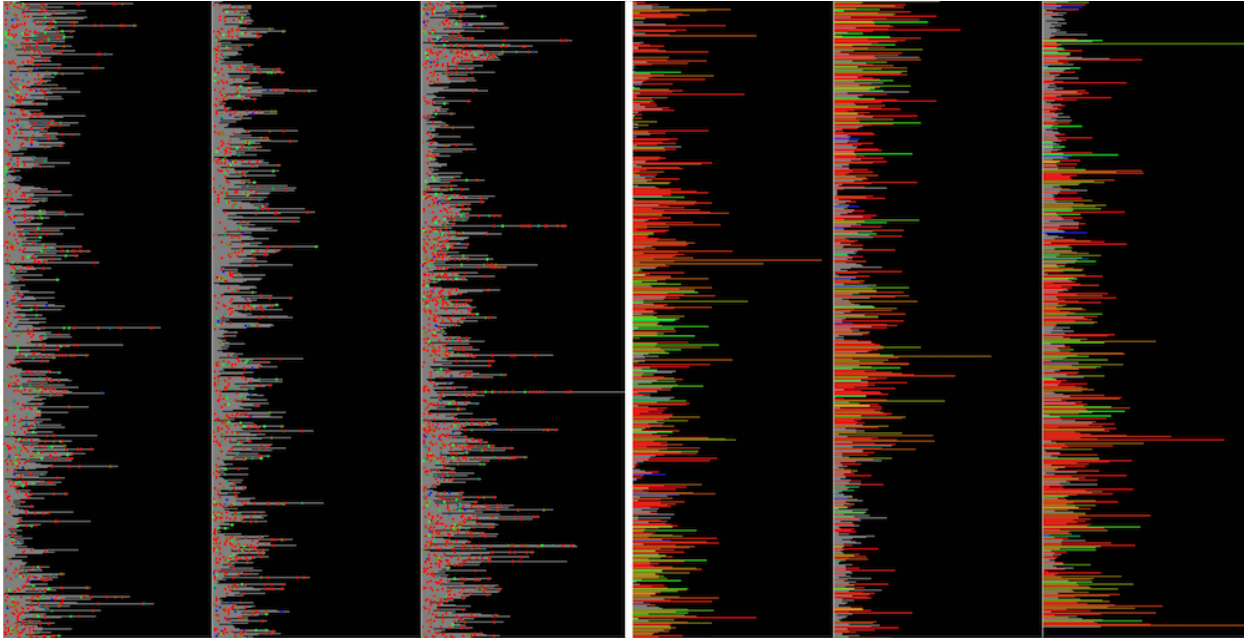


Figure 3.19: The left half of this visualisation of *The War of the Worlds* is generated using the granular dot approach and the right half uses the aggregated, per-sentence approach.



Figure 3.20: *The Story of Doctor Dolittle* with the chapters highlighted.

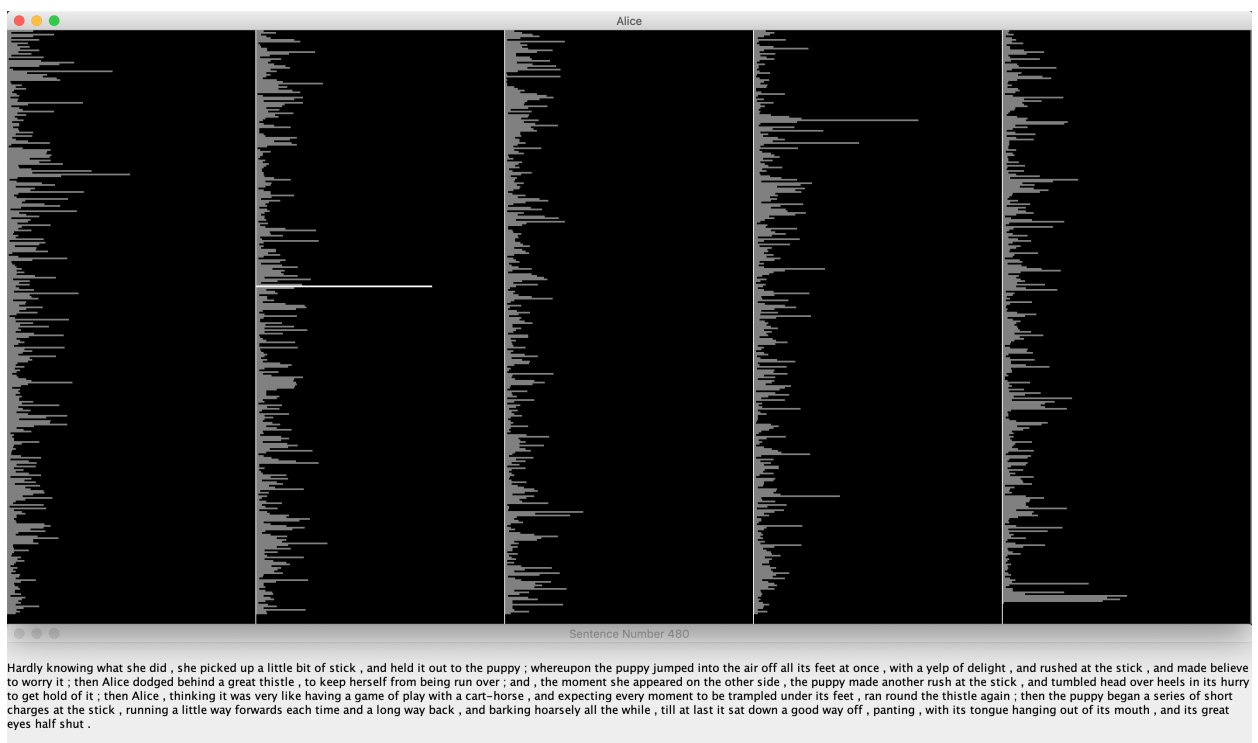


Figure 3.21: *Alice's Adventures in Wonderland* with sentence 480, selected through the interface by clicking on it, highlighted in white.

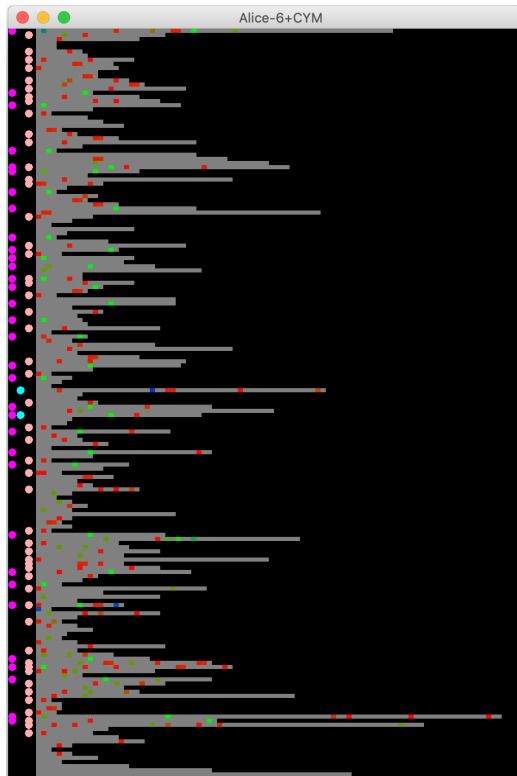


Figure 3.22: Chapter 6 of *Alice’s Adventures in Wonderland* with Alice (pink), the Queen (cyan) and the Hatter (magenta) highlighted in the margin. The granular version of the year information is highlighted in situ.

points for each chapter. This meant I could continue to use the same methods for highlighting the other visual features.

Feedback

The users codified responses for year highlighting are indicated in table **3.5**. The responses to the other features can be found in tables **A.9**, **A.10** and **A.11**.

The users all preferred the horizontal row approach over columns. Their reasoning was that it felt more analogous to reading actual text. User B described it as “like reading, but instead of reading a word at a time, you’re reading a sentence at a time”. The users agreed that it felt more natural

Positive	Negative
Average of years easy to interpret. (E)	Dots view can be overwhelming. (A)
Good to have changing granularity as you zoom in. (E)	
Finer granularity allows you to drill down. (B)	

Table 3.5: Table of codified responses from interview with users’ views on year-highlighting approach.

when zoomed in to revert to a single vertical column though, with similar reasoning. There was unanimous support for the use of the margin approach and they suggested that still having lines connecting the relevant points in the sentence together could be a good addition.

The users agreed that the averaged view of the years was better for an overview, but the granularity available was also useful when analysing further, especially when zoomed-in. They also expressed a desire to be able to tell which words had years assigned to them when highlighting a sentence.

Another significant feedback point was the suggestion of aliasing. For example, in *Pride and Prejudice*, the main character – Elizabeth – is referred to by other monikers, such as Lizzy and Eliza. The users thought a good application of character-highlighting was to see how each alias was used throughout the book, but that it would be useful to incorporate every mention of a character, including aliases, when comparing a character with others.

Finally, the users expressed concern again at the ‘busyness’ of the images being created. More specifically, the difficulty in distinguishing different features from one another due to the interference of the different features’ colours.

3.6.3 Iteration 3 - Final Modifications

Following the interviews I made final, mostly aesthetic, modifications to the tool in response to feedback, also adding significant functionality to the in-

terface. I introduced multiple global boolean options to represent different options available to the user. Most notable: the option to view the year painting as aggregated or individual because this had been highlighted as a benefit. I added another database table to store aliases for each character and introduced a toggle for this feature. When it was ‘on’ the aliases (such as ‘Lizzy’ for ‘Elizabeth’) were highlighted in the same colour as the main name.

I set default settings for the visualisation that reflected the preferences of my users so, when an image was generated, it would be in horizontal format with a margin containing highlighted characters and lines joining the correct place in the sentence (figure **3.26**). This switches to a single column when zooming-in on a chapter and incorporates continuity features to allow you to return to previous settings when zooming-out. I added the ability to change these feature parameters using keybindings. The changeable parameters are listed in table **3.6** with default value and description.

Furthermore, I added the ability to save the current state of the JPanel to a file to save the image you had generated from editing the feature parameters, as well as a new class to print high-resolution versions of the images with user-specified settings.

In order to address the ‘busyness’ of the visualisations I conducted further research into aspects of colour theory and concluded that the ‘busyness’ and lack of discernibility was likely due to the spectrum I had assigned to the age of the words spanning two thirds of the colour wheel[24]. This meant that even when altering hue or tone of colour there was a lot of interference and limited scope of the spectrum to choose from. To combat this I reallocated the spectrum (figure **3.23**). Then, utilising research relating to colour schemes [25], I assigned coherent schemes for characters and locations. I chose an adjacent colours scheme for characters and a triadic scheme for location [26]. I tried to encourage association between the different characters or locations being highlighted and chose a ‘neon-style’ theme for characters (figure **3.24**) and a ‘pastel’ theme for locations (figure **3.24**).

Parameter Name	Description	Default Value
Horizontal	Whether the visualisation is displayed as horizontal rows or vertical columns	True - will be displayed in rows
Character Highlighting	Whether the selected characters should be highlighted	False
Character Lines	Whether the painted characters should have lines joining their occurrences	True
Character Dot	Whether a highlighted character should be represented with a dot to make it more visible, or displayed in-line (used in the zoomed version)	True
Location Highlighting	Whether the related 'location boxes' should be painted for the selected locations	False
POS Highlighting	Whether the selected POS should be highlighted	False
Year Highlighting	Whether the 'age of words' should be painted	False
Aggregated Years	Whether the average per sentence for the year highlighting should be used	True
Margin	Whether the margin should be displayed	True
Number of Columns/Rows	The number of columns or rows that should be used to display the image	Determined per book
Verbose	Whether the full detail for each sentence should be printed (including POS and year, if applicable)	False

Table 3.6: Description and default values for the core changeable parameters of the visualisations.

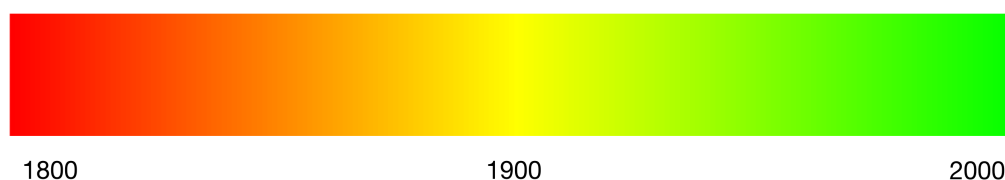


Figure 3.23: The updated colour spectrum assigning each year a colour.

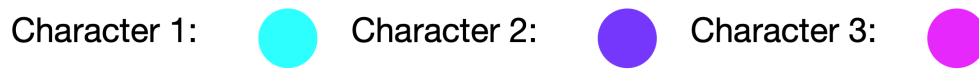


Figure 3.24: Part of the adjacent colour scheme for character-highlighting. There are more than 3 characters possible, but this is the default.



Figure 3.25: Part of the triadic colour scheme for location-highlighting. There are more than 2 locations possible, but this is the default.

3.7 “Finished Tool”

The software I created was more focussed towards generating visually appealing and interpretable visualisations from the text rather than developing a UI. This reflected the desire of the collaborating users to create rich, standalone images which could be used in physical form for narrative analysis (see section ??). I developed a tool which had the capability of ingesting large text files and converting this data into a database format which could be queried from a dynamic front-end. The front-end could be used to create real-time visualisations from the data and, using a basic UI, allowed a user to query the data through input via the mouse and keyboard. It also allowed the user to alter and specify specific visual features of the image to perform explorative analysis of the text visually. Finally, the tool also had the capability of saving a static version of the visualisation which could be used for narrative analysis. These visualisations could be optimised for display on a screen or for hard-copy, large-scale printing.

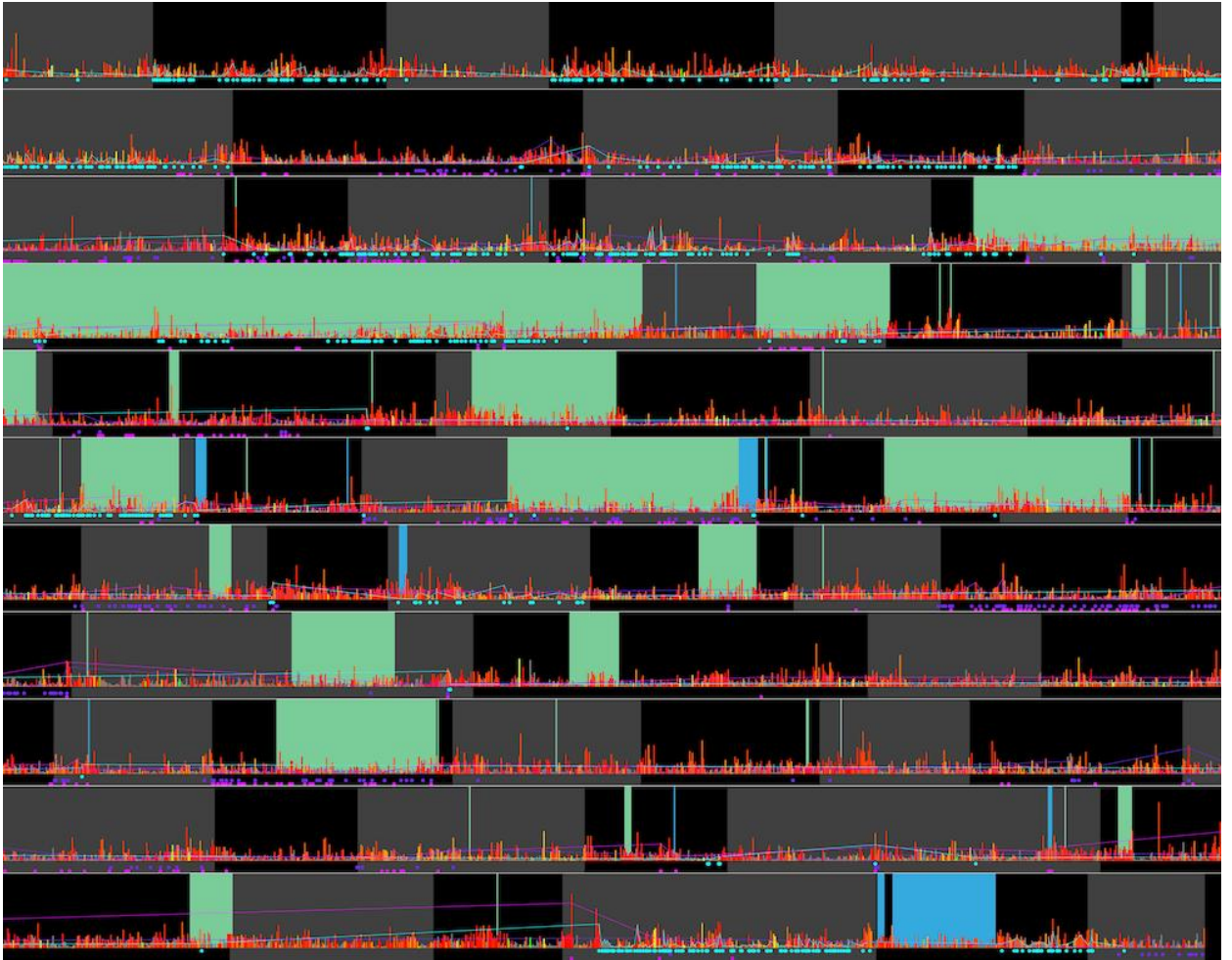


Figure 3.26: Example visualisation of *Great Expectations* with default parameters set and the following information - aggregated year-highlighting, chapter-highlighting, character-highlighting with lines (Joe (cyan), Estella (purple), Havisham (magenta)), location-highlighting (London (green), Blue Boar Inn (blue)).

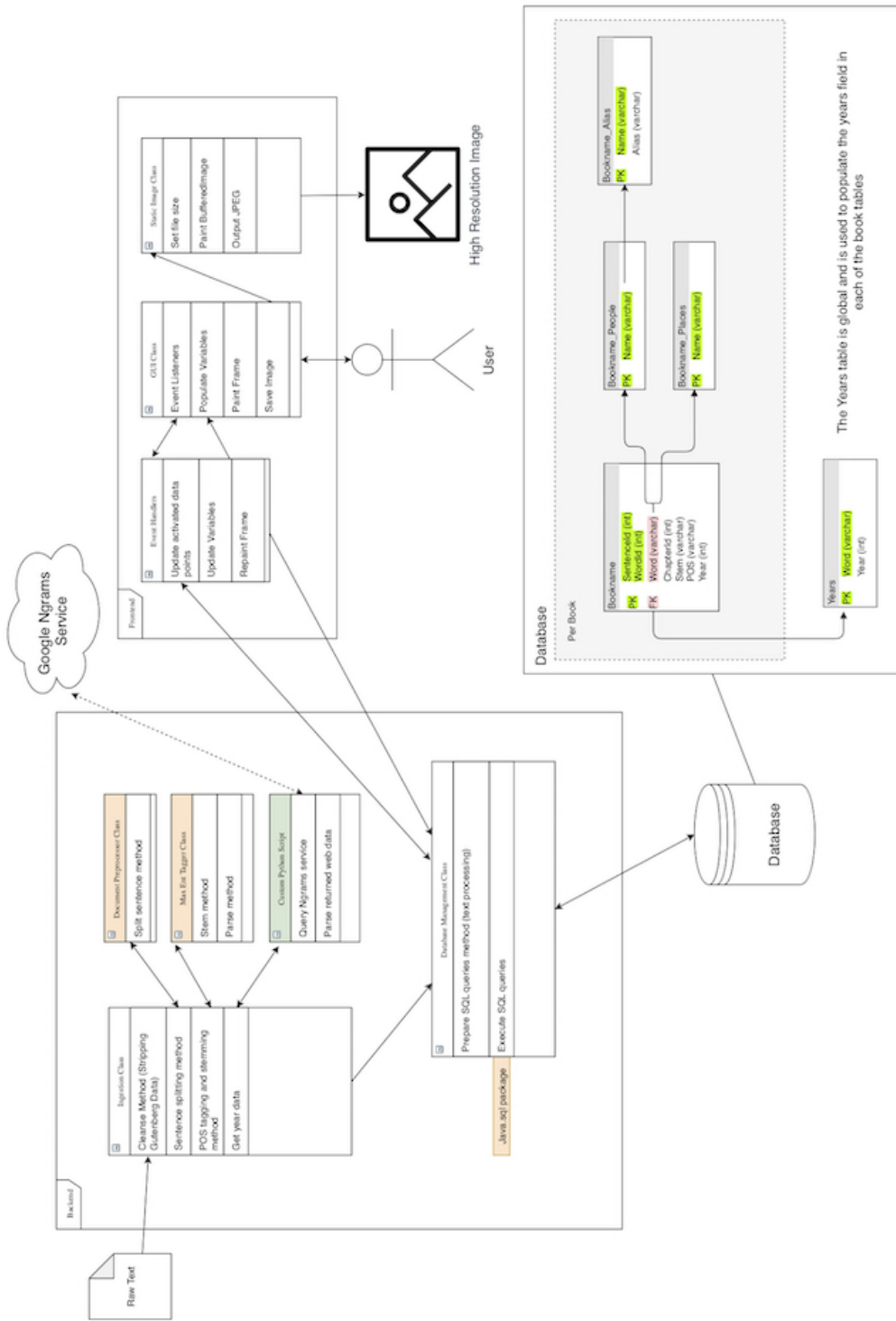


Figure 3.27: The final structure of the system – arrows represent data flow. Green headings are python scripts, orange headings indicate external libraries. There are elements missing from this diagram, but it gives an overview.

Chapter 4

Evaluation

To evaluate the visualisation tool, I must first clarify the goals for the visualisations and tool. Most importantly, I wanted to develop a text visualisation tool which would allow students studying literature to perform exploratory analysis visually. This should be enabled through novel visualisation techniques providing new methods of visual analysis. Moreover, the tool should allow users to generate visually appealing and informative images. Finally, these visualisations should make analysis of the literature less laborious and more engaging than traditional analysis methods.

4.1 Evaluation Plan

Taking inspiration from previous research evaluating information visualisation approaches [27] I developed three forms of evaluation. One focus of the project was to create visualisations which could convey information. The usability of the visualisations was also important. Prior research suggests that significant previous knowledge requirements diminishes the usability of a visualisation [28] [29]. With this in mind, I set out a plan to demonstrate that the visualisations have a tight mapping to the data they represent and so do not require significant prior knowledge.

The evaluation took the form of a controlled experiment (4.2), evaluative interviews (4.3) and two think-aloud studies [4] (4.2.3 and 4.4).

4.2 Experiment

4.2.1 Design and Method

I designed an experiment to determine whether participants could determine which book was represented by a visualisation. The rationale for this was that discernibility¹ was a strong indicator for informativeness.

The experiment was run as an online form which gave participants a minimal introduction to the visualisation techniques and for each of them asked them to identify which book they thought an image was generated from and vice versa. Figure 4.1 shows an example of this for the character-highlighting feature.

I recorded participants' socio-demographic information (age, gender, ethnicity and profession) [30], but the sample size was too small ($n = 31$) for this to yield significant results. The experiment is still active, and this information will be valuable in future evaluation of the results.

The combinations of questions in the form were randomly assigned to limit the influence of individual books on the results. To assert that positive identifications were a result of the visualisation features I had designed the users were tested on images representing 350-sentence excerpts, as well as the whole-text. This strengthened positive indications that the new approaches were conveying information.

The form took less than 10 minutes to complete, indicating that the visualisation approaches did not require significant prior knowledge. A point reinforced as the majority of participants had never seen the images or style of visualisation before.

¹Discernibility here means ability to identify the book from the image.

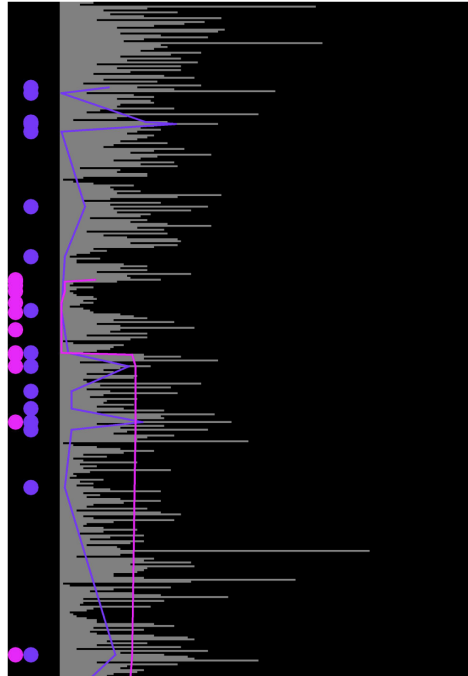
Finally, the form incorporated optional long-text comment sections to invite open-ended comments on each feature. This allowed participants to express views I had not encountered from my collaborating users.

The eventual format of the experiment was as follows:² for every question the user was first presented with a single image and an option of two books to select from; and secondly was presented with two images and asked to select which one they thought had been generated from a selected book:

- Base paint(BP) – 350 sentence excerpt(350)
- Base paint - whole book(WB)
- Base paint with chapter delimitation(ChD) – whole book
- Base paint with character-highlighting(CH) – 350 sentence excerpt
- Base paint with character-highlighting – whole book
- Base paint with location-highlighting(LH) – whole book
- Base paint with year-highlighting(YH) – 350 sentence excerpt
- Base paint with year-highlighting - whole book
- Base paint with chapter delimitation, character-highlighting, location-highlighting and year-highlighting – whole book

²A copy of the form that is still live at https://docs.google.com/forms/d/e/1FAIpQLSe5aE_NdUJ-1_EzhtrnhCaROHoUbSJDmELzfERbvLT5nbI8QA/viewform?usp=sf_link

Figure 10



Select the book you think Figure 10 is generated from *

- 'Strange Case of Dr. Jekyll And Mr. Hyde' (1886) by Robert Louis Stevenson
- 'The War of the Worlds' (1898) by H.G.Wells

Figure 4.1: Example question from section 3 of the questionnaire asking the user to identify the book from the image (the first 350 sentences of the book with the main characters highlighted) - the correct answer is *The War of the Worlds*.

Limitations

Whilst 30 respondents is a valuable sample size [31], increasing this would allow me to draw stronger conclusions from my results and more thoroughly control for external influences.

With further resources the questions could be dynamically assigned based on the users' familiarity with the texts. However, this functionality was not available through most mainstream survey providers, instead I had to manually change the assigned questions at set intervals to ensure their ran-

domisation. This meant that sometimes the correct selection was a book that the participant had never read. To alleviate this impact, I filtered my results to show results for different levels of familiarity and removed results (correct or incorrect) below the chosen familiarity.

4.2.2 Results

The hypothesis I was testing with this experiment, was that the visualisations were discernible – the users could tell the difference between them, and identify the correct book related to an image and vice versa. A result of over 50% would indicate success as the test followed a binomial distribution – there was a correct and an incorrect answer for every question. In order to establish the significance of these results I conducted one-tail binomial significance testing. The decision to use a one-tailed test was made as there was no reason for unexpected negative results. Moreover, initial manual viewing of the results did not suggest any [32] [33]. The following equation was used to determine the significance of each result:

$$p = \frac{\sum_{i=Y}^N \binom{N}{i}}{2^N}$$

Where N is the sample size and Y is the number of correct answers. Where the value of p was less than 0.05 this indicated a statistically significant result at 5% significance level. These are listed in results tables **4.1** and **4.2** where the result is significant.

The purpose of carrying out this experiment was to show the informativeness of the visualisations and that distinguishing features of a book can be conveyed using novel visualisation techniques. In every case where I tested visual features for 350-sentence excerpts and the whole book, results improved in the latter case. This suggests that users were using the length of the book as an informative feature as the relevant lengths of the books was inferable from the images.

Section	Book from Image(BfI) (%)	Image from Book(IfB) (%)
BP - 350	47	39
BP - WB	48	65
BP + ChD	39	63
BP + CH - 350	61	<u>70</u> - $p = 0.047$
BP + CH - WB	<u>84</u> - $p < 0.01$	<u>81</u> - $p < 0.01$
BP + LH - WB	50	<u>77</u> - $p < 0.01$
BP + YH - 350	68	58
BP + YH - WB	<u>78</u> - $p < 0.01$	<u>84</u> - $p < 0.01$
BP + ChD + CH + LH + YH	<u>83</u> - $p < 0.01$	65

Table 4.1: Table of percentage of correct answers for the different sections, both for selecting the book name from the image and for selecting the generated image from a given book. All results are given to two significant figures. Underlined, italicised results indicate significance and are listed with their p-values.

The lower results in the location-highlighting section suggest that this has not been as effective as an informative visual feature. A potential reason for this is due to a lack of understanding of the visualisation feature. This is discussed in the final evaluative interviews (4.3) but the discrepancy between identifying the image from book is evidence to support this, as the increased reference could have helped in their selection. This result is also explained by the varying levels of familiarity with the books. Some of the respondents had only seen film adaptations and as such, the poor results for chapter delimitation are not surprising.

The results for character- and year- highlighting suggest that these are both very informative features allowing even participants with limited knowledge of the books to identify them from features of the generated image. Reasons for these results not being even higher are touched on in 4.2.3.

It is important to observe how the results change with regard to the respondents' familiarity with the literature. Every participant recorded their familiarity with each book using the following scale:

1. Never read it
2. Haven't read it but have a vague idea of what it's about
3. Have seen a film adaptation
4. Have read it once
5. Familiar (studied in school or have read multiple times)
6. Very familiar (have studied recently)

Part of the reason for the initial results not being as high could be due to the vague nature of their knowledge of the books. To illustrate this see table **4.2** which shows the results for each section again, but with a minimum level of familiarity (from 3 to 5 – the first table shows minimum familiarity 2 and no respondents indicated their familiarity as 6).

These results show that as familiarity with the text increases the accuracy of identifying the book does too. This lays credence to the idea of the visualisations being informative and discernible, though the (increasingly small as familiarity increases) sample size makes it hard to highlight significant results and draw firm conclusions from this data. I would like to address this with more resources in the future.

From this data alone it is hard to surmise the full reasons behind the results. The participants lack of familiarity with the images could be a reason for some of the discrepancies. This is reinforced by later questions in the survey receiving better results, though this could also be due to their more informative nature. Those that utilised the optional comment sections identified the initial lack of comparable material as a major challenge. They also pointed to using the length as an identifying feature, reinforcing the effectiveness of this aspect of the images. Finally, users highlighted their lack of familiarity with the texts as a reason for difficulties. In order to understand and confirm some of these shortcomings, I conducted a think-aloud walkthrough of the form.

Minimum Familiarity:	3		4		5	
Section	BfI(%)	IfB(%)	BfI(%)	IfB(%)	BfI(%)	IfB(%)
BP - 350	50	33	60	33	67	33
BP - WB	60	71	70	<u>89</u> $p = 0.030$	71	80
BP + ChD	41	<u>75</u> $p = 0.021$	57	60	60	100
BP + CH - 350	60	62	67	29	67	25
BP + CH - WB	<u>80</u> $p = 0.018$	<u>95</u> $p < 0.01$	<u>100</u> $p < 0.01$	<u>92</u> $p < 0.01$	100	<u>100</u> $p < 0.01$
BP + LH - WB	53	<u>90</u> $p < 0.01$	38	<u>83</u> $p = 0.026$	50	<u>89</u> $p = 0.020$
BP + YH - 350	70	58	N/A	50	N/A	100
BP + YH - WB	<u>80</u> $p = 0.018$	<u>90</u> $p < 0.01$	76	<u>78</u> $p = 0.021$	<u>100</u> $p = 0.016$	<u>100</u> $p = 0.016$
BP + ChD + CH + LH + YH	<u>82</u> $p < 0.01$	<u>73</u> $p = 0.026$	75	67	86	N/A

Table 4.2: Table of percentage of correct answers for the different sections relative to the minimum familiarity with the book, both for selecting the book name from the image and for selecting the generated image from a given book. All results are given to two significant figures. Underlined, italicised results indicate significance and are listed with their p-values.

4.2.3 Think-Aloud Walkthrough

The walkthrough built on protocol laid out in previous work [4]. I opted to perform a concurrent think-aloud study (rather than retrospective) due to research suggesting more accurate representations of thought processes [34]. The walkthrough took the form of a volunteer³ completing the questionnaire whilst explaining his thought processes for his selections aloud. I recorded this walkthrough and cover key points from each section here.

Introductory Section

Whilst reading the introduction the volunteer expressed that sentential length could be a useful indicator for roughly dating the book. Their rationale was that “*modernist authors tend towards stream of consciousness in their writing which could lead to longer sentence length*”. This was an idea that had not previously come up and suggests further informativeness of even the most basic feature. They highlighted serialisation of 19th century literature and the effect this could have on chapter length, providing a potential explanation for the increase in discernibility with chapter delimiting information.

The participant had the following familiarity with the surveyed literature:

- Alice’s Adventures in Wonderland - Have read once
- The Story of Doctor Dolittle - Haven’t read it but have a vague idea of what it’s about
- Great Expectations - Familiar
- Pride and Prejudice - Familiar
- Strange Case of Dr. Jekyll and Mr. Hyde - Familiar
- The War of the Worlds - Have read once

³A final year English student at the University of Cambridge who had not been involved with the project prior to this point.

Base Paint

The user highlighted the lack of referential scale and the difficulty this posed in establishing sentence length. Even using knowledge regarding writing styles they only made tentative guesses. This could be addressed by showing how many words are in each sentence with delimitation and it supports the explanation that users found it hard to discern in early sections due to unfamiliarity with the techniques and images. The user's thought vocalisations showed their attempts to extract information from sentence length, such as the suggestion that "*long interjecting sentences seem to suggest character's vocalisation of thoughts*".

When presented with the full-text images the user utilised the perceived length of the book to make a more informed decision. One thing that became clear was that the attempt to keep the scale the same in order to allow users to infer information about length had not been made clear enough as the volunteer sought clarification on this point. These results explain the performance increase for the whole book compared to excerpts and partially explains why they did not improve more – as user's were not fully comfortable with how the length of the image was analogous to book length (this was an opinion borne out in comments on the survey too).

Chapter Delimiting

In this section the user's initial instinct was to utilise the length of individual chapters rather than the number of chapters. This is not how I envisaged this feature being used but shows an additional transfer of information as the volunteer correctly identified *Pride and Prejudice* through knowledge of the first chapter's length. They also drew associations between the length of chapters and perceived drama in the story – implying that shorter chapters suggested more action.

The user's comments reinforced ideas that the visualisations become more informative the better your knowledge of the book.

Character Highlighting

The participant successfully identified each book in the section, even from excerpts, due to knowledge of when characters appeared. The user voiced observations from the whole book images, such as the episodic nature of *Alice's Adventures in Wonderland*, showing that the visualisations were conveying analytical information about the style of the text, as well as identity information. This explains the poorer results for excerpts due to a requirement for significant knowledge of the text. Moreover, the user repeatedly mentioned how strongly characters sit in the minds of the reader, cementing the use of character-highlighting as a good identifying and informative feature.

They did also highlight the influence the length of the text had on their decision making process too.

Location Highlighting

The user's seeking of clarification for location-highlighting confirmed that the description of how these images were generated was not clear enough. This detracts from the informativeness of this feature, particularly as I was trying to emphasise the lack of required prior knowledge. The user drew connections between differing locations and physical travel, suggesting flaws in the presentation of location information because discussion of a place will still trigger its highlighting.

Additionally, they described the lack of prominence of location and perceived greater prominence of characters. This feedback helps to explain the results in 4.2 of location-highlighting being less informative for identifying a book.

Year Highlighting

The initial lack of referential material was again reported as well as the potential difficulty in interpreting the spectrum due to colour selection. Though

I did address this (**3.6.3**), it warrants further work. Understanding the processes behind the visualisations was raised as the user was unsure how non-sense words from *Alice's Adventures in Wonderland* would be handled.

The whole book images with their added referential information helped the user to decide and confirms my observations of the experiment results. Surprisingly, despite the users' comments on the difficulty of discerning the spectrum, most users (themselves included) were fairly comfortable when it came to selecting the book or image.

Combination of Features

The volunteer's selections in this section raised the influence of previous answers. They discussed how recognising features from earlier sections meant they were more confident in making decisions. This helps to explain the trend of improving performance as the survey progresses. Again the familiarity with the books influenced decisions heavily. Finally, they also showed significant focus on using the character-highlighting information to ascertain the book and using the other visual features as confirmatory information – such as the rough age of the book from the year-highlighting.

4.2.4 Conclusion

The results and discussion from the think-aloud walkthrough were valuable in confirming conclusions already drawn from the experiment results.

Firstly, that familiarity with the book had a large impact on decision making, a point raised by the volunteer and evident in table **4.2**. It further confirmed the informative nature of some of the visualisation features – character- and year-highlighting – and explained the decreased performance for some of the others. This decreased performance was due to a lack of thorough explanation and understanding of the processes behind the visualisation processes. This raises further issues of requiring training and exposure to images before they

become useful, though the rapid improvement suggests this would not be necessary.

The walkthrough raised directions for future analysis into how other features of the book could be inferred from visual features. In general, it confirmed and explained the results presented in **4.2.2**.

4.3 Evaluative Interviews

Throughout the project I conducted well-documented formative evaluative interviews (see **3.2.2**, **3.6.1** and **3.6.2**) which shaped the generated visualisations. I conducted final, summative, evaluative interviews to ascertain the users' opinions towards the tool following guidelines set out in Corbin and Strauss' seminal paper on evaluating qualitative data [35]. In particular I wanted to draw on their argument of repeated themes recurring over the interview process and to see how users felt the tool dealt with issues they had raised.

I sought advice from a social-science researcher in how best to conduct and present the interviews. As per their advice, I have written primarily open and non-leading questions. By doing this, I do not predispose the users to responses, and conformity across interviews is a much stronger indication of agreement. Where the questions are leading, I have acknowledged this in my results. I have separated questions into distinct sections, though, the sequential nature of the interview can affect the answers to later questions. I recorded each of the interviews and later transcribed their responses. You can see the full interview questions and transcriptions in **A.6.1** and **A.6.2**.

4.3.1 Results

I have presented the salient points from each of the four sections, where appropriate linking back to themes from the formative interviews.

Section 1

In this section the user was asked open-ended questions in the general form “*What do you think about ...*”. This allowed the users to express their opinions rather than being led to an answer, thus strengthening cross-interview agreement.

Since the users have all been involved directly with me during the development process they are, at least subconsciously, aware of the aims of the project, and so may adapt their answers accordingly.

That being said, the users agreed that the basic image raised interesting patterns in the text (one user was surprised at the lack of randomness of sentence length), but without additional information did not offer very much. This confirms conclusions drawn from the experiment and walkthrough of the survey about this feature. User B also took the opportunity to echo an opinion raised by multiple users in **3.6.2** that the horizontal presentation mimicked the process of reading a book. The recurrence of this theme suggests a strong visual feature.

The users offered differing feedback on chapter delimiting. Where consensus above helps to draw conclusions and strengthen arguments, divergence in this case shows the versatility of the images created. One user (C) highlighted that combining chapter delimitation with sentential length facilitated identifying which chapters have more speech. Thus, it allowed them to infer more information about the progression of the story from their own knowledge of the book. User A discussed the usefulness of chapter delimitation in “*helping people to locate things in the book to go and analyse more*”. This draws from previous comments from the users about the value of the visualisations as preliminary analysis tools to encourage further investigation. Finally, user E felt the introduction of chapter delimitation made the image more aesthetically pleasing, making the image “*feel more like a book*”. These comments are particularly important as tying the visualisations to the context of the book was an important part of the project.

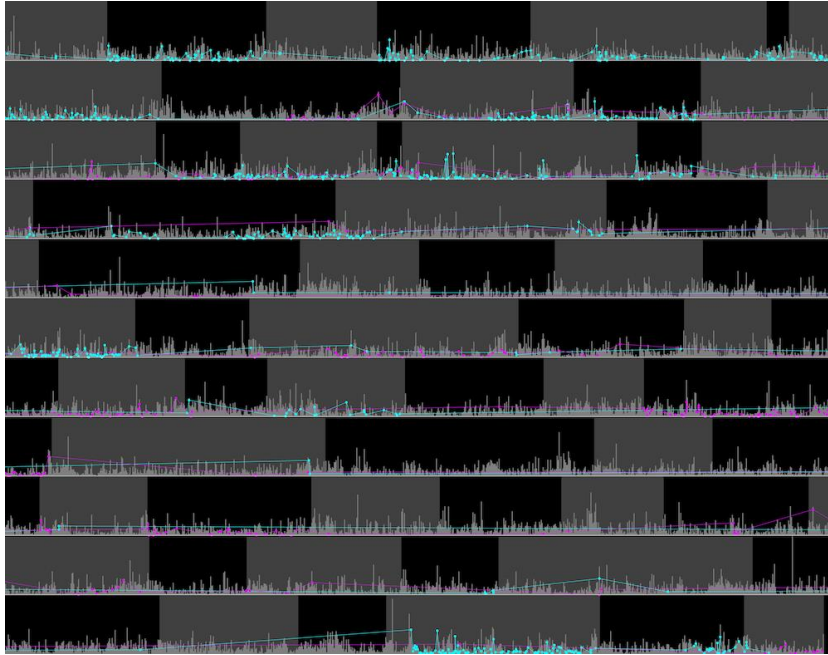


Figure 4.2: Example visualisation of *Great Expectations* showing the appearances of Joe (cyan) and Estella (magenta) throughout the text.

The response to the highlighting of characters' occurrences was unanimously positive. There were differing views as to what made this feature so valuable. User C discussed how the combination of character mentions with other features, such as sentence length, allowed the viewer to infer who was in dialogue with who. User E summarised the divergence of responses well here with: *"I could look at this for a while with different questions and come up with different responses"*. This kind of feedback is really valuable with reference to one of the driving motivations of this project, namely facilitating exploratory analysis. The user continued, giving specific examples from *Great Expectations* where the image presents interesting information about the story. Most notable, the contrast of the opening and ending of the story: the heavy presence of Joe (the protagonist's brother) in the opening, and his absence, and replacement with Estella (the protagonist's love interest) in the conclusion of the story (see figure 4.5).

User B commented on the episodic nature of Alice's interactions in *Al-*

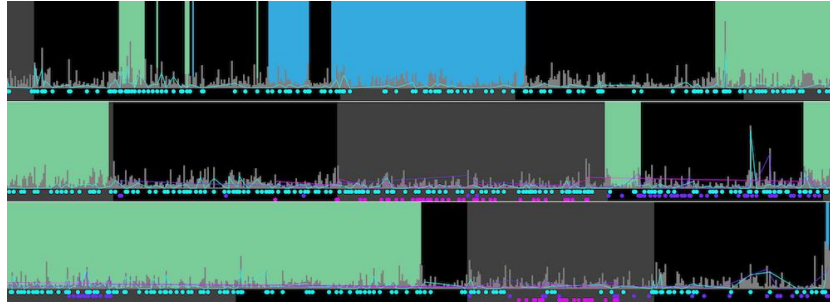


Figure 4.3: Example visualisation of *Alice's Adventures in Wonderland* showing the episodic interactions of Alice (cyan) with the Hatter (magenta) and the Queen (purple), as well as the recurrence of location through the *pool of tears* (blue) and the *garden* (green).

*ice's Adventures in Wonderland*⁴, paired with the repeated returning to the 'garden' location. Both of these features are clearly visualised in figure 4.3.

The most disappointing features from the users' points of view were the POS- and location-highlighting. The feedback towards location-highlighting was more neutral than negative: they did not feel this offered them significant additional information. In their opinions the location information just confirmed what they already knew. As my users were all literature students familiar with the text this could potentially mean that the location information could be a way of imparting this information to less familiar viewers. Users referenced back to their previous comments, and confirmed earlier conclusions (??) that, without understanding of the visualisation processes, it could be hard to interpret the information accurately.

The biggest criticism of the POS-highlighting was a frustration at inability to illustrate points about the text that the user believed to be true. This frustration can be interpreted as a failure of this feature or it could indicate that assumptions students often make are based more on subjective reading of the text rather than authorial decisions. The users thought this feature had a lot of potential, but, due to the proliferation of almost every part of speech, it was hard to extract valuable information from the images. More than one user suggested this might be a feature more useful to linguistics

⁴Just as the volunteer had done in 4.2.3.

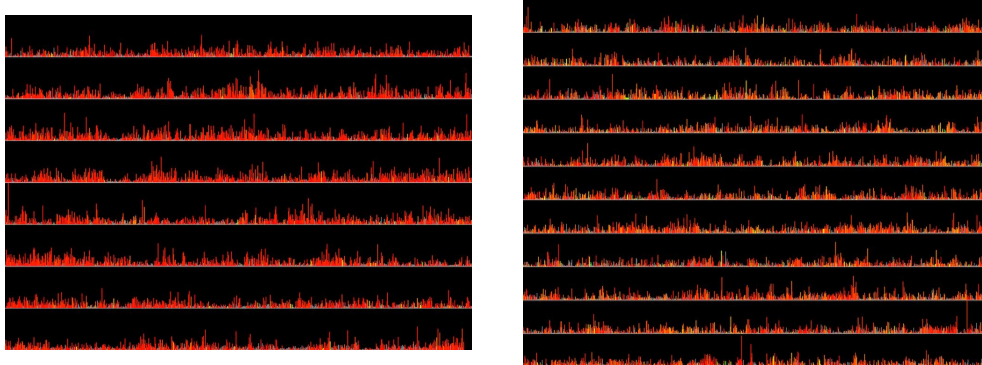


Figure 4.4: The image on the left is *Pride and Prejudice*, published in 1818, the image on the right is *Great Expectations*, published in 1861. Despite only slightly more than 50 years difference in publication date, there is a clear visual difference between the two images.

students, as their familiarity with NLP could allow them to obtain more information and select more informative parameters.

With regards to year-highlighting, one user raised the concern that had appeared in the survey walkthrough (4.2.3), that the selection of colours made it hard to discern. In general the users were surprised by the perceived accuracy of this feature in indicating the age of the book – this can be observed in figure 4.4. They expressed disappointment that it did not allow for particular analysis of the text – such as the use of modern language around urban locations – but were positive about the combination of aggregated and granular versions of this feature. Again, user A emphasised the value of this as a starting point for further in-depth analysis.

Finally, the users agreed that the combination of features into one image worked well. There was a tendency for the image to be overwhelming, but building it up through layers helped to appreciate each aspect of the image independently. They also liked the use of margins as, in the words of one user, they “reminded me of annotating books by hand”, as well as nicely separating the visual domains of the features.

Section 2

Section 2 discussed the visualisations with regards to previous approaches discussed in **3.2.2**. These questions are leading and the answers likely to be pre-disposed in favour of the visualisations I have created. However, I do not think this detracts from the validity of their comments, as they offer firm reasoning behind *why* they believe my approach to be an improvement. All of the users agreed that the context my visualisations provided was invaluable in making the images more informative. They also commented on the effectiveness of visualising multiple aspects simultaneously. An opinion expressed by user B was that “*you can see frequency [of occurrence] as a whole book but you have to infer it for yourself*”. They viewed this as a positive for the visualisations as it encouraged interpretation of the image and to consider reasons behind the visual features that appeared, leading to better analysis.

The users also agreed that the images alone could be used for analysis of a text, with the caveat that the viewer is already familiar with the text. User E discussed the power of the visualisations in confirming theories as to do so without visualisation is laborious and time-consuming. Each of the users independently referred to the images being most useful as supplementary analysis tools.

When discussing similarities between their usual analysis of text compared with the visualisations, the users highlighted two main similarities. Firstly, the analogy to reading due to the design decision of presenting the image in horizontal rows. In tandem with this, the use of margins as annotations resembled how they would approach analysis with a physical book (see figure ??). User A discussed the importance of spotting patterns when building arguments and how “*the visualisation is perfect for that*”.

They did feel that the visualisation analysis also offered significant differences compared to reading. They drew the primary positive of the image allowing you to “*take in*” the book in one visual field in contrast to normally requiring intense knowledge of a text to make statements on it as a whole. They also noted the usefulness of the visualisation as “*a map of where to go*”, or

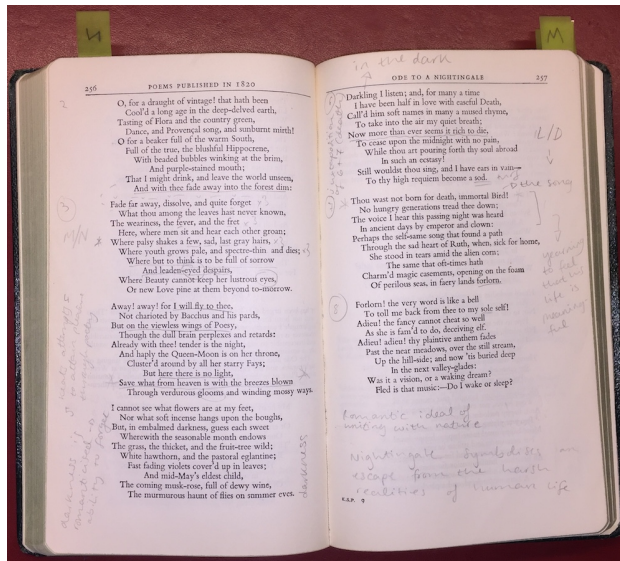


Figure 4.5: Example of physical annotations from user E's studies.

the idea, prevalent throughout formative interviews that the visualisations provide guidance to interesting sections of the story for further analysis. The users do not draw any strong negative differences in the use of visualisation to analyse. This is most likely due to the moderately leading nature of this section.

Section 3

In section three, the users were all confident that they could identify books from their images, especially when provided with the experiment reference sheet (figure 4.6). They all noted that the length of the book, shown by the visualisation, was an important feature in allowing them to identify the image. This, combined with year-highlighting, allowed them to pair the factual information (date of publication, length of book) with their interpretations of characters and locations to identify the book using all the visual features available. The prevalence of characters' names allowed them to make inferences about whether the book was written in the first or third person, and give them more confidence in identifying the visualised text.

Book details

Name	Author	Date Published	Highlighted Characters	Highlighted Locations	Authorial Voice
Alice's Adventures in Wonderland	Lewis Carroll	1865	Alice, Queen, Hatter	Garden, Pool of Tears	3rd person
Doctor Dolittle	Hugh Lofting	1920	Doctor, Polynesia, Jip	Africa, Puddleby	3rd person
Great Expectations	Charles Dickens	1861	Joe, Estella, Havisham	London, Boar Inn	1st person
Jekyll and Hyde	Robert Louis Stevenson	1886	Utterson, Jekyll, Hyde	London, Soho	3rd person
Pride and Prejudice	Jane Austen	1818	Elizabeth, Darcy	Longbourn, Netherfield	3rd person
The War of the Worlds	H.G.Wells	1897	Curate, Ogilvy, Henderson	Woking, Mars	1st person

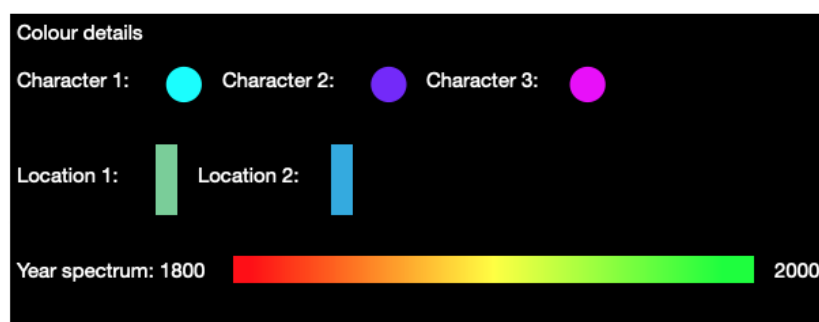


Figure 4.6: The reference sheet used in the online form.

An example of the informativeness of the year-highlighting approach is in *Pride and Prejudice* where Austen uses the word “relationship”. This word appears bright green as its usage is much more common in modern texts. It is extremely interesting to see this as a word not in keeping with the lexicon of the period. User E elaborated on this point by pointing to how “the idea of romantic relationships is very favourably compared to Austen’s work” and expressed their surprise at how “relationship” was not a common word when she wrote and reminds you “*that her work was groundbreaking*”.

Section 4

Section 4 addressed the potential usage of the tool. This was leading, but, when combined with previous discussion and engagement through project,

the feedback is still valuable. All of the users thought that the tool would have been useful during their studies, users A and E highlighting the idea of visualisation of the whole context of a book allowing them to identify patterns for further exploration. They also described the approach to visualising authorial decisions directly as more objective than relying on subjective analyses they may generate themselves.

The users proposed two main views on the primary use of the tool. Firstly, the use of the visualisations to “*objectify things that are hunches*”. This reflects discussion from the previous paragraph of being able to show directly the author’s repeated use of character names (as an example) over the whole text rather than having to identify multiple fragments to support your claim. Secondly, user B discussed the use of the images as an interrogative tool. This is important to highlight as in **3.2.2** they discussed their primary interest as creating visual art. This was still a big factor, but they now also recognised the benefit of the tool as an analytic device. This is testament to the informativeness and novelty of the images created.

4.3.2 Conclusion

The results of the final summative interviews are powerful in drawing conclusions about the success of the project. Taken in combination with the formative interviews the recurrence and adaptation of the users’ opinions over time becomes clear. Most notable is the repeated allusion to the benefit of this tool as a supplementary aid for analysis, in particular drawing the viewer in to interesting aspects of the image relating to interesting aspects of the text. In addition to this, the more “*distant perspective*” of this novel form of analysis offered by visualising the whole book, allowed the users to make statements confidently about patterns throughout the text.

4.4 Think-Aloud Study

As the final part of evaluation, I conducted a study of a user (E) creating their own narrative visual analysis. I built on cognitive walkthrough techniques with the web [36] but the eventual structure was more closely aligned to a think-aloud study [4], having the user describe aloud what they wanted to visualise. Though I describe the user as ‘using’ the tool here, due to the focus on the visualisations rather than usability, I had the user describe aloud what they wanted to do and I performed the actions to generate the images. Whilst this creates some cognitive dissonance between the thought and act of creating the visualisation, the recorded comments allow for insight into how a genuine user of the tool wanted to create visualisations. The walkthrough enabled me to see how motivations drove the decisions the user made, and whether they could create an informative narrative analysis through the exploratory nature of the tool I had created. The user chose the book *Great Expectations* due to their familiarity with this text.

The walkthrough served as further confirmation of several points, in particular the power of having the whole context of the book available in one visual frame. The user found identifying patterns using character-highlighting very useful. This was emphasised by the fact that they began with this feature and added the other features to the image to see how characters interacted with them. The user chose to show the interaction of Herbert and Estella with the main character: Pip. By highlighting these three characters (figure **A.1**) they were able to draw attention to sections of volume two where Herbert and Estella occur frequently but in almost exclusive isolation from one another.

The user also chose to highlight the locations London, Satis (Miss Havisham’s house) and the Blue Boar (an inn in Pip’s home town) as well as using chapter-highlighting (figures **A.2** and **A.3**). They were frustrated at the lack of immediate clarity, though still found it interesting to see the prevalence of London in volume two and the coinciding of brief interludes of the Blue Boar with absence of interaction with Herbert or Estella.

The user was frustrated when adding POS-highlighting to the image at their inability to identify meaningful patterns (figure **A.4**). They attributed this to their usage being too frequent, coupled with their own lack of expertise of analysing language use. They thought this could be a useful addition for a linguistics student but was difficult to use without the requisite knowledge. Another important observation here was the worry of “*falsely privileging information*”. This arose when highlighting the use of comparative adjectives. By selectively viewing this data the user was worried about impressing their views of the text onto the image, rather than extracting information about the text from the image.

When experimenting with year-highlighting the user reinforced the usefulness of this as a comparative tool. They also raised here disappointment at not being able to see aspects suggested as potential uses of this feature – such as, the use of modern language around different environments. However, they suggested that this could lead to challenging the validity of these arguments using more objective data from the visualisations rather than subjective interpretations.

The user turned on aliasing and added the characters Joe, Mr. Jaggers, Magwitch and Orlick to the visualisation, combined with the location data and chapter delimiting. They created the image below and highlighted key points about it. The image on first view can be quite overwhelming – for further understanding of its constituent parts, see **A.5**.

Key Observations of Figure 4.7

- Joe’s presence at the opening and closing of the story.
- Joe’s fading presence and replacement with Estella/Herbert when Pip moves to London.
- Competition of Herbert/Estella for Pip’s attention – seen through occurrences of each character in turn from row 6.
- Prevalence of London as a location in Volume Two – about three quar-

ters of the way along row 3.

- Brief interludes of mentions of Blue Boar, and thus home, coinciding with absence of Herbert/Estella.
- Build-up of tension through character mentions at key points of the book – row 10 (Magwitch’s demise).
- Dissipation of tension after these events with presence of familiar characters – row 11 with Joe in the ante-penultimate chapter.
- Unsatisfactory nature of ending finishing with Estella, where Joe would make more sense – oft criticised and not ending Dickens’ originally wanted [37].

4.4.1 Usability

During the walkthrough the usability of the tool was discussed at length. This is vital for future work with the software but not for my evaluation. The most important points raised were the benefits of the adaptable granularity, especially when zooming in. They found this most interesting with the year-highlighting, allowing them to identify specific interesting words. They were also able to glean more information from the POS-highlighting when zoomed in.

We discussed at length the order of presenting additional characters, locations and parts of speech to the user and whether this should be done alphabetically or by frequency. The user could see the benefits of both, as characters in alphabetic list makes it easier to find slightly more obscure characters, and may encourage visualisations of characters you may not consider. Conversely, with parts of speech this is not something the user necessarily knows off the top of their head so to show relative frequency is perhaps more valuable.

Finally, the user discussed the potential of adding a key. When generating the images for use as narrative analyses they thought that this would be a valuable addition. They did mention how it could detract from the visual

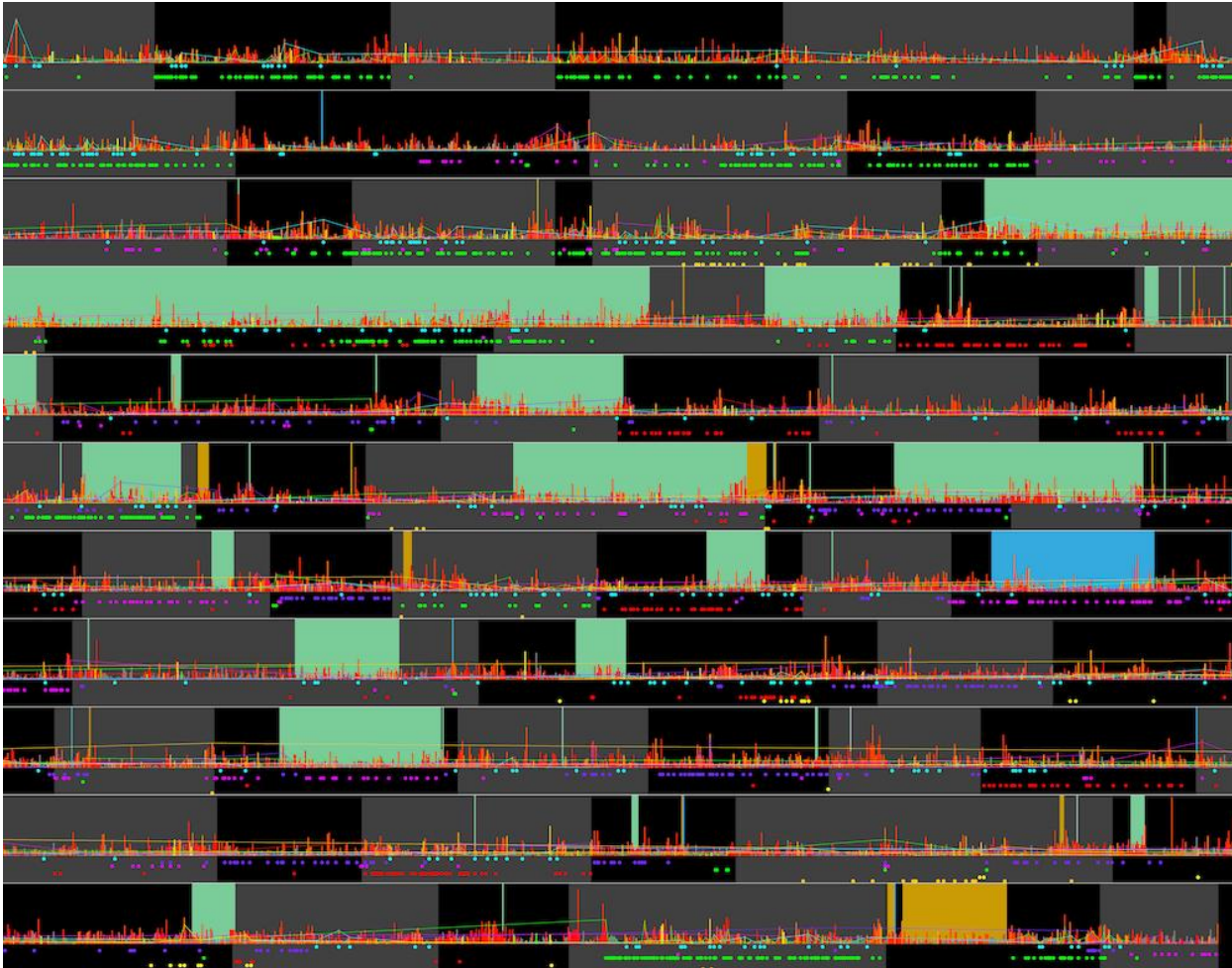


Figure 4.7: Visualisation of *Great Expectation* generated by user E in the cognitive walkthrough with default parameters set and the following information: aggregated year-highlighting, chapter-highlighting, character-highlighting with lines (Pip (cyan), Herbert (purple), Estella (magenta), Joe (green), Jaggers (red), Magwitch (yellow), Orlick (orange)); location-highlighting (London (green), Satis House (blue), Blue Boar Inn (orange)).

appeal of the image when viewing in more artistic settings. They also discussed that when generating the image themselves it did not feel necessary as they had just made the design decisions.

4.4.2 Conclusion

This walkthrough represents a significant success for the project: the user generated an informative narrative visualisation using the tool. This was enhanced as the visualisation told a narrative I had not considered, and echoed comments from user A in **3.2.2** about viewing characters as “*agents of plot development*”. Under scrutinisation it offers up limitations of the visualisations and how these could be improved.

Chapter 5

Conclusions

5.1 Improvements

A main area for improvement was the experiment. The results suggest confirmatory statistical evidence for claims made in the project, but due to insufficient recruitment, and reliance on familiarity with the novels, the sample size is generally not large enough to draw conclusive, statically significant findings. Instead, I must use the results from the experiment as supporting evidence for my qualitative analysis.

The lack of UI is explained throughout the project. This is compensated by having the functionality to interact with the tool in real-time in place meaning new projects only has to link this low-level functionality to a well-designed front-end. An in-depth UI would allow for further exploration of the visualisations from the users, in their own time.

5.2 Successes

Broadly, the work has been successful in accomplishing the aim of allowing explorative analysis of text. The users frequently referenced how the novelty

of visualising the whole book has allowed them to observe characteristics of the book they would struggle to see using traditional methods. This highlights an improvement of this work on Tim Regan’s work [15] as the images are not constrained to an individual medium and can be viewed digitally and in print form, utilising a row approach to present more information in a usable manner. This work distinguishes itself from artistic approaches [7] by presenting the visualised features accessibly. The project represents a progression from work such as TextArc [6] as it enables exploratory analysis of features determined by genuine potential users for real academic analysis. This is connected to the project’s clear aim of empowering literature students.

These points are exemplified by the production of figure 4.7 during the think-aloud study (4.4): a genuine example of a student creating a narrative analysis using the visualisation features I have designed and implemented.

There are certainly aspects that can be refined, such as the visualisation of location- and POS-highlighting, and the use of colour in the spectrum for displaying the word ages. However, none of the things I have approached in the project have been complete failures, and future work building on the initial investigations can address these flaws.

The decision to not pursue an interface is vindicated by the strong positive results I obtained from the images – both from collaborating users and the experiment. Though the majority of feedback and evaluation has been the opinions of a small selection of students – these users represent the target consumers of these visualisations. Moreover, I have addressed their feedback using proper interview technique and sought advice from professional researchers when assessing their responses.

Throughout the project the users have helped to shape images that have repeatedly been reported as informative and discernible – by them, and unfamiliar users. They have routinely expressed the value they can see in using these visualisations in tandem with their studies, in particular as starting points for further analysis, and to test theories on a whole text. The con-

clusions from my evaluation (**4.2.4**, **4.3.2** and **4.4.2**) together assert the success of this project in creating novel visualisation techniques to enable exploratory analysis of literature, and facilitating the generation of visual narrative analyses.

5.3 Future Work

The most important aspect of future work is to develop a UI to allow users to interact with the visualisations in a more in-depth and informed manner. By enabling users to perform exploratory analysis in their own time, without supervision could open up the tool to more, unforeseen, uses.

Further collaboration with students more finely interested in linguistics would be valuable to understand how better to draw visual information from the NLP techniques I have employed. This could be particularly valuable in work which looks at the use of language in social media [38] and the visualisations could be well applied in this field.

As highlighted in **3.6.2** the year-highlighting could provide interesting avenues for future work, especially the investigation of words without clear peak usages. Additionally I would also fine-tune the existing mechanisms and allow more customisation, such as colours and size of visual markers, in order to allow users to generate visualisations to their own aesthetic tastes and investigate the accessibility of the tool.

Appendix A

A.1 Initial Interview Responses

Positive	Negative
Could be useful for tracking frequency of words. (C)	Hard to see usefulness. (C)
Gives slightly more context. (C)	Difficult to interpret. (C)
Let's you see where things occur in the story. (B)	Without additional information (such as book/chapter number) not helpful. (A)
	Ugly. (B)
	Cannot necessarily interpret time as linear. (B)

Table A.1: Table of codified responses from initial interview with users' views on linear timelines.

Positive	Negative
Useful for key overall themes. (C)	Doesn't actually tell you very much. (C)
Good for repeated ideas. (C)	Stupid. (B)
Good for quick analysis of how frequently words that are discussed in literature are actually used. (C)	Not to do with book, but language. (B)
	No context. (B)

Table A.2: Table of codified responses from initial interview with users' views on word clouds.

Positive	Negative
Being able to split up into separate books in one image is very helpful. (C)	Length of sentences more interesting in verse poetry than prose. (C)
Length of sentences is interesting. (C)	For classical work, translation could have a big impact on this. (C)
Useful tool of analysis to see the shape of a text and how it moves. (A)	Could be difficult to understand without explanation. (A)
Useful for writing to ensure you have a good balance of sentence length. (A)	
It's cool. (B)	

Table A.3: Table of codified responses from initial interview with users' views on novel idea of structural interrogation.

Positive	Negative
Useful for some texts. (A)	Struggle to understand concept quickly. (C)
Really useful for analysing plot movement. (A)	Less useful for texts which are based in a sole location. (C)
Useful for looking at eco-criticism. (A)	
Good. (B)	

Table A.4: Table of codified responses from initial interview with users' views on novel idea of location-highlighting.

A.2 End of Iteration 1 Interview Responses

Positive	Negative
Easier to see lots of information in one go with columns. (A)	Loses some of the aesthetic with columns. (A)
Horizontal way of viewing is like a new way of reading. (A)	
Horizontal display allows for more natural reading style. (B)	
Horizontal display allows for quicker consolidation of information. (B)	

Table A.5: Table of codified responses from interview with users' views on structural interrogation presentation.

Positive	Negative
Even in speech, places are created in imagination, and discussing them is a way of being there. (A)	False negative highlighting when place is being discussed. (A)
Discussion of place and intervening places is interesting. (A)	
Good way of ascertaining culmination of actions in the book into key locations. (B)	

Table A.6: Table of codified responses from interview with users' views on location-highlighting approach.

Positive	Negative
Useful to see how author uses nouns and verbs etc. (C)	Could be difficult to interpret. (B)
Really useful for some forms of analysis. (A)	
Potential to show character's agency through the book. (A)	
Could show whether certain types of speech are associated with gendered characters. (A)	
Cool to see how things interact. (B)	
Combination of POS highlighting with characters and location could provide interesting insight. (B)	

Table A.7: Table of codified responses from interview with users' views on the potential POS-highlighting approach.

Positive	Negative
Use as a comparative tool could be very valuable. (C)	Not sure how well it would work. (A)
In combination with location could show age of language associated with rural/urban locations. (C)	
Interesting to see how close to the time of publication the language used is. (A)	

Table A.8: Table of codified responses from interview with users' views on the potential year-highlighting approach.

A.3 End of Iteration 2 Interview Responses

Positive	Negative
Horizontal more visually appealing. (E)	Not as analogous as genuine reading. (A)
Similar to skyline and ties in with location tracking. (E)	May not be one option that suits all media. (B)
Horizontal is a more natural way of reading. (A)	
More visually appealing. (A)	
You're reading, but one sentence at a time instead of one word. (A)	
Horizontal approach is a natural way of reading. (B)	

Table A.9: Table of codified responses from interview with users' views on vertical vs. horizontal presentation.

Positive	Negative
Easier to follow in horizontal version because fewer breaks in lines. (E)	
With lines it is easier to see groups together. (A)	
Draws the eye to it. (A)	
Allows you to link things together more easily. (B)	

Table A.10: Table of codified responses from interview with users' views on lines connecting character mentions together.

Positive	Negative
When looking with other information, works well. (E)	No best option as in situ, don't need margins. (E)
Makes clearer where words are appearing. (A)	
Good way of adding information whilst leaving bulk for main info. (B)	

Table A.11: Table of codified responses from interview with users' views on the margin approach.

A.4 Overview of Literature

A.4.1 Alice's Adventures in Wonderland (1865) by Lewis Carroll

“The story centres on Alice, a young girl who falls asleep in a meadow and dreams that she follows the White Rabbit down a rabbit hole. She has many wondrous, often bizarre adventures with thoroughly illogical and very strange creatures, often changing size unexpectedly (she grows as tall as a house and shrinks to 3 inches [7 cm]). She encounters the hookah-smoking Caterpillar, the Duchess (with a baby that becomes a pig), and the Cheshire Cat, and she attends a strange endless tea party with the Mad Hatter and the March Hare. She plays a game of croquet with an unmanageable flamingo for a croquet mallet and uncooperative hedgehogs for croquet balls while the Queen calls for the execution of almost everyone present. Later, at the Queen's behest, the Gryphon takes Alice to meet the sobbing Mock Turtle, who describes his education in such subjects as Ambition, Distraction, Uglification, and Derision. Alice is then called as a witness in the trial of the Knave of Hearts, who is accused of having stolen the Queen's tarts. However, when the Queen demands that Alice be beheaded, Alice realizes that the characters are only a pack of cards, and she then awakens from her dream.”¹

¹<https://www.britannica.com/topic/Alices-Adventures-in-Wonderland>

A.4.2 The Story of Doctor Dolittle (1920) by Hugh Lofting

“The Story of Doctor Dolittle is the first of his Doctor Dolittle books, a series of children’s novels about a man who learns to talk to animals and becomes their champion around the world.”²

A.4.3 Great Expectations (1861) by Charles Dickens

“Great Expectations is the story of Pip, an orphan boy adopted by a blacksmith’s family, who has good luck and great expectations, and then loses both his luck and his expectations. Through this rise and fall, however, Pip learns how to find happiness. He learns the meaning of friendship and the meaning of love and, of course, becomes a better person for it.”³

A.4.4 Strange Case of Dr. Jekyll and Mr. Hyde (1886) by Robert Louis Stevenson

“Strange Case of Dr Jekyll and Mr Hyde by Robert Louis Stevenson is a narrative about the complexities of science and the duplicity of human nature. Dr Jekyll is a kind, well-respected and intelligent scientist who meddles with the darker side of science, as he wants to bring out his ‘second’ nature. He does this through transforming himself into Mr Hyde - his evil alter ego who doesn’t repent or accept responsibility for his evil crimes and ways. Jekyll tries to control his alter ego, Hyde, and for a while, Jekyll has the power. However, towards the end of the novel, Hyde takes over and this results in their deaths.”⁴

A.4.5 Pride and Prejudice (1818) by Jane Austen

“Pride and Prejudice, romantic novel by Jane Austen, published anonymously in three volumes in 1813. A classic of English literature, written with

²<https://etc.usf.edu/lit2go/221/the-story-of-doctor-dolittle/>

³<https://www.gradesaver.com/great-expectations/study-guide/summary>

⁴<https://www.bbc.com/bitesize/guides/z88wjxs/revision/1>

incisive wit and superb character delineation, it centres on the turbulent relationship between Elizabeth Bennet, the daughter of a country gentleman, and Fitzwilliam Darcy, a rich aristocratic landowner.”⁵

A.4.6 The War of the Worlds (1897) by H.G. Wells

“The story, which details 12 days in which invaders from Mars attack the planet Earth, captured popular imagination with its fast-paced narrative and images of Martians and interplanetary travel. The humans in *The War of the Worlds* initially treat the invasion with complacency but soon are provoked into a defensive state of war.”⁶

⁵<https://www.britannica.com/topic/Pride-and-Prejudice>

⁶<https://www.britannica.com/topic/The-War-of-the-Worlds-novel-by-Wells>

A.5 Progression of Visualisation in Cognitive Walkthrough

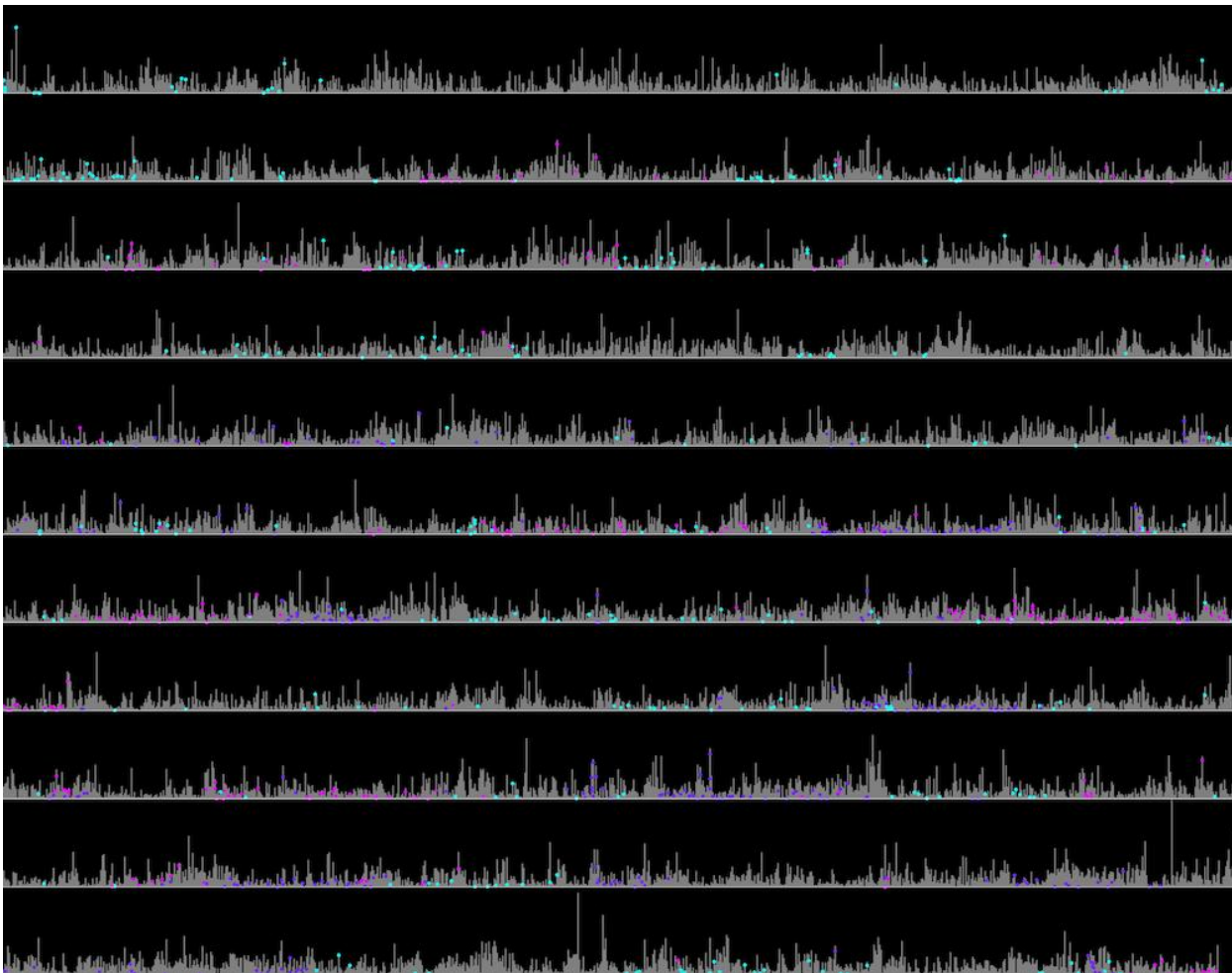


Figure A.1: First visualisation created by the user highlighting characters Pip (cyan), Estella (magenta) and Herbert (purple).

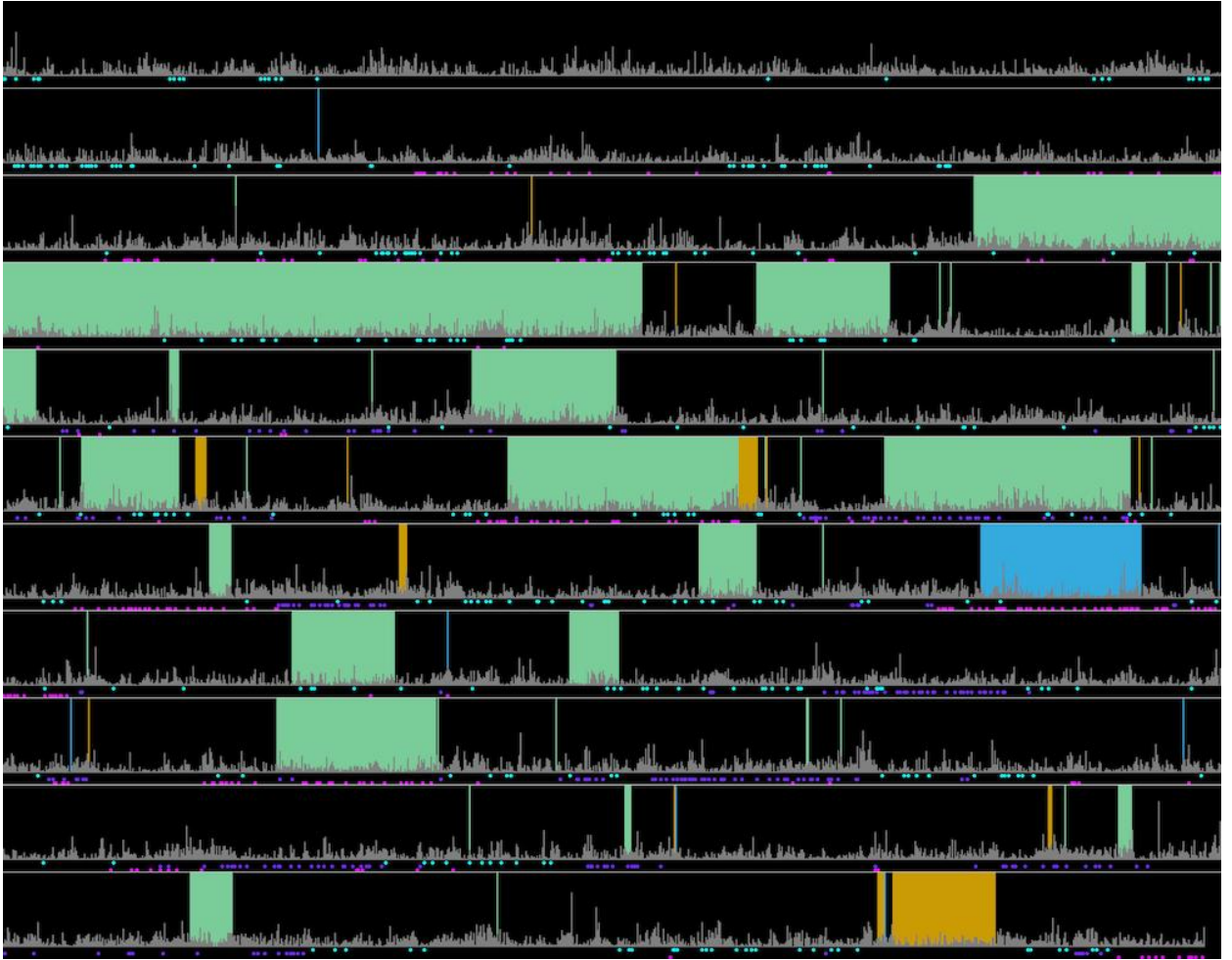


Figure A.2: Second visualisation created by the user highlighting characters Pip (cyan), Estella (magenta) and Herbert (purple); locations London (green), Satis House (blue) and Blue Boar Inn (orange).

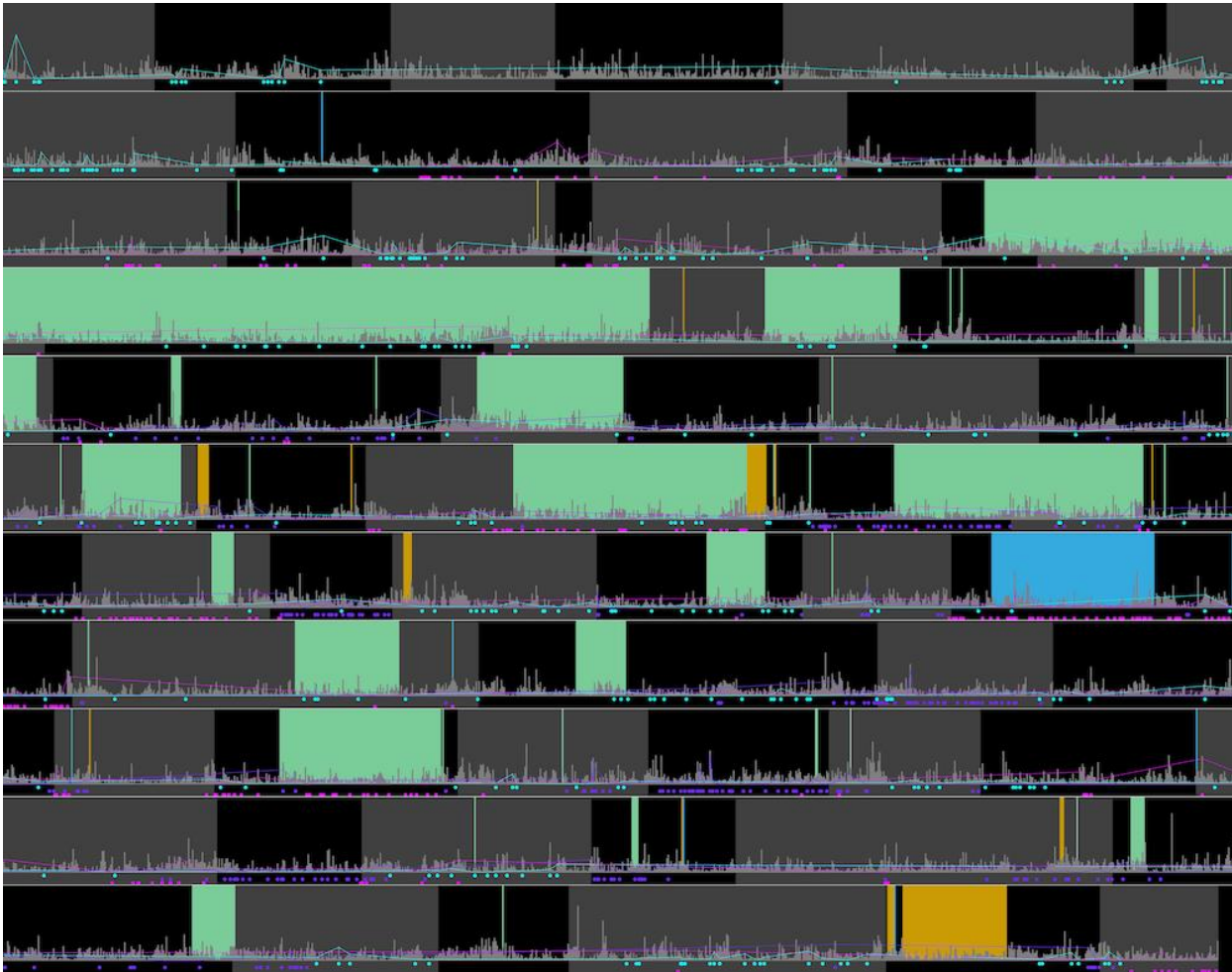


Figure A.3: Third visualisation created by the user highlighting the same as above and turning on chapter highlighting.

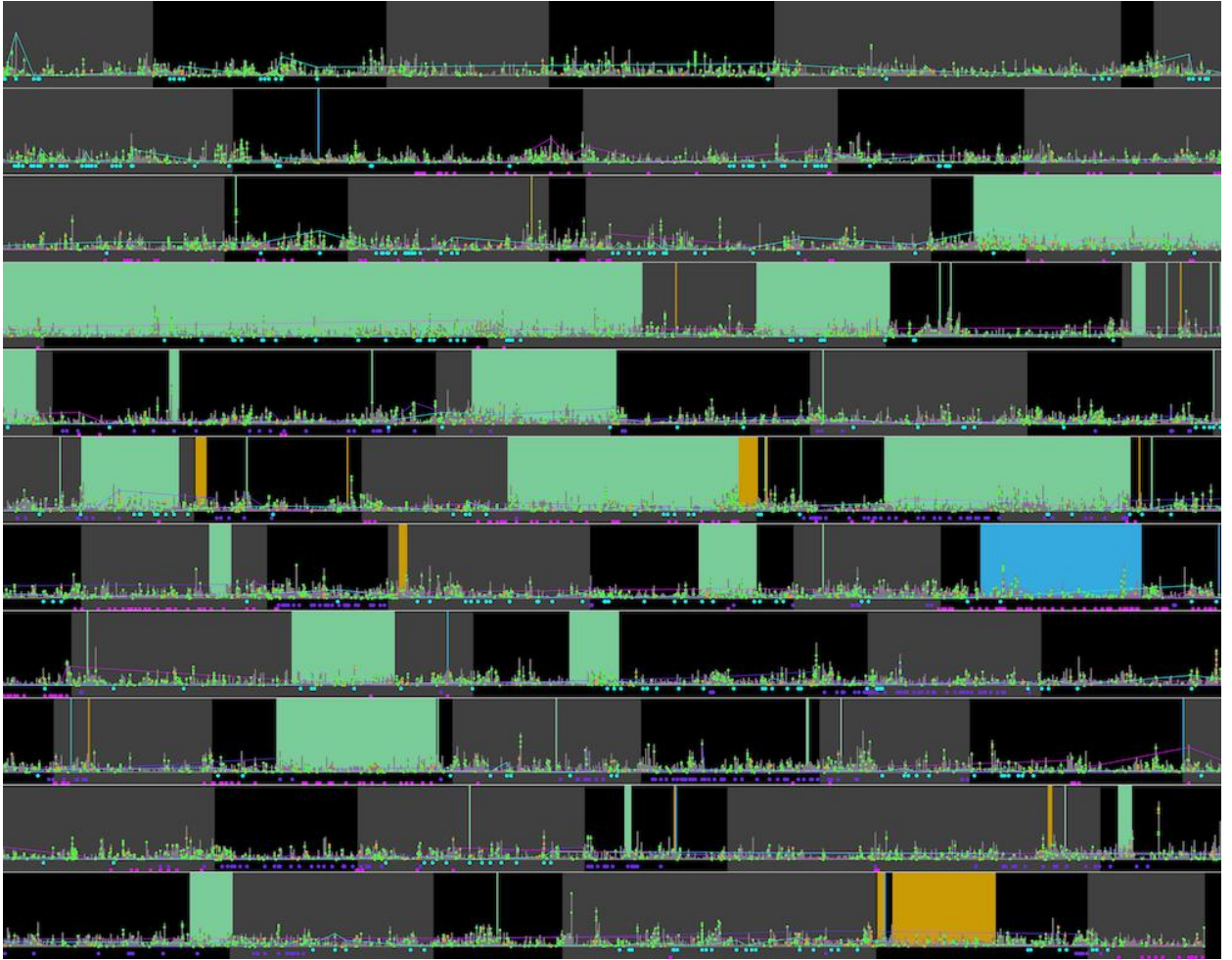


Figure A.4: Fourth visualisation created by the user highlighting the same as previously and attempting to highlight POS through personal pronouns (green) and comparative adjectives (orange).

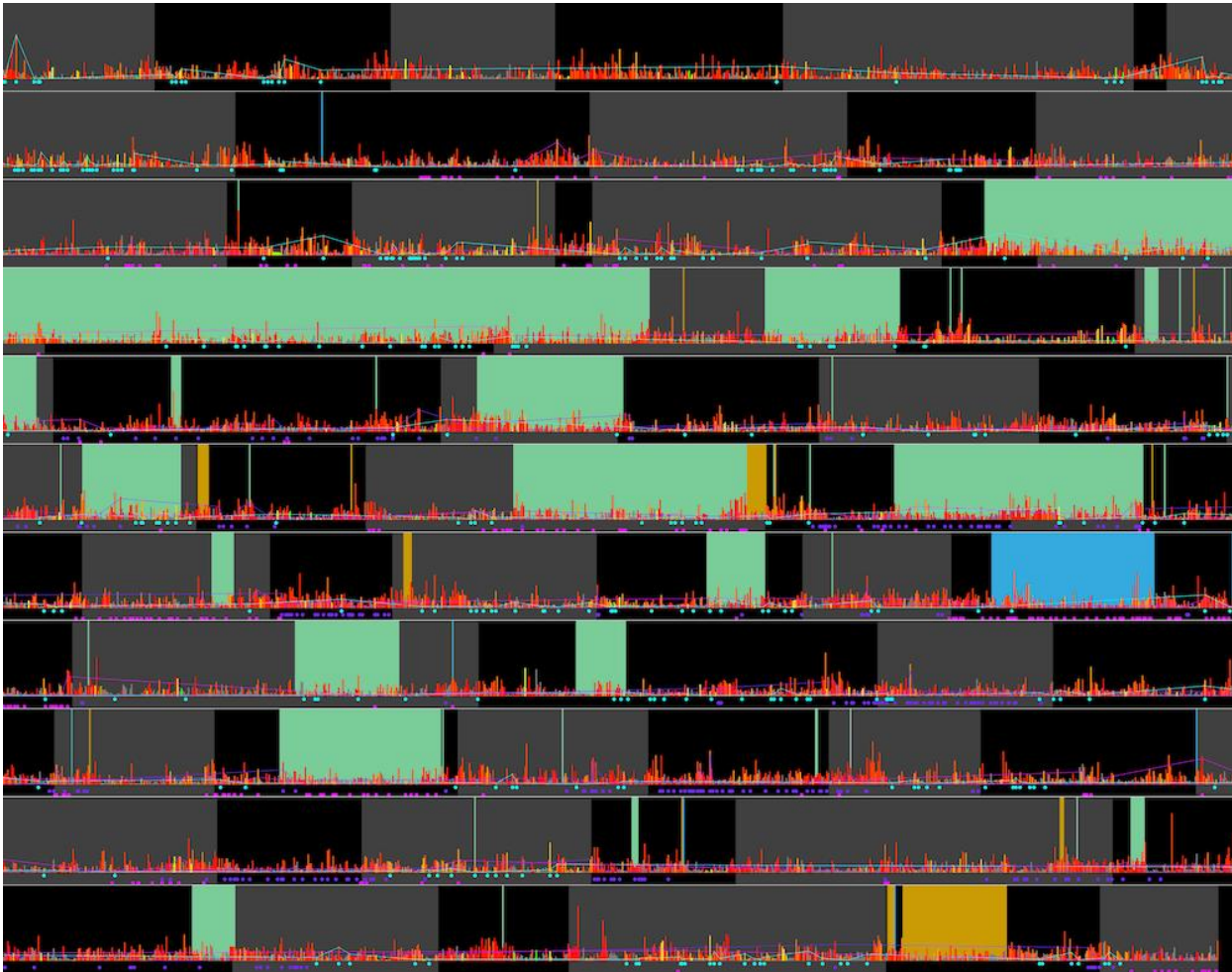


Figure A.5: Fifth visualisation, removing POS-highlighting and adding aggregated year-highlighting.

A.6 Final Summative Interviews

A.6.1 Questions

Section 1

What do you think about the basic image?

What do you think about chapter-highlighting?

What do you think about character-highlighting?

What do you think about location-highlighting?

What do you think about year-highlighting?

What do you think about POS-highlighting?

What do you think about the use of margins?

What do you think about the combination of features into one image?

Section 2

How does the visualisation I have created compare to previous ‘traditional’ visualisation techniques?

Can you analyse a text purely using the visualisation?

In what ways is analysis using the tool similar to how you would normally analyse text?

In what ways is it different?

Section 3

Can you identify a book just from its visualisation? If so, what features help you to do so?

Can you give any specific examples of features which surprised you in visualisations?

Section 4

Do you think this tool would have been useful during your studies?

Can you see yourself using a tool like this in your studies in the future?

What do you see as the primary use of it?

A.6.2 Transcriptions

User A

Date : 19.05.2019

Interviewer(I)

Respondent : User A (A)

This interview was not carried out in person due to unavailability of the respondent. Therefore, the interviewer was provided with the questions and example images and responded with a written document.

Section 1

I: What do you think about the basic image?

A: [*No specific comment.*]

I: What do you think about chapter highlighting?

A: Chapter highlighting would be especially useful as it would help people locate things in the book to go and analyse more.

I: What do you think about character highlighting?

A: [*No specific comment.*]

I: What do you think about location highlighting?

A: [*No specific comment.*]

I: What do you think about year highlighting?

A: Only thing with year one is maybe a little bit unclear because of all the colours but still good for getting an overview.

I: What do you think about POS highlighting?

A: [*No specific comment.*]

I: What do you think about the use of margins?

A: It is good with the margins as it makes it a bit clearer.

I: What do you think about the combination of features into one image?

A: Really like all the highlighting stuff, think they're all really effective

Section 2

I: How does the visualisation I have created compare to previous 'traditional' visualisation techniques?

A: I think your visualisation is actually much more helpful than the other ones because it can give a bit more detail and include multiple aspects at once.

I: Can you analyse a text purely using the visualisation?

A: I think the visualisation would be a really helpful starting point for analysis because it gives you an overview of interesting things and patterns which you could then do more in-depth analysis about.

I: In what ways is analysis using the tool similar to how you would normally analyse text?

A: Spotting patterns is a really key part of building an argument and the visualisation is perfect for that.

I: In what ways is it different?

A: It gives insight into the whole book whereas you might normally just focus on one section. It also gives a good sense of how different things like location and character interconnect which you might miss otherwise.

Section 3

I: Can you identify a book just from its visualisation? If so, what features help you to do so?

A: You could identify a book just from the visualisation based on locations, character movement etc. but I guess you would have to be pretty familiar with it beforehand.

I: Can you give any specific examples of features which surprised you in visualisations?

A: One thing that surprised me about the character analysis is that there were some quite long gaps where none of the characters seemed to appear - I thought this was interesting as they are some of the most significant characters yet a lot of the story seems to happen without them.

Section 4

I: Do you think this tool would have been useful during your studies?

A: I think it would have been super useful to me. Wish you'd come up with it a couple of years ago.

I: Can you see yourself using a tool like this in your studies in the future?

A: [*No specific comment.*]

I: What do you see as the primary use of it?

A: Interesting if you wanted to analyse agency or something. When starting out on an essay and identifying patterns you could explore further. It would be a good way to back up an argument with statistical data rather than just relying on more qualitative analysis.

User B

Date : 17.05.2019 Interviewer(I) Respondent: User B (B)

Section 1

I: What do you think about the basic image? Showed user Pride and Prejudice image.

B: It's interesting because the lines of the whole image are the way you read but then sentences are like the opposite way to the way the sentences would be.

I: Does it tell you anything interesting from the sentence length?

- B: Yeah, well it is interesting to see the patterns, but without other information it doesn't really tell you that much.
- I: What do you think about chapter highlighting?
- B: I guess with the length of sentence for speech you can see which chapters have more speech and how it is broken up.
- I: What do you think about character highlighting?
- B: You can kind of see who is speaking to each other which is interesting. Can also see who the focus of the book is at that point. You can see that it is a third person book.
- I: What do you think about location highlighting?
- B: [*No specific comment.*]
- I: What do you think about year highlighting?
- B: [*No specific comment.*]
- I: What do you think about POS highlighting?
- B: [*No specific comment.*]
- I: What do you think about the use of margins? Especially the dots in the margins
- B: Makes things more visually interesting and clearer because they're connected but also not. And you can look at the dots as rows of themselves or related to the sentence lines.
- I: What do you think about the combination of features into one image? Does it work? Or is it too busy?
- B: No think it does work especially as the visual domains are kind of separated.

Section 2

- I: How does the visualisation I have created compare to previous 'traditional' visualisation techniques?
- B: There's more details and you can tell more from it. And just having a word written bigger for example is not as visually interesting as the images you have.

- I: What do you think makes the new images more informative?
- B: You can see the frequency like as a whole book and you work it out for yourself.
- I: Can you analyse a text purely using the visualisation?
- B: In what sense?
- I: Do you think you can make observations and points about a book just from the images?
- B: Like you can make interesting observations about the book but to go any further you have to have read the book. It's more to be used in addition to having read it. Doesn't necessarily mean you have to have read the book, for example someone could try and look at how much women speak in the book and you wouldn't need to know the text and this would be useful for that. It is versatile and you can take a lot of interpretations from it, but still important to have read the book.
- I: In what ways is analysis using the tool similar to how you would normally analyse text?
- B: It is a kind of reading because you are interpreting the shapes instead of words. And similar use of the margins for annotations.
- I: In what ways is it different?
- B: It makes you think about in a completely different way. I wouldn't consider it in terms of numbers normally. It is less work. This is like a map of where to go and you can further analyse. It draws your attention to things you would not have noticed before.

Section 3

- I: Can you identify a book just from its visualisation? If so, what features help you to do so?
- B: You can use all of it, but using it in combination and using the guide helps you work out which locations and therefore which book. You can identify how many main characters there are. You can see the main character's omnipresence in a third person book, like in Alice in Wonderland. Can see coming and going of secondary characters, and when they are interacting with other

characters. Prevalence of location and use of it as a theme is interesting. In Great Expectations repeated use of London in volume 2 is valuable. And the length of the book is quite apparent from the image which helps to identify.

I: Can you give any specific examples of features which surprised you in visualisations?

B: The word usage was really interesting, didn't think it would be so easily visible. The other aspects make sense but the usage was surprising. Containment of characters to their scene is really interesting.

Section 4

I: Do you think this tool would have been useful during your studies?

B: Tended to go for more weird analysis than number analysis but can see how it would be useful for some people. I would have found it interesting either way, but I haven't really focussed on the things that you can see from it. Can definitely see the use of it as a starting point. I would also have thought of it more as an image that can represent things. Alice in Wonderland has illustrations and this is very different but still visual. The logic of Alice in Wonderland is very calculated and based on logic games, so is interesting if you're doing computery things.

I: Can you see yourself using a tool like this in your studies in the future?

B: Won't be studying in the future, but yes because of the reasons above. Utilising it more as a representation of the book.

I: What do you see as the primary use of it?

B: Still agree with use as visual art, but now it is more clear that you can interpret it more easily and there is more information on it it is also useful for things like finding parts of books you might investigate. Since it has started it has become something a user can interpret much more easily. And you can still pile lots of information on to make pretty pictures.

User E

Date : 18.05.2019 Interviewer(I) Respondent: User E (E)

Section 1

I: What do you think about the basic image? Showed user Pride and Prejudice image.

E: Without anything overlaid doesn't mean that much to me, without the context. I feel like I don't have that much to say about this one. Sentence length is surprising because expect to see it all over the place, and interesting to see it squished down. More uniform that is ? especially in Great Expectations.

I: What do you think about chapter highlighting?

E: Can compare the lengths, again, without knowing the context of where in the story the chapters are happening, I find it difficult to glean much meaning. It is a lot more aesthetically pleasing though I would say. It makes it feel a bit more like a book which you are looking at.

I: What do you think about character highlighting?

E: Already I feel like I could look at this for a while with different questions and come up with different responses. I don't think it makes me 'feel' a certain way about the book, but I don't think any of the pictures do. All of a sudden there's like potential. (Referring to Great Expectations) . And immediately it is striking that you have Joe at the beginning and right at the end. Interesting how much more weighted it is towards the beginning. From my memory of the book feels very important at the end, but actually it is more of a flash. And then obviously you end with Estella instead of Joe, which you start with. Maybe hints at the fact that the ending with Estella is very satisfying. People don't think of Great Expectations as this great romance. It is good at highlighting that we are more emotionally attached to Joe, for good reason. Estella properly vanishes.

I: What do you think about location highlighting?

E: Doesn't really tell me much that I didn't already know, just from my memory. It's mainly just London that you can see (in Great

Expectations) it doesn't really change much for me. We know he's in London in the middle. In a way I feel this is the feature that could be misinterpreted the most without knowledge of the book or the techniques used. I get it because I have seen it and understand it.

I: What do you think about year highlighting?

E: This is one that depends a lot more on the specific book you're look at. For Great Expectations it is not massively interesting or significant to look at, but as soon as you have something like War of the Worlds which is explicitly engaged with new technology then it throws up a lot more interesting ideas. This is one where looking at it as a whole, I don't think - oh, this means that. If I were to look at it on a smaller scale would be interesting to see the bits where the language suddenly becomes a lot more modern. It is really good for comparison between books.

I: In Pride and Prejudice she basically only uses words of the era apart from a couple of points where she uses modern words, such as 'someone' and 'relationship'

E: That's really interesting, because with out an in-depth knowledge of language you wouldn't necessarily know that these words were not of the time, especially as you probably won't be reading lots of books from that era. This is really important, because when we romanticise relationships in modern day media it still so often comes from something like one of Jane Austen's novels and it's so easy to forget that actually for her time was pretty modern. That is the only thing that might be slightly misleading from the averaged view. Words like relationships being a newfangled thing is pretty interesting. Can also see how an author's vocabulary changes over time.

I: What do you think about POS highlighting?

E: It's like what you would expect to see vs. what actually happens and it's just that these parts of speech are actually quite common, and especially in a so-called realistic medium, like novels. When you have the long full sentences it is really difficult to make any kind of judgement based on that information. Especially as I don't have a great working knowledge of the parts of speech.

I: What do you think about the use of margins?

- E: It's helpful when you are looking at more information than just the characters.
- I: What do you think about the combination of features into one image?
- E: Yes, it is quite a lot to handle, if you could look at each section instead of absorbing as one whole thing it would be good. It is pretty easy to understand, especially if you have it laid on bit by bit. If I wanted to look at chapters or characters, you don't have to strip away the other information because it is like clearly demarcated.

Section 2

- I: How does the visualisation I have created compare to previous 'traditional' visualisation techniques?
- E: More information and the immediate context of having the entire book there in some form which makes a big difference I think.
- I: Can you analyse a text purely using the visualisation?
- E: No, I think for me I don't see how I could make basic sense of it without knowledge of the plot.
- I: Do you think you can analyse the image with prior knowledge of the book?
- E: Yes, definitely. There's a lot there that we've already talked about. I think the occurrence of characters and how they combine with characters is something I am definitely drawn to on a personal level. Maybe not in *Great Expectations* but in other ones locations could be interesting, and you can see how it interacts with the building up of tension. With *Great Expectations* that's definitely seen in the way that characters bunch up together towards the end. There's definitely stuff that has come up or been revealed in a very quick and obvious way from the images. Wouldn't have felt confident making comments about Herbert and Estella before and that was more of a hunch and something I was interested in and then you can really quickly see that yeah that is the case. That is something you can only really do with visualising it. Even if you read it you would still have to like note down every time the name came up.

- I: In what ways is analysis using the tool similar to how you would normally analyse text?
- E: I guess it's like reading a page superficially, but doesn't really translate to actually being comparative to how I would usually analyse text.
- I: In what ways is it different?
- E: Fundamentally quite different, gives you access to like a different scale. Gives you a much more distant perspective, you can see it as a whole, in a way in which you just can't usually unless you know that book inside and out. For me, with novels that is something I am really not good at.

Section 3

- I: Can you identify a book just from its visualisation? If so, what features help you to do so?
- I: With the caveat that you are moderately familiar with the book.
- E: Yes, definitely with some exposure to the images and knowledge of the books. The length of the book helps to distinguish it. With War of the Worlds I know there will be more colours in the words. I would look for that if trying to distinguish it from other long novels. Locations help as well. The interaction of characters is also very useful. Can identify if a book is third person too from the use of character names, with Alice in Wonderland you are expecting to see just her.
- I: Can you give any specific examples of features which surprised you in visualisations?
- E: I guess I was definitely surprised at the efficiency of the word age thing. I think that is one you have to see it to understand it.

Section 4

- I: Do you think this tool would have been useful during your studies?
- E: Yeah. I think something I use all the time for text, particularly something like Shakespeare, you can get the full text or all his

work and then search for usage of words. To be able to visualise that would make a massive difference, because you don't know if sometimes you're reading into it. But if you could look at several plays and really quickly see that in this play he uses it a bunch and in others he doesn't. If you're only looking at one play it's difficult without the context. Even when you have the computer searching it's still not as efficient as visualising, and you are more likely to bring your bias to it.

I: Can you see yourself using a tool like this in your studies in the future?

E: If I were going to study any more then yeah definitely.

I: What do you see as the primary use of it?

E: It is difficult to choose one aspect of it. The characters is probably what is most interesting to me. In general objectifying things that are hunches is really powerful. Being able to see the whole context is really valuable too and is a good way of consolidating ideas about the book as a whole.

Bibliography

- [1] Eleanor Drago-Severson, Deborah Helsing, Robert Kegan, Maria Broderick, Nancy Popp, and Kathryn Portnow. Three developmentally different types of learners. *Focus on basics*, 5(B):7–9, 2001.
- [2] Abigail J Sellen and Richard HR Harper. *The myth of the paperless office*. MIT press, 2003.
- [3] Thomas Kvan. Collaborative design: what is it? *Automation in construction*, 9(4):409–415, 2000.
- [4] Riitta Jääskeläinen. Think-aloud protocol. *Handbook of translation studies*, 1:371–374, 2010.
- [5] Kostiantyn Kucher and Andreas Kerren. Text visualization techniques: Taxonomy, visual survey, and community insights. In *2015 IEEE Pacific Visualization Symposium (PacificVis)*, pages 117–121. IEEE, 2015.
- [6] W Bradford Paley. Textarc: Showing word frequency and distribution in text. In *Poster presented at IEEE Symposium on Information Visualization*, volume 2002, 2002.
- [7] Stephanie Posavec. Writing without words. *Stefanie Posovec*, 2006.
- [8] Boris Mueller. Visualpoetry-generative graphic design for poetry on the road. In *ACM SIGGRAPH 2009 Art Gallery*, page 41. ACM, 2009.
- [9] Natalia Bilenko. *The narrative explorer*. PhD thesis, Master's thesis, EECS Department, University of California, Berkeley, 2016.
- [10] Catherine Kohler Riessman. *Narrative analysis*, volume 30. Sage, 1993.
- [11] Chong Ho Yu. Exploratory data analysis. *Methods*, 2:131–160, 1977.
- [12] Clarence Larkin. *Dispensational truth*. Delmarva Publications, Inc., 2014.

- [13] Ben Fry. The preservation of favoured traces. <https://fathom.info/traces/>. Accessed: 2019-05-25.
- [14] Charles Darwin. *On the origin of species, 1859*. Routledge, 2004.
- [15] Tim Regan and Linda Becker. Visualizing the text of philip pullman’s trilogy his dark materials. In *Proceedings of the 6th Nordic Conference on Human-Computer Interaction: Extending Boundaries*, pages 759–764. ACM, 2010.
- [16] Martin Wattenberg, Fernanda B Viégas, and Katherine Hollenbach. Visualizing activity on wikipedia with chromograms. In *IFIP Conference on Human-Computer Interaction*, pages 272–287. Springer, 2007.
- [17] Lucie Flekova, Florian Stoffel, Iryna Gurevych, and Daniel Keim. Content-based analysis and visualization of story complexity. *VISUALISIERUNG*, page 185, 2018.
- [18] Kostiantyn Kucher, Teri Schamp-Bjerede, Andreas Kerren, Carita Paradis, and Magnus Sahlgren. Visual analysis of online social media to open up the investigation of stance phenomena. *Information Visualization*, 15(2):93–116, 2016.
- [19] Craig Larman and Victor R Basili. Iterative and incremental developments. a brief history. *Computer*, 36(6):47–56, 2003.
- [20] Mats Alvesson. *Interpreting interviews*. Sage, 2010.
- [21] Blue letter bible. <https://www.blueletterbible.org>. Accessed: 2019-05-27.
- [22] Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60, 2014.
- [23] Peter Willett. The porter stemming algorithm: then and now. *Program*, 40(3):219–223, 2006.
- [24] Daniel Cohen-Or, Olga Sorkine, Ran Gal, Tommer Leyvand, and Ying-Qing Xu. Color harmonization. In *ACM Transactions on Graphics (TOG)*, volume 25, pages 624–630. ACM, 2006.
- [25] Gunther Kress and Theo Van Leeuwen. Colour as a semiotic mode: notes for a grammar of colour. *Visual communication*, 1(3):343–368, 2002.

- [26] Stephen Westland, Kevin Laycock, Vien Cheung, Phil Henry, and Forough Mahyar. Colour harmony. *JAIC-Journal of the International Colour Association*, 1, 2012.
- [27] Muzammil Khan and Sarwar Shah Khan. Data and information visualization methods, and interactive mechanisms: A survey. *International Journal of Computer Applications*, 34(1):1–14, 2011.
- [28] Martin Davies. Concept mapping, mind mapping and argument mapping: what are the differences and do they matter? *Higher education*, 62(3):279–301, 2011.
- [29] Patricia A Chalmers. The role of cognitive theory in human–computer interface. *Computers in human behavior*, 19(5):593–607, 2003.
- [30] Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann. Online controlled experiments at large scale. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1168–1176. ACM, 2013.
- [31] Mario Luis Small. How many cases do i need?’ on science and the logic of case selection in field-based research. *Ethnography*, 10(1):5–38, 2009.
- [32] Graeme D Ruxton and Markus Neuhäuser. When should we use one-tailed hypothesis testing? *Methods in Ecology and Evolution*, 1(2):114–117, 2010.
- [33] Herbert D Kimmel. Three criteria for the use of one-tailed tests. *Psychological Bulletin*, 54(4):351, 1957.
- [34] Marsha E Fonteyn, Benjamin Kuipers, and Susan J Grobe. A description of think aloud method and protocol analysis. *Qualitative health research*, 3(4):430–441, 1993.
- [35] Juliet M Corbin and Anselm Strauss. Grounded theory research: Procedures, canons, and evaluative criteria. *Qualitative sociology*, 13(1):3–21, 1990.
- [36] Marilyn Hughes Blackmon, Peter G Polson, Muneo Kitajima, and Clayton Lewis. Cognitive walkthrough for the web. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 463–470. ACM, 2002.
- [37] Edwin M Eigner. Bulwer-lytton and the changed ending of great expectations. *Nineteenth-Century Fiction*, 25(1):104–108, 1970.

- [38] H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Lukasz Dziurzynski, Stephanie M Ramones, Megha Agrawal, Achal Shah, Michal Kosinski, David Stillwell, Martin EP Seligman, et al. Personality, gender, and age in the language of social media: The open-vocabulary approach. *PloS one*, 8(9):e73791, 2013.