

Text analysis and visualisation: creating deep interfaces to read textual document collections

Jaume Nualart Vilaplana

May 2016

A thesis submitted for the degree
of Doctor of Philosophy in Communication
University of Canberra

||*||

Contents

1. Introduction	8
1.1. Context and motivation	8
1.2. Aims and goals	9
1.3. Practice and research	10
1.4. Dissertation overview	11
2. Context and background	13
2.1. Digital libraries and text collections	15
2.2. Text analysis and text collections	20
2.2.1. Text clustering and interfaces	24
2.2.2. Topic models	25
2.2.3. Topic model labelling	27
2.2.4. Topic model evaluation	31
2.3. Interfaces and text collections	33
2.3.1. Interface paradigms	34
2.3.2. Standard and trending practice in text collection interfaces .	37
2.3.3. Data visualisation and interfaces to text collections	40
2.3.4. Interfaces and reading	44
2.4. Summary	48
3. Gaps, research questions and contribution to Knowledge	50
3.1. Gaps	50
3.2. Research questions and Contribution to knowledge	54
4. Methodologies and methods	56
4.1. Methodologies	56

4.2. Methods	58
4.2.1. Tools and techniques	58
4.2.2. Evaluation and validation	59
5. Results: the artefacts	62
5.1. Visference: Journal of Machine Learning Research	64
5.2. Crossreads I: Eugeni Bonet exhibition	68
5.3. Crossreads II: "In your computer", by D Quaranta	72
5.4. Diggers Diaries I: WWI Diaries	75
5.5. Diggers Diaries II: WWI Diaries	79
6. Discussion	82
6.1. The miraculous analysis	84
6.2. The fascination for the data	92
6.3. The interface and its codes	96
6.4. Deep interfaces	101
7. Conclusions and future	104
Bibliography	111
A. appendix	125

Abstract

This research brings together data analysis with software engineering and visualization, with a specific focus on text mining and large document collections. My aim is to devise new, rich and simple visualization interfaces, which I call deep interfaces.

By deep interfaces I introduce the idea rich contents as a product of the statistical analysis combined with human curation of labels, and is interpreted as a flow of subjectivity, complexity and diversity between reader and interface and vice versa.

The focus of these interfaces is not the representation of textual document collections as in Moretti's distant reading, but to revisit traditional reading from the point of view of state of the art methods of textual analysis. Thus, the proposed interfaces can help us discover and explore text document collections by reading their contents. This is a practice-led research project that develops theoretical issues through the generation of practical artefacts. The research process is cumulative, following a reflexive methodology. The key outcomes of the project are embodied in an interface to a large collection of ANZAC war diaries: Diggers Diaries — <http://diggersdiaries.org> .

Acknowledgements

To Jordi, Kiesse, and Amelia. To Maria Rosa and Jaume. To the rest of my family. To my best friends.

Thanks to Dr Gabriela Ferraro, for her support, collaboration, and friendship; to Dr Joelle Vandermensbrugghe for her guidance to PhD students, her advice, and her friendship; and to Dr Giulio Zambon for his advice, conversations, and friendship.

I'd like to thank my colleagues from the Machine Learning Research Group of NICTA for their help, support, and interest shown in relation to my research project. Special thanks to my former supervisors: Dr Wray Buntine, Dr Mark Reid, and Dr Hanna Suominen; and to Dr Bob Williamson, the leader of the group, for his support and trust.

I'd like to thank the Faculty of Arts and Design staff for their help and professionalism.

Thanks to Dr Rob Fitzgerald for his warm support and the interesting and challenging project we shared.

Finally I'd like to thank Dr Mitchell Whitelaw, my thesis supervisor, who has shown and taught me how academic research works. Mitchell has introduced me to the academic discipline of digital humanities, a territory where we have shared our passion for data visualisation.

FORM B

Certificate of Authorship of Thesis

Except where clearly acknowledged in footnotes, quotations and the bibliography, I certify that I am the sole author of the thesis submitted today entitled –

Text analysis and visualisation: creating deep interfaces to read textual document collections

I further certify that to the best of my knowledge the thesis contains no material previously published or written by another person except where due reference is made in the text of the thesis.

The material in the thesis has not been the basis of an award of any other degree or diploma except where due reference is made in the text of the thesis.

The thesis complies with University requirements for a thesis as set out in the *Examination of Higher Degree by Research Theses Policy*. Refer to <http://www.canberra.edu.au/current-students/current-research-students/hdr-policy-and-procedures>



Candidate's Signature

15 / May / 2016
Date



Primary Supervisor's Signature

16 / May / 2016
Date

NOTE: The wording contained in Form B must be bound into the thesis, preferably as the third page, and signed by the author of the thesis and the supervisor. However, the exact layout does not need to be duplicated.

1. Introduction

In 2015 there was a growing amount of accessible data, we know its growing speed (Cisco, 2016, Bono, 2015), we know its size (Pappas, 2016), we can plan future information infrastructure needs to store it, but what we do not know is how we are going to use this amount of data. This thesis aims to contribute new methods we can use to access that data, focusing on textual documents that are grouped into digital collections.

This chapter briefly presents four questions that are fundamental to understand the nature of this research project. These are firstly, the context and motivation that informs this research on interfaces to textual document collections. Secondly it outlines the project's aims and goals. Thirdly, this chapter introduces the practice-based and practice-led nature of the research and its corresponding methodology and process. Finally this chapter offers an overview of the whole project, setting out the structure and topics of the following chapters.

1.1. Context and motivation

We have moved from information overload in the 1990's to big data. Information overload was presented as the biggest problem of the digital age and the necessity of computational help to deal with it was pointed out (Maes et al., 1994). Nowadays, in the age of big data, this excess is seen as an opportunity to index, mine, analyse, and discover insights, patterns, and ultimately solutions to information related problems (Mayer-Schönberger and Cukier, 2013).

This project is developed in the context of big and accessible data, and the unsolved problems we deal with due to these large amounts of data. Specifically, this research focuses on how we can deal with large collections of textual documents. These texts can be found in major public institutions that host digital collections,

usually accessible online with a text-based interface. This project explores the possibilities of implementing new techniques to improve the way we interact with these large collections of text.

According to that goal this project has some necessary "preconditions", such as the availability of digital document collections contents. While the collections are catalogued online, the contents of these collections are not always accessible. different contents from different knowledge domains, have different degrees of access and are used for different proposes. This project works with text collections from academic and cultural domains, where full text documents are increasingly available. Full text access to the contents of text collections is required also, in order to apply state-of-the-art techniques of text analysis that can enrich what we know about the collections. Text analysis techniques such as topic modelling are increasingly mature and well developed. Due to new easy-to-use software, these techniques are increasingly being used in digital humanities (DH) and other fields. In addition, web browsers are increasingly powerful and the technical standards are solid, hence, it is feasible to develop rich interactive displays.

From a more personal point of view, the urgent need of tools to deal with the new digitised collections is related to education, culture, and eventually, personal and collective freedom. In that sense, my experience over the last fifteen years working with free media, public libraries, and free knowledge, is so far integrated in the project's goals, as to make digital collections, and especially text collections, more accessible using a combination of analysis and visualization techniques.

1.2. Aims and goals

This research project seeks to bring together data and text mining with software engineering and visualisation, with a specific focus on text mining and large document collections. Its aim is to apply text mining and analysis to the creation of powerful new interfaces to digital document collections.

As a broader goal, this project seeks to bring to life the stored knowledge that lies behind the poor, text-based interfaces of textual document collections. Its aim is to develop new interfaces to these collections, in order to make them more accessible and useful as tools for knowledge discovery and transmission in education,

research, and other cultural contexts. This project introduces "deep interfaces" which combine text analysis with visualisation elements.

The practical outcomes of this project are web applications with interfaces focused on exploration-for-reading. Exploration-for-reading tools focus on the challenge of reading large-scale collections; of where to start with tens of thousands of pages. Deep interfaces help by revealing the content of the collection at the level of the interface. Many representations of digital collections focus on high-level overviews. In contrast, exploration for reading aims to help us encounter the collection through reading.

1.3. Practice and research

According to Candy (Candy et al., 2006), there are two kinds of research projects based on practice: practice-led research, and practice-based research. This project involves a combination of both types. The main focus of this research project is to advance knowledge about practice; it seeks to learn about practice through making artefacts. The process of making combined with reflection and iteration has led to new theoretical contributions. As an example of this, the notion of deep interfaces has arisen as a consequence of the practice. Thus by Candy's definition it can be considered practice-led research.

At the same time, in practice-based research the contribution to knowledge is demonstrated through the existence of created artefacts, and this is also the case of this project, which has produced and published several interfaces to significant text collections (see chapter 5 Results). These artefacts embody the transformed reading experience presented in this dissertation. They also show that the techniques are feasible, and how these approaches can be implemented in practice in a reusable and accessible form. All source code for the project is accessible through open repositories (Nualart, 2016).

1.4. Dissertation overview

This dissertation is presented in seven chapters, including this introduction. The chapters follow a traditional dissertation outline.

Chapter 2 "Context and background" reviews relevant literature and practice in three related fields. First, it addresses digital libraries (DL) in relation to text collections, showing how DL are reaching a post-human scale and, in contrast, the reviewed major public institutions are not using state-of-the-art tools and techniques in the interfaces offered to explore and discover the collections they host. Second, it reviews existing work on text analysis applied to text collections, showing that there are a variety of techniques —document clustering, topic modelling, topic labelling, and evaluation— available today, but no extended application of them in interfaces to text collections. Finally this chapter reviews interfaces to text collections. This third review revisits some paradigms in the interfaces to DL and the use of data visualisation techniques to represent collections, as well as a brief review of e-reading and interfaces oriented to read texts.

In Chapter 3 the gaps, opportunities, and lessons learned from the literature review are described and analysed. Leading on from these gaps, the project's key research questions and contributions to knowledge are outlined. Following the structure of the reviews in Chapter 2, the gaps found include a lack of advanced interfaces to text collections in the reviewed major institutions. In contrast, innovation is found outside the institutions, in practitioners and research demo sites. At the same time, there is an opportunity to use state-of-the-art techniques of analysis and apply them to text collections. Finally there are techniques applicable to deal with the problem that readers have no time to read the vast amounts of available texts and materials. The proposed solution is contained in the concept of crossreading, that is, that large collection of texts can be segmented into small pieces and then a sample of the pieces can be read giving, thus, a view, or at least a taste of the collection.

Chapter 4 describes the methodologies applied during the research project, including practice-led and practice-based research, reflexive and empiricist methodology. Secondly it describes the methods, tools, and techniques used in the development of artefacts. This includes techniques for: data gathering —APIs, spiders—

text segmentation, topic model analysis, topic model labelling and evaluation, and interface development with modern javascript libraries. The chapter closes presenting the methods used to evaluate the artefacts: online questionnaires and semi-structured interviews.

Chapter 5 presents the results of this research project, introducing each artefact in order of creation:

- Visference, an improved interface to the papers of an academic conference,
- Crossreads (I and II), a way to read text collections as fragmented narratives,
- Diggers Diaries (I and II), an interface to support reading of a historical collection of diaries from Australian soldiers in World War I.

Each artefact is described following a similar schema, that is: a figure showing a detail of the generated interface, the artefact name and description, the dataset(s), the collaborators, the data process, the data analysis, interface development, user evaluation studies, project outcomes, and a brief narrative about the artefact.

Chapter 6 presents the discussion. This chapter takes the significant gaps found in Chapter 4 —the text analysis, the reading task, and the interface— and, after analysing each of them in detail, defines the concept of deep interfaces. This concept refers to a group of techniques and strategies to improve the creation of new interfaces to text collections. Deep interfaces encompass techniques that help to answer the main research question: How might we create interfaces to text document collections that let us explore the collection by reading its contents? The idea behind deep interfaces refers to interfaces that look similar to standard web sites, but they offer options that take the user to an interpretative browsing of the contents, what is called the deeper meaning and structure of the collection. The so-called depth is a product of the integration of text analysis elements into the collection interface.

Chapter 7 recaps the whole dissertation, before discussing the limitations, applications and possible future work in relation to the proposed methods and the artefacts presented. In closing this chapter presents some final thoughts as the conclusions of the research project.

2. Context and background

Most of the data we generate everyday has an unstructured form, that is, mainly, text (Grimes, 2008, Heer, 2010). Part of this text is born digital (social networks, emails, etc.) and is stored in hard disks; part comes from printed paper and, after a digitization process, is stored, mainly, in digital libraries (DL), grouped into collections. As the number of text documents available digitally grows every day, the urgency for better information systems to access them grows too. In the age of data, every advance in a field of knowledge, sooner rather than later, will affect every other field, especially in an interdisciplinary domain such as the study of data. As this field grows in size —research centres, publications, industry, studies—it grows in diversity —biology, social networks, health, government, security, entertainment, etc.

This chapter reviews current research, and development of rich interfaces to DL, focusing on textual document collections. The related fields combined in this review are: data analysis, data visualisation and interface design. This project tries to work at the intersection of these three fields, and also tries to contribute to generating an interdisciplinary space where a combination of methods and points of view can generate good results.

This chapter starts with the review of DLs from important institutions that host large digital collections and a huge amount of digital objects. All of the reviewed cases are public institutions, except the Internet Archive, a San Francisco-based non-profit digital library. All the cases provide free public access to their collections through the web. The review finds that the big institutional DL studied are also meta-libraries, that is aggregators. This aspect contributes nowadays to the posthuman scale of digital libraries:. there are too many texts, to read even a small portion of them. In order to "read" such collections we need computational tools to guide and support our reading. As will be seen in this second chapter, these

DL are not using state-of-the-art interfaces to their massive amount of contents. The review shows that big DL such as Trove (Trove), Europeana (Europeana, 2008) are based on classic interfaces: text-based web pages. The review also shows some innovative proposals from institutions, demonstrating the gradual emergence of new features to help us reading, discovering and visiting digital libraries.

The second section reviews text analysis. It shows that DL catalogs based on standard metadata are unable to represent the complexity of these collections. It reveals that text analysis is important in a list of knowledge fields related to data: information retrieval, data mining, natural language processing, and computational linguistics. Therefore, text analysis seems necessary to represent the mentioned complexity. In this direction, established text analysis tasks are listed and described, such as: text categorisation, entity and concept extraction, taxonomy induction, sentiment analysis and document summarisation. The review then goes to text similarity techniques, such as clustering. Included is a brief history of clustering as a technique used since the 1990's to categorize or group documents within collections. It categorises three kinds of text collections according to the kind of cataloging: raw metadata, extended metadata, and object analysis. Finally, as an approximation to the rhizomatic complexity of the data, and especially texts, the concept of Capta is reviewed —defined by Johanna Drucker, as an active position to counter the non-critical attitude and position of acceptance we have to data. Closing this second block a review of topic modelling follows, including the labelling process and its evaluation.

The third section of the chapter presents a review of interfaces to digital libraries and collections of textual documents. It begins with a review of philosophical and conceptual paradigms from different authors for the last twenty years, including Sheiderman, Marchionini, Greene, Stamen, Dörk and Whitelaw. standard and trending practices in interfaces to digital collections are reviewed, showing that mainstream DL institutions,, tend to use classic interfaces, that is text based web pages. review of the role of data visualisation and interfaces to digital collections, shows how visualisation can be used as an element of the interface. Finally, closing the chapter, a brief review about interfaces dedicated to assist in reading text is presented. This last review shows that standards for reading texts on digital devices are established and offered for most devices and interfaces in contemporary

systems.

After this review, in the next chapter, the gaps and the opportunities found are presented, as well as the contributions to knowledge that the research offers.

2.1. Digital libraries and text collections

The term Digital Library can refer to a collection of digital objects, and, among other things, to the management of the collection, or to the institution or service offered by librarians. Here we use DL as a collection of digital objects. A DL can include objects such as video, images, text, and audio, all in an electronic format. DL also contains metadata, that is, data about the digital objects.

A digital collection can contain, according to Shneiderman, seven types of data, where text is a 1-dimensional data-type that includes "textual documents, program source code, and alphabetical lists of names which are all organized in a sequential manner" (Shneiderman, 1996). Textual documents can be kept together, forming a text document collection. As in other collections, textual documents can have diverse structures and contents. In this dissertation we refer to textual document collections that give access to the full text of the documents, so the texts can be read, analysed and exhibited.

When a set of textual documents that are part of one or more DL are grouped, they form a collection of textual documents. The cohesion factor that generates a collection can be broad, and, therefore, can affect the way the collection is presented and accessed. For example, the cohesion factor of a digital collection can be historical. This happens with, e.g., collections of press articles from a period of time and a defined location(Trove, 2010, library of congress, 1800). But a collection can be also a product of a curated work (American Anthropological Association, 2016). In both cases the nature and history of the collection can affect the design of the interface to access it.

Digital Libraries were born in the second half of the twentieth century, in parallel to the computer, and the digital era. The initial dream of a digital library as an integral source of knowledge, accessible from any place, is a reality today (Bush, 1989, Besser, 2004). One of the considered first DL is dated in 1976, the Oxford Text Archive (of Oxford, 1976). It contains literature and language resources, this

is, textual document collections. The term DL emerges in 1994 with the "Digital Libraries Initiative" (Fox, 1999). For some time the terms "virtual", and electronic are also used to refer to DL. Today "virtual libraries" usually refers to distributed libraries, known also as aggregators (Fox and Sornil, 2003).

As Borgman [1999] argues, "research and practice in DL has exploded worldwide in the 1990s". Borgman identifies several factors feeding the growth of DL, including the increased availability of networked computing and the availability of targeted research funding. Rehear (2004) explains that US, Europe, and Asia governments invested important amounts in DL research projects in the second half of the 1990s. Initially some humanists and librarians were sceptical about computer scientists researching DL, and not applying the results to real libraries. Hurtle in an editorial of D-Lib Magazine in 1999 says: "Rightly or wrongly, the DLI-1 grants were frequently criticised as exercises in pure research, with few practical applications" (Hirtle, 1999). Besser, about this period, says "we will call this the *experimental* stage of digital library development". In the period 1995-2000 international conferences, journals, and online news publications about DL were created in both scientific and humanities disciplines (Besser, 2004).

Today digital libraries can be found across a number of disciplines and domains. The scope and range of current DLs, includes:

- Public collections: institutional archives, public libraries, E.g., Library of Congress.,
- Private not-for-profit archives and collections (e.g. The Internet Archive)
- Commercial databases, online book shops (e.g. Google Books, Amazon, etc.)
- Public scholarly collections (e.g. university repositories)
- Commercial scholarly databases and journals (e.g. JSTOR, EBSCO, etc.)

A key point of DL is the access to their collections and digital objects, that is, the copyright of their contents. Sometimes, institutions can offer only access to metadata of specific digital objects due to license restrictions. Visualisation approaches that are reviewed in this chapter refer to freely accessible online DL. An example of this is the Million Book Project (Linke, 2003). This option looks like

one of the solutions to make the contents of DL accessible to the public, but there are still unsolved financial problems, since open access contents are expensive to maintain, and no clear business model is associated with them (Seadle and Arms, 2012). In any case, there are creative and revolutionary proposals in academia, like Shamos that argues: "copyright does not protect facts, information or processes, we propose to scan works digitally to extract their intellectual content, and then generate by machine synthetic works that capture this content, and then translate the generated works automatically into multiple languages and distribute them free of copyright restriction" (Shamos, 2005).

Examples of DL that offer open access to their contents include the Internet Archive, whose collections contain digital objects that are out of copyright and donated by public libraries, as well as works under copyright but published under license from the copyright owner. Archive.org acts as a publisher (archive.org). Europeana is a metasearch engine that harvests metadata from other institutions. Another case of, also called, aggregators (of metadata) is Trove (Trove, 2009), the biggest Australian source of digital collections . Trove accepts contributions from remote collections, but it is also a "growing full-text digital resource" where press articles, personal diaries and letters, books, journals, and biographies can be freely accessed. In the research world a notable open access repository of scientific papers, arXiv.org, was established in 1991 by Paul Ginsparg (arxiv.org). This is a repository for preprint versions of scientific papers, as Seadle says "in essence anybody can post an unreviewed paper claiming that it is original research" (Seadle and Arms, 2012). Arxiv is a modern concept for the distributed curation of digital library content by the community of its users.

The scale of these archives is notable, especially in the context of this project, which focuses on reading. For example Trove offers over two hundred million digitised newspapers. Arxiv.org offers over a million e-prints of scientific papers. The Open Library of Archive.org offers one million free ebook titles available to read. Modern digital libraries have reached what might be termed a "posthuman" scale, where they are too large to be read in the conventional sense.

Today some DL have become meta-libraries, also referred to as virtual libraries, and as aggregators. Aggregators collect metadata from source DLs and index it in new databases. This is the case of Europeana, an aggregator that collects

metadata from more than two thousand cultural institutions across Europe (Europeana, 2008). The Australian Trove "brings together content from libraries, museums, archives and other research organisations and gives you tools to explore and build" (Trove, 2009). Other popular aggregators include public libraries, and university repositories. These repositories are also sources collections for aggregators, allowing localization of titles through meta-search (OpenDOAR, 2005).

Since 1999, digital libraries use systems to catalog their contents based on metadata (Milstead and Feldman, 1999). Metadata practices, as a descendent of the cataloging, grew out of traditional library management —indexes, card catalogs etc. Metadata repositories can store data about physical and/or digital objects. Best practices in repositories recommend the use of standards for metadata vocabulary terms (Duval et al., 2002). A widely adopted set of standards in this direction is Dublin Core Metadata Initiative (DCMI, 2001), that defines a vocabulary of terms to be used for physical and digital resources. Standardization of metadata allows easy and more effective sharing of resources among machines.

With the evolution of online resources, the relation between data and metadata becomes diffuse, thus, e.g., Munzner, in her recent manual "Visualization analysis and design" (Munzner, 2015) argues that "the line between data and metadata is not clear, especially given that the original data is often derived and transformed". Therefore Munzner, decides to not "distinguish between them, and refer to everything as data". In this text we differentiate between data and metadata in the sense that data is unstructured (since we study textual documents) and metadata has the structure of a list of property-value pairs (e.g. in JSON {"Title": "Tree of Science", "Author": "Ramon Llull", "year": 1596 } etc).

Focusing on textual collections coming from paper, and the access to their contents, there is one aspect that differentiates textual documents from other media, that is the digitization process (Seadle and Arms, 2012). The difficulty in having full text access to the documents of collections is that, once the documents are digitised, a transcription is needed, and this task requires human resources. OCR systems can help in the process of transcription but success depends, mainly, on the quality of the copies. Sometimes the digitised collections are published online, as in the collections of the University of Washington Libraries (U.Washington), and the State Library of New South Wales (SLNSW, 2010). The digitised texts

are not always available as machine readable text. When the contents are readable by machines, then we can get all the advantages of information retrieval, and text analysis. Several institutions have developed strategies to make them available as text: Trove (Trove, 2010) uses "crowdsourced" corrections to an automated OCR transcription. The Bentham papers archive uses crowdsourced transcription but is now also developing a machine learning system for transcribing handwritten text (see e.g. Causer and Wallace, 2012, UCL, 2000, TranScriptorium, 2013) Finally the State Library of New South Wales (SLNSW) funds manual transcription of documents such as the WWI diaries collection (SLNSW, 2014)

When the full text of collections is available, then we have metadata and data, all as text. Additionally, the textual documents contain some kind of structure: paragraphs, sentences, and/or chapters, sections. This extra data can be extracted and stored as metadata too, thus, enriching the structure and adding extra dimensions to the collection.

The "library" model for digital collections of documents brings with it specific conventions. The traditional library stores books, and makes them findable (via index or catalog). One key function for the librarian is to "look inside" the collection, advising users on its content. With digital textual documents the collection contents can become transparent. At the same time new challenges of scale arise. Thanks to computers it is possible to efficiently analyse large amounts of text; this analysis can bring structure to collections of documents, and a deeper knowledge of the collection contents. This new information can help with generating overviews, relationships, richer models of content, and, eventually, improved understandings of a collection and its contents. The next section presents a review of text analysis and text document collections.

This section has reviewed big and public institutional DL. Most of them show the tendency to become big aggregators, this is, concentrations of information (metadata) to access to local and remote resources (digital collections, and digital objects). This aspect makes the DL grow even faster, to what we call a posthuman scale. This scale brings urgency for new tools to visualize, explore, and to read the text documents of the DL. The review shows the contrast between this posthuman scale, and the lack of active innovation in the online versions of the institutional digital collections

2.2. Text analysis and text collections

This section introduces the concept of text analysis and its role and use in the development of DL and text collections. As shown in previous section, the number of available DL and text collections grows everyday. The urgency for tools that help us to deal with it also grows. Text analysis seems to be a *sine qua non* as the main ingredient of these tools, therefore this section reviews the use and opportunities in text analysis for DL and text collections. The review starts with an introduction where text analysis is defined and established methods are listed and described. The following section reviews the role of text analysis and metadata in cataloging digital objects. Five major digital libraries are compared in their use of methods to catalog collections. Finally this section discusses the concept of data and its complexity and subjectivity, in contrast with the simplification that often arises when using standard metadata properties to represent the rich contents that digital collections contain.

When first encountering a specific collection of textual documents, before any analysis takes place, the kind of metadata properties expected for a textual document could be, for example: title, author, dates, locations, etc. These fields are standard metadata fields and are the primary source for cataloging collections. Digital libraries can readily increase the number of indexed items in their databases due to the easy and inclusive metadata formats used. This is the case in all five DL reviewed (see table 1). The use of standard fields of metadata brings homogenisation, therefore it is an approximation to the abundance and complexity of the collection. Metadata inevitably, brings a simplification or summary of the content. Standard metadata simplifications were primarily designed to support search and retrieval. An example of a collection of texts represented almost exclusively by metadata is “Mapping the republic of letters”, a meta-project to visualise letters from 1600 to 1800 by recognized intellectuals (Findlen et al., 2011). The visualisations listed in the project represent the metadata of the letters, such as: dates, locations, senders, recipients. However in most cases there is no access to the letters themselves.

Beyond the metadata there is the analysis of the documents of the collection with computer-aided (procedural) text analysis. Procedural text analysis is a key

point in the foundation of the fields of information retrieval (Liu, 2009, Salton and McGill, 1986), data mining (Han et al., 2011, Kantardzic, 2011), natural language processing (Manning and Schütze, 1999) and computational linguistics (Grishman, 1986, Hausser and Hausser, 1999). Today, most of the methods used in those fields use techniques grouped under the label of machine learning (ML).

Established and developed text analysis tasks include:

- Text categorisation: this is the technique of assigning documents of a collection to two or more predefined categories. A typical example of text categorisation is to classify emails as spam or not-spam. This is a typical ML task. An example and study of a visualization of document collection classification is Di Munzio's work (Di Nunzio, 2006)
- Text clustering: a generic name for a variety of techniques that deal with the task of grouping documents according to their similarity. See the following discussion for examples of text clustering.
- Entity and concept extraction: as part of information extraction, it is a term that includes the task of identifying and interpreting meaningful parts of a text. These parts can be names —entity recognition extraction— or concepts defined in a customised or pre-existing thesaurus (Tseng, 2002).
- Taxonomy induction: a techniques for organising terms of a document in a hierarchical way using external lexical resources (Fountain and Lapata, 2012). The most used resource is WordNet (Fellbaum, 1998). Other resources have been used such as Wikipedia (Ponzetto and Strube, 2011)
- Sentiment analysis: is the analysis of the attitude of the author of a message in relation to one or more topics. It is also known as opinion mining (Leetaru et al., 2013, Pak and Paroubek, 2010).
- Document summarisation: the process of shortening a text while showing the main points of the whole text (Jones et al., 2002, Gong and Liu, 2001).

To realise the huge evolution of the field, notice that in 1999 text analysis was limited to text categorisation and clustering, information extraction, and summarisation. Clustering of documents was used for document visualisation, and

categorisation(Tan et al., 1999). In this dissertation, we focus on the most popular text analysis techniques applied to visualisation of collections of texts, which are text clustering techniques. Such techniques can produce multidimensional properties of the objects of the collections, and this generates a richer classification, and therefore, a wider possible range of interpretations of the collection.

Text analysis is a broad concept that includes multiple techniques, methods and fields. In order to decide which aspects of the text analysis in relation to collections of texts to review, three dependent levels of text analysis in collections can be defined:

- Level I: Raw metadata. Text collection without text analysis, that is, using only standard metadata properties.
- Level II: Extended metadata. Text collection with extended metadata. This process enriches the initial metadata generating derived properties, e.g. from the fields "age", to generate "groups of age". This case can include some analysis of the object, e.g. counting the length of the text of each document.
- Level III: Object analysis. Text collection with text analysis of the textual documents of the collection. This process generates new properties, e.g. a list of topics for each document, a classification of documents according to contents, authorship analysis, recommendation system according to text contents or style, etc.

In the cases studied, and in general, in institutional DL there is a lack of application of techniques from the trending field of big data and data analysis, and in particular text analysis. To find projects that use text analysis as a technique that can help to understand and get to know a collection, and to build tools to explore and discover the collection, we need to go to independent data visualization practitioners and researchers. Section 2.3.3. Reviews data visualization and interfaces to text collections, and includes a list of cases that meet the criteria for level III.

CHAPTER 2. CONTEXT AND BACKGROUND

Reviewed site	Level	Comments	URL
Europeana	Aggregator: level depend on the source	An aggregator that indexes millions of online resources.	-
State Library of New South Wales	I	Interface for reading It offers crowd-transcription of documents	http://transcripts.sl.nsw.gov.au/project/World%20War%201%20Diaries
Internet Archive	I and II	Complete interface for reading	https://archive.org/details/MBLWHOI
Trove	I and II	Indexes millions of online and offline resources. newspapers collection is full text searchable. It offers crowdsourced transcription of documents.	http://trove.nla.gov.au/newspaper/article/13279967?searchTerm=&searchLimits=1-australian=y
Library of Congress (USA)	I and II + transcriptions	Newspapers collection is fulltext searchable. Complete interface for reading	http://chroniclingamerica.loc.gov/lccn/sn83045389/1916-04-03/ed-1/seq-1/
Non institutional projects	I, II and III	There is a gap on the use of Text Analysis techniques in institutional digital collections.	See 2.3.3 Data visualization and interfaces

Table 1. Comparison of five major digital libraries according to how they use text analysis. None of compared cases meet the criteria for level III.

Traditional natural sciences define data as a representation of reality with a quantified, and usually controlled error of measurement. The results, with natural sciences methods, are expected to be objective. The concept of data, however, is less deterministic from a humanities point of view. In this direction, Drucker (2011) differentiates between data and capta. Since data, from the Latin "given", assumes a passive position of the observer in front of the real phenomena, capta, from the Latin "taken", puts the observer in an active position, being able to interpret the measurement and, therefore, critically extract conclusions. Drucker says: "From this distinction, a world of differences arises. Humanistic inquiry acknowledges the situated, partial, and constitutive character of knowledge production, the recognition that knowledge is constructed, taken, not simply given as a natu-

ral representation of preexisting fact". Text analysis treats the document as data and attempts to use computational techniques to reveal features of the texts. The complexity of the data, its subjectivity, instead of being an obstacle in this research project, is a motivation for the review of cases that use procedural methods, as well as curated ones for the text analysis of the collection. Rather than treat the results of text analysis as "data" that objectively represent the analysed text, Drucker's concept allows us to frame them as "capta" for critical human interpretation.

2.2.1. Text clustering and interfaces

Text clustering deals with grouping a set of documents according to their similarity. Clustering is a text analysis techniques that groups documents according to some defined properties. Representations of collections of texts with clusters are very common, in 2D (Weskamp, 2004, Paulovich and Minghim, 2008, Masad and Nayar, 1011, Lagus et al., 2004), and in 3D as well (Chalmers and Chitson, 1992, Wise et al., 1995, Carter and Capretz, 2003, Andrews et al., 2002)..

The most popular technique to measure the similarity among documents is cosine similarity. Every textual document belonging to a collection of textual documents can be represented by a vector using methods based on term frequency-inverse document frequency (TF-IDF). TF-IDF is a statistical value intended to score the importance of a word to a single document in a collection. TF-IDF provides a quantitative measure of the occurrence of specific words (relative to the document as a whole). By comparing significant words from different documents we can generate a measure of document similarity (Lee et al., 2005). This analysis generates a network of documents. The cosine of the angle between two document-vectors is a measure of the similarity of the two documents.

The similarity among the documents of the collection generates a network of relationships. Visualisations of these generated networks include Andrews' work InfoSky (Andrews et al., 2002), which proposes a space travel metaphor where planets are documents and proximity among them is proportional to similarity. Today the use of TF-IDF techniques is related to big data processes (Leskovec et al., 2014). Another example of similarity broadly understood, is the networks of chapters in Grimm's Fairy Tale Network by Jeff Clark (2009), where similarities

are calculated by comparing the vocabulary used in each tale of the collection. Similarity among documents can sometimes be seen in comparing representations of the documents. This is the case of Word Storm by Castella and Sutton (Castella and Sutton, 2014), an interesting variation on word clouds: each paper from a conference is represented by a word cloud. The algorithm that places most frequent words in a 2D area has been modified in order to place recurring words in the same location, making it easier to compare word clouds. Methods to calculate document similarity are used in information retrieval systems for the task of term weighting. Term weighting, defined by Zhang et al, "is the job to assign the weight for each term, which measures the importance of a term in a document" (Zhang et al., 2011). Examples of the visualisation of search result document similarity through 2D network diagrams include Kartoo (Baleydier, 2001), and Touchgraph (TouchGraph, 2001). Since Kartoo and Touchgraph were commercial projects, the techniques used to measure similarity are not published.

2.2.2. Topic models

Topic modelling is a group of statistical methods that analyse collections of textual documents extracting the thematic composition of each piece of the corpora. A topic model is represented by a set of words that have high probability of appearing together in a document within a corpus of textual documents. These words, usually called terms, "look like topics because terms that frequently occur together tend to be about the same subject" (Blei, 2012). The number of topics to generate is usually set as a parameter to the analysis process. Each document gets a normalised score for every topic. Topic model algorithms are not based in semantic analysis, but purely statistical co-occurrence. In machine learning terms: "topic models offer an unsupervised, data-driven means of capturing the themes discussed within document collections" (Aletras et al., 2014).

Topic models were introduced by Blei, Ng, and Jordan in 2003 as a "generative probabilistic model for collections of discrete data such as text corpora" (Blei et al., 2003). LDA developed from previous works including the early work on Latent Semantic Indexing (LSI) (Deerwester et al., 1990) and probabilistic Latent Semantic Indexing (pLSI) approach by Hofmann (1999).

Topic models can be used in many ways, e.g. documents can be grouped based on their combination of topics. Part of the collection can be filtered and selected for further reading according to specific topics. One example of grouping documents according to their similar topics is the Stanford Dissertation Browser (Ramage and Chuang, 2012) where an interactive visual browser allows the exploration of theses, showing similarities along different domains of knowledge.

The relationship between topic models and digital humanities (DH) is an example of interdisciplinary collaboration between two relatively young research communities. An influencer of this collaboration is Franco Moretti and his distant reading tools. As a general strategy , Moretti recommends that literature researchers read less and, instead, represent more (Moretti, 2005), Thus, the texts can be studied from overview representations, timelines, graphic comparisons, and story diagrams, among others. Topic model methods fit very well into Moretti’s idea, and have been used in the humanities largely as we review in the following paragraphs.

Literature and scientific texts have been a perfect terrain for experimentation with topic models in recent years. Currently there are tools (like Mallet) easy to use for non-experts in math and statistics, to use especially to welcome DH researchers. Critiques and opinions from DH arise too about the complexity, the interpretation, and in some cases the value of topic models. For example Lisa M. Rhody analyses the behavior of topic model analysis in poetry, and finds that the reduction of the complexity of figurative texts (such as poetry) that topic model analyses produce demonstrates the very complexity of the source texts. She suggest that ”topic modelling poetry works, in part, because of its failures”, and talks about ”an interpretive space (...) between the literary possibility held in a corpus of thousands of English-language poems and the computational rigour of Latent Dirichlet Allocation (LDA)” (Rhody, 2012). Schmidt in his ”Words Alone: Dismantling Topic Models in the Humanities” in 2012 was quite sceptical in the use of topic models by humanists, and argued: ”simplifying topic models for humanists who will not (and should not) study the underlying algorithms creates an enormous potential for groundless — or even misleading — insights.” (Schmidt, 2012). Goldstone and Underwood have, independently, applied topic models to the study of the history of literary studies. They used different software, stop-word lists, and

numbers of topics. The results have overlapped and diverged in different places, demonstrating the non-universality of the methods. Despite results that denote again the complexity of the interpretation of topic model analyses, they "reached a shared sense that topic modelling can enrich the history of literary scholarship by revealing trends that are presently invisible" (Goldstone and Underwood, 2012).

Today, topic model applications to DH, and to text collections is supported and recommended as a powerful tool to analyse large collections of documents. Recent works, like Jockers' Macroanalysis: Digital methods and literary history (Jockers, 2013) are optimistic about the use of computing analysis, and advocate the revolutionary potential of large-scale literary analysis. Two cases that use topic model analysis of a text collection are discussed in 2.2.3: "Mining the dispatch", by Nelson (Nelson, 2010a), and "Topic Modeling Martha Ballard's Diary" by Cameron Blevins (Blevins, 2010)

One more remarkable case of topic model exploration and management is MetaToMAT (Metadata and Topic Model Analysis Toolkit), by Snyder et al. In their words this approach is "a visualization tool that combines both metadata and topic models in a single faceted browsing paradigm for exploration and analysis of document collections" (Snyder et al., 2013).

The efforts to make topic modelling available to a community of non-experts is a key point in the growing use of these techniques by the humanities community. There are multiple free and commercial software tools that calculate, represent, and manage topic models. To cite the most popular ones: Mallet (McCallum, 2002), LDAvis (Sievert, 2015), dfr-browser (Goldst, 2014), Termite (Chuang et al., 2012), Serendipity (ale). Most of these tools also offer representation of the topic models. They represent the relative scores of each topic, their distribution in time, the comparison of terms in topics.

2.2.3. Topic model labelling

Topic model labelling is a method that is used to complement topic model analyses, making topics human readable. As mentioned, a topic model is an abstract statistical construct that may or may not be equivalent to a "topic" or a theme for a human reader. Therefore, in order to be interpreted and used by humans,

topic models need effective representations, adapted to the nature of the collection and oriented to specific tasks. One way to represent a topic model is through the number of terms that co-occur frequently along the documents of the collection. As shown in Figure 1 the list of topics that a topic model software such as Mallet outputs by default, have no label. By default we obtain a list of topics that we can number as topic1, topic2, topic3, ... Efforts in improving this basic representation of topic models have been published, indeed this is an active research question.

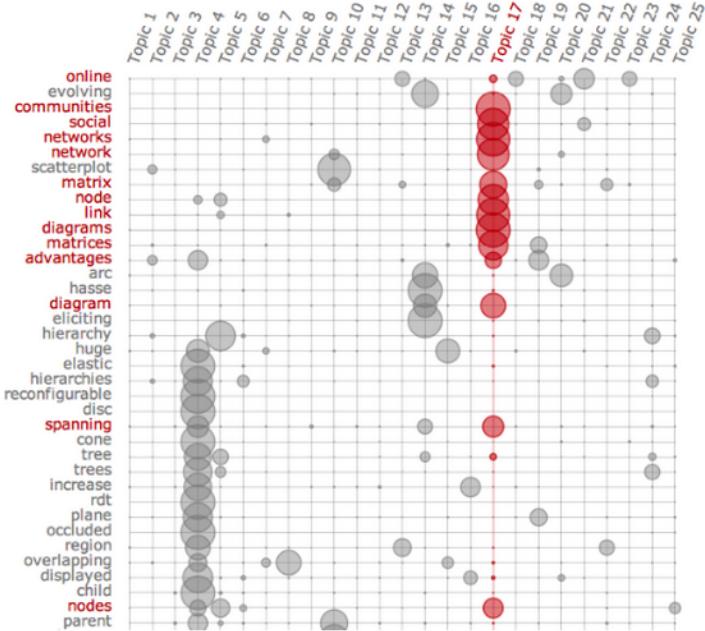


Figure 1. Topic visualisation with “Termite”. By default the the generated topics are named: topic1, topic2, topic3,...

Aletras et al (2014), show that text based labels are easily interpreted by humans in evaluation studies. Text labels are more effective than image labels. A text label can consist of a list of terms, a single top term, a phrase, or a sentence (Mei et al., 2007).

A multitude of methods for text labelling of topic models are published, including procedural and curated methods. Procedural methods seek to automatically assign labels to topics. This is used in cases where the number of topics generated is large, and cases where the whole process is automated. Some methods recombine the

n-terms and then reapply the model, getting a ranked list of terms within a topic; then the highest ranked term can be used as label for the respective topic (Lau et al., 2010). Other methods involve external data sources, this is the case of Lau et al (Lau et al., 2011) that obtain text labels from the n-terms, from titles of Wikipedia articles containing the terms, and from sub-phrases extracted from the Wikipedia article titles. Since topics can overlap each other, some other works produce labels from topic hierarchies based on parent-child relationship between topics (Mao et al., 2012). The mathematical review of these procedural methods is beyond the scope of this dissertation.

Curated methods of labelling topics, that is human generated topics, are common in DH, where experts can label according to a defined goal (Schmidt, 2012). Some authors from science present the subjectivity of curated labelling as a problem that "can easily be biased towards the user's personal opinions" (Mei et al., 2007). Meeks (2011) comes to the conclusion that small corpora of texts produce broad topics. Since he is interested in topic networks, a bigger corpora creates networks too complex to be visualized in 2D. He finally assigned topic labels asking a group of experts for "a comparison of labelling of topics (...) based on their interpretation of the topic's connection to various words, as well as to various papers".

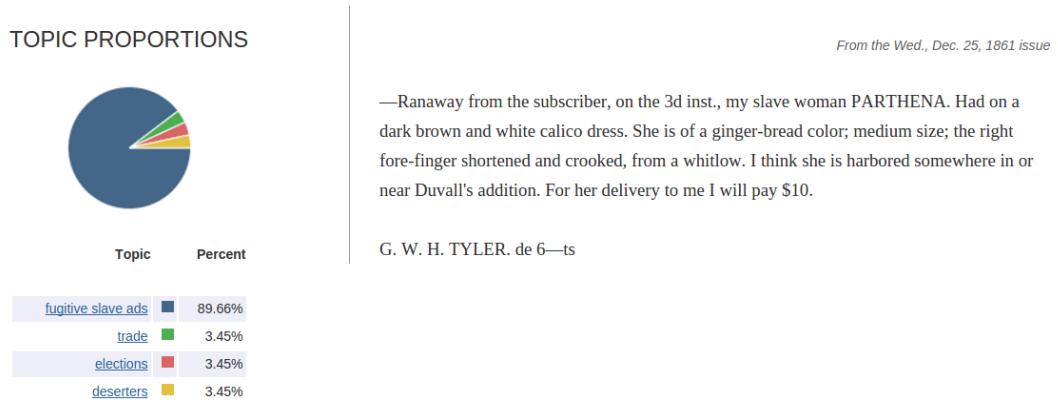


Figure 2. "Mining the dispatch" by Robert K. Nelson (2010). Detail of one news text of the collection and its topics pie chart.

There are two DH projects that inspired our research: "Mining the dispatch" and "Topic Modeling Martha Ballard Diary". "Mining the dispatch", by Nelson (Nelson,

2010a) is a topic model analysis of a collection of news texts from the U.S. Civil War during the years 1860-65 in Richmond, the capital of the Confederacy's newspaper of record (see Figure 2). This project shows line charts graphing the change in topics along time. Nelson explains how the labelling of topics was done, and his impressions about the results: "I have given each topic a label based upon my reading of pieces drawn from that category; these labels are informed judgement calls and are imperfect". For each topic, the interface offers a numbered list of excerpts of news texts ordered by relevance. One more click on each excerpt shows the full text for the article.

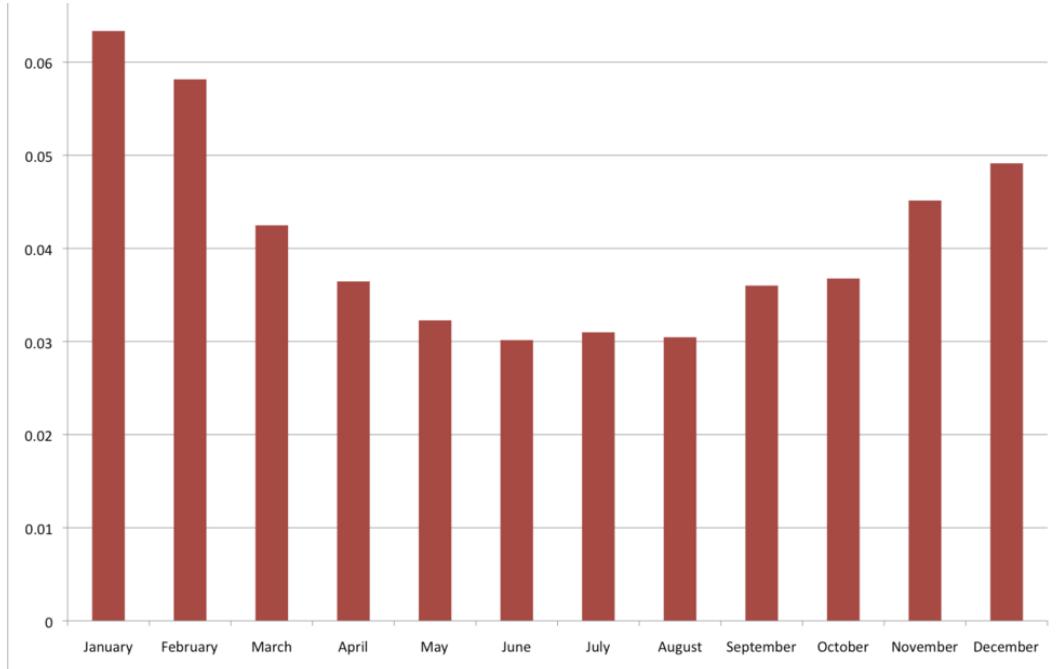


Figure 3. Scores for topic “cold weather” grouped by month, from “Topic Modeling Martha Ballard’s Diary”, by Cameron Blevins (2010).

”Topic Modeling Martha Ballard’s Diary”, by Cameron Blevins (Blevins, 2010) analyses the diaries of Martha Ballard, a New England midwife who kept a daily diary for over twenty-seven years from 1785. The collection consists of one thousand four hundred handwritten pages. The topics generated are manually labelled with “descriptive titles”, such as: midwifery, church, death, gardening, shopping,

illness, cold weather, etc. This project outputs charts with topic evolution, and provides striking evidence of the effectiveness of topic modeling in representing document content (see Figures 3, and 4).

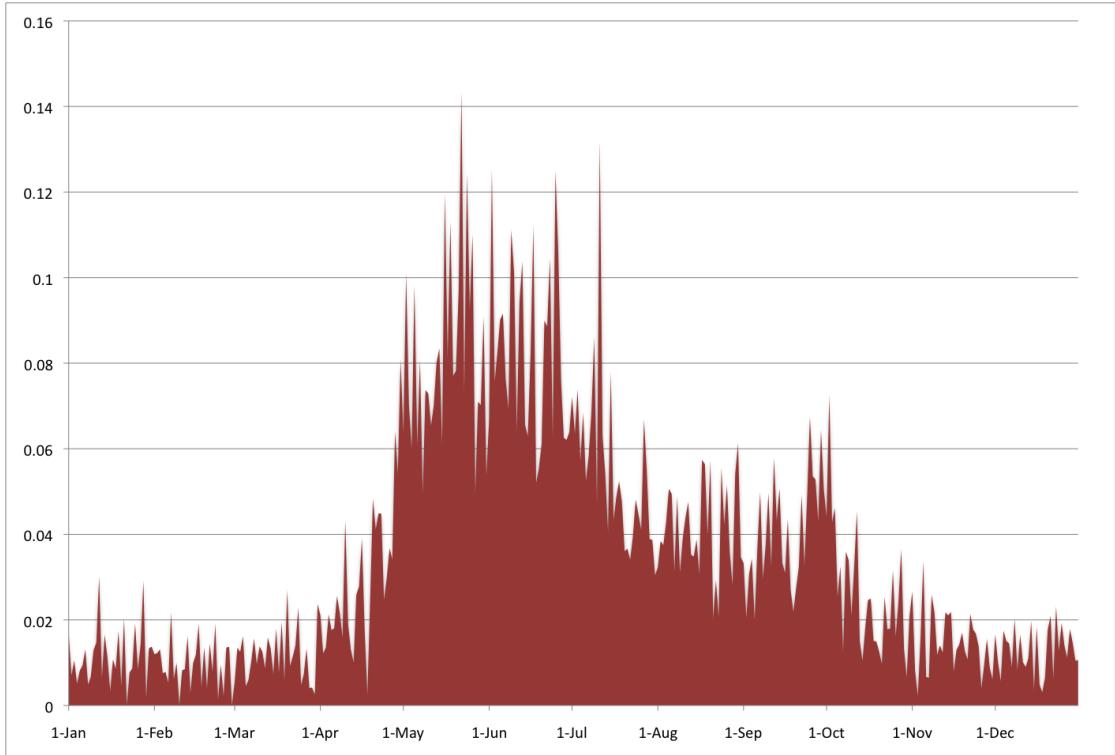


Figure 4. Scores for Topic gardening by day, from "Topic Modeling Martha Ballard's Diary", by Cameron Blevins (2010).

2.2.4. Topic model evaluation

Topic model interpretation is a not well-defined task. According to Schmidt (Schmidt, 2012), topic models are "like words, they are messy, ambiguous, and elusive". Therefore it is difficult to establish methods to evaluate and compare the quality of topic model analysis. A number of questions arise initially: How good are the results of topic model analysis? How might we compare different analyses? How might we measure their consistency from a human point of view? All these questions guide research into topic model evaluation.

In DH we find that evaluation of the results of topic model analysis is undertaken

by experts in the contents of the analysed documents. This is the case of Goldstone and Underwood, who compare their independent analysis of the same corpora and validate, as experts, which topic models diverge and converge in a more appropriate way (Goldstone and Underwood, 2012). There are also cases where evaluation consists of finding evidence in the corpus; this is the case of Ballard’s diary, where entries related to cold weather are higher in winter, and, complementaryly, entries in the diary talking about gardening are higher in summer (see Figures 3 and 4). These charts validate, the use of Mallet for topic model analysis, as the author says (Blevins, 2010).

Computer generated topic model labels can be evaluated by humans. This is the case of Newman et al, who in 2010 that published a scoring model that predicts human scores. The experiment produces a ranked list of terms of a topic, where the best word is used as a label for representing the topic. This list is evaluated by humans by comparison, that is, the humans do a ranking and then the rankings (curated and procedural) are compared. It is based on point-wise mutual information (PMI) of word-pairs —PMI measures the association between two words according to large text corpora. The model uses external sources such as Wikipedia, Google n-gram data set and Medline to adjust the initial calculated ranking of terms. (Newman et al., 2010).

Visualisation tools are used to overview and evaluate topic models, usually by visual comparison, like Termite (Chuang et al., 2012), or with a dashboard that allows multi-comparison and editing of the topics on-the-fly, like Serendipity (ale). Despite these and other works, the subjectivity of topic models when evaluated by humans makes it difficult to establish reliable measures of topic model quality.

This section has shown that while text collections are increasingly available as full text documents, current DLs use only basic metadata to represent and provide access to their collections. The text analysis techniques reviewed here can be applied to characterise large collections. The main text analysis methods with examples are mentioned, specially text clustering, as a general technique to group similar documents and, deeply, topic modelling, as the prominent modern technique for text analysis of large collections of textual document in the humanities. The success of topic models is linked to the process of labelling the topics, thus making them human-readable. This importance has encouraged innovation in eval-

uation of topic model labelling methods. In some reviewed experiments procedural labelling can be comparable to curated labelling, though in DH the reviewed cases are curated by experts. In the application of topic models we found a gap: most of the reviewed cases use the output of the analysis to practice “distant reading”, that is, to represent the collection as a whole. The use of topic models in textual document collection interfaces to help exploration and reading is not common in the reviewed cases.

The next section presents a review of interfaces to text collections. It includes a review of paradigms in interface design strategies, a review of standards and trends in in-production interfaces to text collections, a review of the role of data visualization in interfaces to text collections, and a brief review of interfaces where the main aim is to assist in reading textual documents.

2.3. Interfaces and text collections

The development of systems to store, catalog, and interact with digital collections has outstripped the development of user interface. In recent decades the interfaces to text collections in the domains of cultural heritage and scientific publication have not developed much, especially when compared to commercial interfaces, such as online shopping sites, social network sites, and even online banking (Mario Perez-Montoro and Jaume Nualart, 2015). A broad definition of interfaces to text collections would include any software and hardware that connects humans and digital data. For the sake of this research project we refer to interfaces to text collections as software mostly developed as web applications, and accessible online. Since this is a very big field, this section presents reviews of four interrelated topics: interface paradigms, standard and trending practices in text collection interfaces, data visualization, and interfaces for reading.

Firstly several interface paradigms or metaphors that founded strategies for innovation in interface design for the last twenty years are reviewed. This review will help the reader understand the approach and innovation of the artefacts. Secondly standard practices in contemporary interfaces to text collections are reviewed. Several examples of image document collections are included, as inspiration and a source of applicable ideas for text collection interfaces. As a third step in this

multi-topic review, we consider data visualisation approaches used as interfaces to text collections. This review of data visualization helps to define its role in modern interface design. Finally, due to the importance of textual documents in this dissertation the final review examines what we call “interfaces for reading” that bring textual collections closer to potential readers.

2.3.1. Interface paradigms

This section reviews general strategies and concepts to advance the design of interfaces to text collections, and digital libraries. This brief review includes selected works that are relevant to the topic, and have influenced this research project, especially during the conceptual development of the artefacts.

Probably, the most influential work in interface development and evolution for the last twenty years is ”The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations” by Ben Shneiderman —which has more than 3,500 citations in the literature (Shneiderman, 1996), and his Visual Information Seeking Mantra: “Overview first, zoom and filter, then details-on-demand”. In Scheiderman’s words: ”To sort out the prototypes and guide researchers to new opportunities”, he proposed ”a type by task taxonomy of information visualizations”, for collections of items, with multiple attributes each. These tasks are: overview the collection, zoom it, filter it, get an item’s details-on-demand –these three are included in the mantra— as well as get related items, access to the exploration history, and allow extraction of subsets of the collections. The mantra was presented as a design principle based on Schneiderman’s experience in ”information collections” and libraries.

In this direction, in 1998 Marchionini et al (Marchionini et al., 1998) analyse early interface developments in a collaboration between University of Maryland designers and staff from the Library of Congress. The goals were “that users should maximise their interactions with information resources and minimise their attention to the system itself”. Three kinds of tools were developed: “Overviews of collections, object previewers, and object gatherers”. The interface design is based on simplicity and clearness: “A standard toolbar on the left of the screen lists the functions available”. The justification for the necessity of interfaces is

defined in relation to the size of the collections: “very large amounts of materials in many different formats with varying levels of descriptive metadata makes searching difficult and browsing more important”. In these early years some structural elements were introduced: hierarchical table of contents in sidebars, navigation bars, quick and advanced search forms.

In 2000, Greene, in collaboration with Schneiderman and others, (Greene et al., 2000) establish a methodology to apply the idea behind the Visual Information Seeking Mantra. With the aim of aiding ”designers of digital library interfaces” they ”present a framework for the design of information representations in terms of previews and overviews”, as well-delimited elements to be used as pieces and features of the interface. Previews are defined as an analogy to bibliographic records, they ”acts as a surrogate for, a single object of interest”. A bibliographic record suggests text —but a preview can be a thumbnail image, or other representation of the object: ”An effective preview is an information surrogate that communicates to the user, at the appropriate time, sufficient information about the primary object it represents to support users in making a correct judgement about the relevance of that object to the user’s information need”. Overviews, in contrast, as an analogy to catalogs. In that sense Greene et al, supported by Marchionini’s idea that information seekers ”engage to change their knowledge state” until they have satisfied their information needs, present previews and overviews as representations that support browsing and scanning for exploration. ”Previews and overviews” from this point of view, can be seen as a formalization, and an elaborated development of Schneiderman’s previous works, such as the Information Seeking Mantra.

In 2009 a new paradigm was introduced by Stamen Studio, an influential team of practitioners that argue that an overview could include all the items of the collection. This idea is known by the term ”show everything”. The items are not hidden any more, the visitor gets an idea of the whole collection as opposed to the tiny search box that initially hides the collection. Despite the assertion of ”show everything” when the collection is too big to fit technically and/or perceptually on the screen, then alternatives are required, such as zoom, and filters, which introduce intermediate states in between Greene’s previews and overviews. Stamen’s slogan is seen as a combination of pragmatism (Whitelaw, 2015b) and provocation (Whitelaw, 2015a). Beyond what Shneiderman and Greene provide, ”show

“everything” is significant in that it introduces this idea into a contemporary web context and begins to show how rich collection interfaces can work in the modern web browser.

The Shneiderman-Stamen combination could be seen metaphorically as a walker who arrives to the top of a hill (home page), and from there can see a valley (overview), and decide where to go and which paths to take (detail).

In 2011, Dörk et al introduce the paradigm “Information Flaneur” (Doerk et al., 2011). They presents “explorability as a new guiding principle for design and raise research challenges regarding the representation of information abstractions and details”. The information flaneur proposes to break the rigid axes of information seeking systems based on the concept of offices and corporative buildings by transcending the traditional hierarchy of files and folders to a new perspective using the city as the metaphor for the information seeker. The seeker —the flaneur— walks the city following, and motivated by curiosity, emotions, and personal experience. This model reframes the user, rather than the user as driven by task and function, here the user can be motivated differently. The Information Flaneur emphasises individual interpretation and subjective experience, and it also changes the criteria for “success” in interface design, recognising that there is more than “task” and “information retrieval”. This idea puts complexity, individual interpretation, and subjectivity in the centre of the design.

In this direction, arises the concept of Generous Interfaces, by Mitchell Whitelaw (Whitelaw, 2015), referring to interfaces that show the abundance of digital collections, in contrast with the classical interface that “is ungenerous” and “withholds information, and demands a query”. A GI would emphasise browsing and visual exploration, and would support experimentation with visual elements. Most of the cases presented in Whitelaw’s paper are applied to digital collections of images. Nevertheless, the introduced ideas and artefacts in it are applicable to any collection of digital objects including text collections.

This review shows the conceptual development of trends and strategies in interface design. The list of authors reviewed seems to follow the same path from Schneiderman’s “Mantra” to Greene’s “Overview and previews”. These efforts shed light into the beginning of the digital era and the evolution of information displays, twenty years ago. These initial proposals function as guidelines for researchers and

designers of interfaces to DL. The paradigms of Stamen, Dörk and Whitelaw incorporate subjectivity and challenges, and are less explicit, but more flexible in their proposals. Far from eluding complexity, this brief journey through paradigms and metaphors —Schneiderman, Greene, Stamen, Dörk, Whitelaw— take us to a point that I develop in chapter 6 Discussion, where too big, too complex, or too undefined does not mean necessarily a complex interface, and, therefore problems in explaining how to use a new feature, or, ultimately, a whole new interaction concept. Simplicity in the interface does not require simplification of collection contents, but a simplification of visual language and interaction features. The complexity and diversity of cultural collections, as defined in generous interfaces, is what should be revealed through interfaces.

2.3.2. Standard and trending practice in text collection interfaces

We introduce and comment on a list of what can be considered standard practices in DL and, particularly, in text collection interfaces. Interfaces to digital collections have developed in parallel to digital libraries, that is since around 1990, therefore to review interfaces to textual document collections it is necessary to review interfaces to DL in general, because we will see that in classic (text-based) interfaces there are no significant differences between DL and text collections.

The DL reviewed here from institutions, archives, and libraries have, mostly, the form of a text-based interface. Text-based means that the interfaces to access the collections are standard text pages, where the two features always present (Tedd and Large, 2004) in the interface are the navigation and the search systems (Mario Perez-Montoro and Jaume Nualart, 2015). Navigation is a simple tree of categories and subcategories, and the search system is a traditional full text query with the option of advanced search that includes a set of filters. The results, in both cases, navigation and search systems, are flat lists of items with a short description of each item, and sorted by categorical criteria (that is, by date, by name, etc.), or search system ranking, usually not transparent to the user (Mario Perez-Montoro and Jaume Nualart, 2015).

Archive.org, one of the referents for online non-profit DL, (archive.org) has in-

troduced faceted lists of collections (for textual documents and other multimedia collections alike), and modern graphic design prior to search. The faceted lists 5) are rudimentary: sizes depend on to the length of the titles, but not on the size of each collection, and section of collection. Only a basic by-type (video, image, text, audio) filter is generalised. No overviews are offered, neither metadata visualisations, nor visual analysis. A faceted view is a kind of "poor" overview which reveals the contents of the collection but in a very limited way.

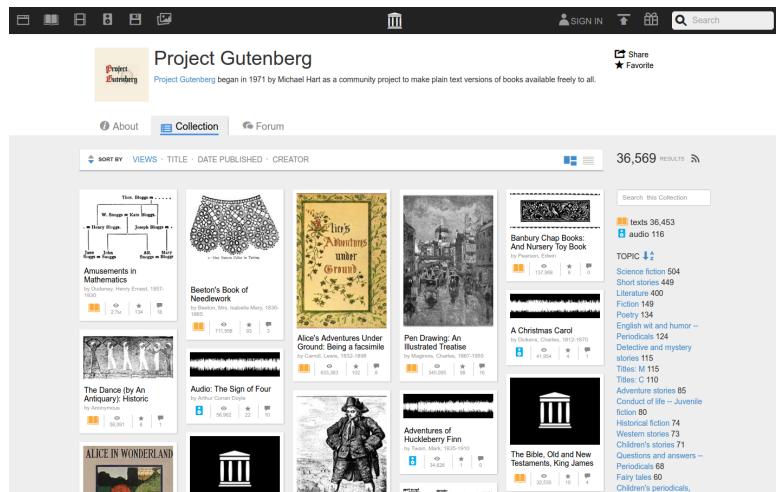


Figure 5. Screenshot of Project Gutenberg interface at of archive.org.]

The Library of Congress offers several digital collections that range from performing arts to legislative texts. The collections that include textual documents are: Historic Newspapers, United States Legislative Information, and Web Site Archiving. All of them offer a search box with filters, and a taste of the collection, e.g. in the case of historical newspapers, the home page of the collection offers images of the newspapers from one hundred years ago of the day of the visit of the page (see Figure 6). The search results are a faceted list that includes images of the newspaper page with the query highlighted in it.

CHAPTER 2. CONTEXT AND BACKGROUND

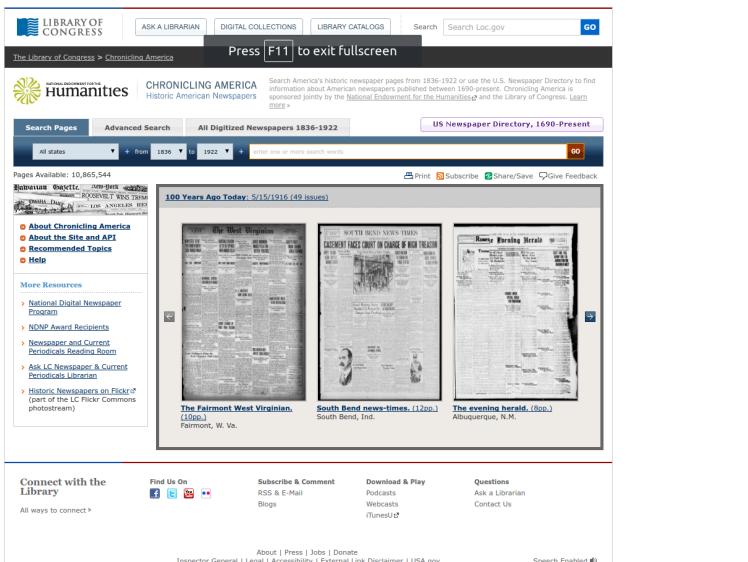


Figure 6. Screenshot of the “Women’s history” collection interface, at the Library of Congress (USA)]

A differentiated case is the Europeana site (Europeana, 2008). While the site itself uses conventional approaches, they are innovating via other platforms mainly through data sharing. Europeana is an aggregator and a meta-search engine that incorporates external resources in order to add features to its contents. The strategy to offer not only text-based navigation and search in Europeana has two main directions. On one hand, Europeana uses external services like exhibitions at Google Cultural Institute, and they have active accounts in Pinterest, Facebook, Google+, and Twitter. On the other hand, Europeana maintains and promotes a rich API to encourage others to work with their materials. Related to standards, Europeana is implementing interface standards in faceted lists of digital collections, but also promoting the reuse of the contents, as mentioned above, in commercial free services.

In addition to the standards, this year, 2016, the New York Public Library (NYPL) integrates innovative ideas, such as the recently published its public domain collections in a revolutionary way. NYPL created a Git repository that published the digital objects of their public domain collection of images. They publish a snapshot of the collection regularly. That encourages people to work on the collections and propose innovation. Like other institutions, the NYPL inno-

vates through an internal group called NYPL Labs. To present the collection of multimedia objects, (more than 186,000 initially released), the interface presents the collections and the items of each collection as a standard faceted list (NYPL, 2016a). NYPL also presents several experimental interfaces for the collections, like a long mosaic of small images representing 180,000 images under the public domain (NYPL, 2016b).

Back to the NYPL several works involving collections of textual documents can be found: "Navigating The Green Book", (Foo, 2013) and "Gutenberg authors" (NYPLlabs, 2015) that offers visual ways to explore the text and the authors of Gutenberg books. The "Green book" has a derived interface that shows the contents of the book in a map, as the book is a travel guide.. More generally "The Networked Catalogue" is an interactive space-metaphor exploration tool of the NYPL catalog, by topics and similarities (REF). It uses only Library of Congress Subject Headings (standard descriptive metadata), but looks for co-occurrence of subjects in metadata items, that combined with the number of items per category, creates a graph of relatedness between subject terms. Also remarkable the "NYPL Archives and manuscripts" network generator that allows exploration by dynamic building of a network of relationships (REF). The interactive view offers access to each item. These NYPL initiatives offer promising ideas and prototypes that in the future could be more widely adopted.

This sub-section has presented a brief review of what big DL such as Internet Archive, Library of Congress, Europeana, and New York Public Library, offer as interfaces to their digital collections, and especially to their collections of textual documents. The practices of these institutions are, here, considered as standard practices for text collection interfaces. The review also mentions more innovative proposals from the NYPL. The following sub-section focus on data visualization approaches used for collections representation and exploration interfaces.

2.3.3. Data visualisation and interfaces to text collections

The simplest definition of data visualization could be "the science of visual representation of data" (Friendly and Denis, 2001). Indeed, any interface is a concrete representation of an abstract data structure (including a list of search results, for

e.g.). As shown in the following section, when building interfaces to digital collections, data visualization can be used as an element of the interface, and as a full interface by itself.

This review shows that visualisation techniques (designed and systematic relations between data features and visual features) are increasingly integrated into text collection interfaces in non institutional sites.

The case of the German digital library visualization is remarkable. The library hosts the cultural and scientific heritage of Germany in digital form. It contains a continuously growing list of collections of all kinds of digital objects. At the time of writing, the collections contained more than eighteen million objects. While the official collection site uses a standard text-based interface (DDB, 2016), researchers from Potsdam University of Applied Sciences have developed an experimental visualisation: "Deutsche Digitale Bibliothek Visualized" (Bernhardt, 2015). This visualization interface offers four options to explore the collection: Periods & Sectors, Keywords, Places & sectors, and Persons & Organizations, respectively based on timeline, tags, map, and network representations. The visualisation is based on the faceted metadata exposed by the DDB API (see <http://infovis.fh-potsdam.de/ddb/>). The proposed visualisation dynamically shows changes in the facets, which helps to show the contents of the collection. The "navigation" aspect of this visualisation involves linking to the official collection and triggering a faceted search —returning a list of documents. As well as showing visualisation techniques, this example shows how the relationship between DL and interface is now flexible —that we can have multiple interfaces, and both official and experimental can co-exist online.

The above mentioned four options to browse the collections of the German digital library combined with the type of objects (Archive, Library, Media, Research, Museum, Monument protection, other), offer interactive data visualization tools that are elements integrated in a standard interface—standard in the sense we defined in section 2.3.2. The interface of this DDB visualization does not need special directions or tips to be used for a first time visitor because of the clear use of established conventions, i.e. timeline, word cloud, use of interactive widgets and tool tips that help to explain the visualisation during interaction. The visualisation element actually occupies 90% of the page, but it does clearly sit within a conven-

tional web page structure (e.g. header navigation, introduction page, etc). It is also worth noting that the interface is delivered using web standards and does not require any special software at the user end.

This practice of integration of visualization elements can be found also in "Explore Australian Prints + Printmaking" (Ennis and Whitelaw, 2014). This project represents the prints collection of the National Gallery of Australia that contains "54,158 works, 26,612 images, 19,949 artists, 3081 galleries, 8726 exhibitions and 9090 references". As in the DDB visualization, the site offers on the home page, in addition to the conventional search interface, several ways to discover the collection: by Subject, Timeline, Works and Networks, Decade Summary, and All Artists. Again in this example, elements and techniques of data visualization, mostly based on metadata representation, are integrated into a web interface.

Another example in this direction is the case of Eugenics Archives by the Social Sciences and Humanities Research Council of Canada (Collective, 2016). This rich archive collects information from Canada related to the practice of surgical sterilisation for those deemed "mental defectives", in order to improve the genetics of future generation. This practice was in effect until 1972. The home page of the archive offers "twelve interactive tools to explore this archive" (Encyc, World, Game, Connections, Our Stories, Timeline, Players, Institutions, Interviews, Pathways, Media, Database) that "reflect the collaboration of scholars, survivors, students, and community partners in challenging eugenics". The range of options goes from pure narrative, to charts, videos, and network relationships.

According to reviewed cases, the use of the home page of DL of collections as a portal of interfaces is growing. In all cases the classic catalog-based view of items can be found, but next to it is offered a list of possibilities. The proposals try to avoid the necessity of an initial explanation on how to use and interact with the interface by using established and well-defined techniques. Having in mind this idea of data visualization elements integrated within an interface, in this section, a review of techniques follows for the visualization of collections of digital objects, and in particular, textual document collections.

Spot by Jeff Clark (2009) is a real-time representation of tweets grouped by Banner (by similarity), Timeline, by User, by Word (word with related tweets around), by software used, by groups (tweets with same words). This case again

uses data visualization elements integrated giving several options to discover the collection. In this case the groupings come from the analysis of tweets (such as similarity among tweets, and tweets related to a specific word) combined with their metadata (author, date, software client). This same idea of multi-representation can be found also in static, almost infographic representations of large collections of items, like the case of "X by Y" by Moritz Stefaner (2010) that represents almost 40,000 project submissions to the Prix Ars Electronica, from the early beginnings in 1987 up to 2009. The items are grouped in several ways (by country, by topic, by year, by prize, and by category of the prize) so the overview of the collection is multiple at a glance.

Other strategies to represent collections of texts are not a list of data visualization elements integrated within an interface, but a full proposal that becomes an interface by itself. This is the case of Newsmap by Marcos Weskamp (2004) that offers a dynamic treemap with the latest Google News ordered by categories (colours) and weighted by area sizes according to the number of related articles that the Google News aggregator included as relevant. This representation of the collection of latest news is dynamic and helps to show the changes in the stream of news.

Another example of data visualization as an interface is the case of Grimm's Fairy Tale Metrics by Jeff Clark (2013). Clark represents the sixty-two stories of the "Grimms's Fairy Tales". It designs a multi-sortable table that allows comparison of the tales in several ways: physically (by length, by lexical diversity), and by topics (analysing frequencies of keywords). Each cell gives a graphic score (a bar) of the respective property. This is a powerful tool to compare documents, in this case, from the domain of literature.

In other cases the strategy is not innovative in the presentation of the visualization, so the users will understand the interface with no need of introductory explanations. This is the case of Word Storm by Quim Castella and Charles Sutton (2014), that represents a collection of journal papers from a conference as a special faceted list. Each paper is represented as a word cloud, but the authors modified the algorithm that generates the word cloud in a way that the words appear always in the same position for all the documents of the collection. This admirable trick makes word clouds visually comparable, and, therefore, the traditional static list

of documents becomes substantially enriched.

This review shows that data visualisation can be used as element(s) of text collections interfaces. Sometimes the data visualisation elements can represent the whole interface, sometimes the elements are integrated within standard interfaces. The classic tasks of search and detail-on-demand via a catalog entry are always present in the reviewed cases of institutional DLs. A trend seen across the reviewed cases is a home page that offers several overviews and/or options to interact with the collection. At the same time the reviewed cases try to avoid features that require an initial explanation on how to use the interface.

The second part of the review focuses on data visualization projects that represent textual documents. Comparing to the reviewed institutional DL interfaces, the data visualization interfaces —usually done outside of the official institutional sites— show state-of-the-art techniques, experimentation, and innovation, including the use of text analysis. For example, Spot is a sort of twitter dashboard to follow trends, people, keywords. The data is dynamic; the input tweets are constantly analysed, indexed, and visualised. By contrast there is a lack of text analysis and data visualisation tools in institutional DL interfaces. Since the aims of this research project are to create interfaces that support and invite to read text collections, the next sub-section reviews the concept of e-reading, and the standards of reading interfaces.

2.3.4. Interfaces and reading

Since the first particular aim of this research project is to propose practical techniques of interface design, and development for text visualisation and exploration techniques oriented to read, this sub-section reviews the efforts to design interfaces whose primary goal is to support text reading. The scope of this review is the integration of texts from collections into interfaces that support and prioritise reading. This sub-section presents a brief review of the concept of active reading, and briefly considers the wider concept of readability. To close the background and context chapter, a review of interfaces to text collections that allow online reading is presented.



Figure 7. Bookwheel, a mechanical reader designed by Agostino Ramelli from his book *Le diverse et artificiose machine*, 1588.

Since the beginning of electronic displays of texts, the study of text reading on screens and on paper is a productive field that evolves rapidly. A demonstration of that is that in 1997, a highly cited paper of O'Hara and Sellen (O'Hara and Sellen, 1997) presents the multiple advantages of paper versus screens when reading texts. The cited advantages of paper include: supporting annotation while reading, quick navigation, and flexibility of spatial layout. These "advantages" today look outdated since modern reading interfaces from applications for e-readers for tablets, mobiles, phablets, etc, include them by default. The action of reading, and, simultaneously, underlining, highlighting and commenting the text was defined as "active reading" in 1972 (Adler and Van Doren, 1972). Today, a variety of software with these features can be installed on any digital device. So, in fact any device used as an e-reader is, potentially a device ready for active reading.

In this context it is necessary to visit, briefly, some concepts in the umbrella of text readability. Text readability has been measured in a wide diversity of ways

and for a wide diversity of purposes. The measure can take into account elements of graphic design, typography, human perception and psychology, legibility, vocabulary, syntax, etc. In consequence a number of methods to measure and calculate readability (Tinker, 1963, Fry, 2006) have been published.

According to the scope of this research project, we are interested in how to increase readability, that is, how to make a text more accessible through reading. Independently of the method used to calculate readability, there are two differentiated ways to improve it. On the one hand, there are methods to increase readability by modifying the contents of the text by using text analysis techniques to make a specific text easier to understand to people, including non-native speakers and those with special needs, poor literacy, aphasia, dyslexia or other language deficits (Ferraro et al., 2014). On the other, readability can be improved through modifying the visual aspect of the text and the interaction with it, through a specific hardware-software interface. Adding graphic elements to a text can improve comprehension of the text (Ferraro et al., 2014, Suominen et al., 2014).

In 2015 there are many online services that are offered to a reader embedded in the browser. Some of the services are commercial SAAS (software as a service), allowing online management of documents, with annotation, text-highlighting, share button, bookmarks. This is the case of Issuu, Google docs, Scribd, PDF.js, etc. In other cases, the publishers offer e-readers as software, in addition to e-readers as a device; examples include the Amazon Kindle app and web page, Barnes & Noble's NOOK book app, Google play, etc.

Back to DL and text collections, Table 2 reviews some institutions' online readers. In this table we find two DL that offer powerful online readers with multiple features, these are the Australian Trove, and the Internet Archive. Both offer a book metaphor, share, read-this (English), presentation mode, zoom, 1-2 pages viewer, and print. Trove has extra features such as: tools for social tagging, commenting and transcription, transcription reader, download, Citation formats. The Library of Congress offers multiple formats for the texts, mainly HTML pages. Europeana, due to their aggregator nature, redirects to the reader at the referenced institution page.

Reviewed site	Reader features
Internet Archive	Book metaphor, share, read-this (English), presentation mode, zoom, 1-2 pages viewer, print.
Library of Congress (USA)	Several formats: from HTML to PDF. When available it offers read-this (English).
Trove	For newspapers: tools for social tagging, commenting and transcribe, transcription reader, scanned image, zoom, print, download, citation formats,
Europeana	Redirects to the reader of each institution

Table 2. Reader features of some major institutional DL

Another aspect related to the way we read digital text is about fragmented communication. The number of messages we receive per day has increased, as well as the number of channels delivering short multimedia messages. In the past decades several works have explored the possibilities of breaking the linearity of a text. The philosophers Deleuze and Guattari have described the rhizomatic structure of knowledge, which inspires this project too: "In a book, as in all things, there are lines of articulation, segmentarity, strata and territories; but also lines of flight, movement, deterritorialisation and de-stratification" (Deleuze and Guattari, 1987). The complexity of the knowledge produces a multiplicity of narratives. In the novel *Hopscotch* by Cortázar (1966), the author proposes two reading orders for the chapters; the text starts with: "In its own way this book is many books, but mostly it's two books". Another relevant work is the Project Xanadu from 1960 (Ted Nelson, et al, 1960) considered the first hypertext project in the digital era, and it was a visionary definition of standards for the WWW that were mostly not included in the standard protocols. One of Xanadu's rules states: "Every document can consist of any number of parts each of which may be of any data type". The open Xanadu project encourages non-linear navigation of text. The aim of Xanadu's demo is to demonstrate the possibilities of hypertext. These are the main examples that prompt this project to investigate the effects of reading in alternative ways in combination with normal reading.

esides that, the review presented several cases where narratives can be read in more than one way. These cases encourage the idea, developed in detail in CH???, of segmenting long texts as a strategy for reading. This idea of narrative

multiplicity is developed in detail in chapter 6.

2.4. Summary

This chapter presented three diverse and complementary brief reviews..Many big public DL institutions are becoming meta-libraries that catalog digital objects not only hosted locally, but from other institutions, making it possible to search across multiple collections. The extreme of this, in the studied case, is Europeana, that only works as an aggregator, not hosting its own objects. One way or the other, DL collections are now far too large to be readable in any traditional sense. This posthuman scale reinforces the urgent need for tools to deal with this volume of materials. At the same time, there is a contrast between this posthuman scale, and the limitations of the standard interfaces based on lists, facets, and standard metadata that we find in the official pages of the reviewed institutions.

The availability of large digital collections in DL offers significant opportunities for computational text analysis. The examples reviewed here show that techniques such as topic modelling can help represent document content. A simple classification of levels of data in indexing digital objects of a collection is proposed: Raw metadata (which uses only metadata), Extended metadata (which enriches metadata with simple operations), and Object analysis (which includes new information about the collection through analysis of the text of the documents of the collection). The DL reviewed do not use Object analysis to enrich the data about the collection. Examples of text analysis in text collections can be found in practitioners and researchers working outside institutional DL. In 2.2.3 several cases of data visualization that include text analysis are reviewed. This third level which includes text analysis can add subjectivity in the interpretation of the contents of the collection, but at the same time, it can improve the representation of the complexity and abundance of the collection —as seen in 2.2.3 "Mining the dispatch" and "Topic Modeling Martha Ballard's Diary". In this direction, Johanna Drütcher has influenced this research with her concept of *capta*, in opposition to *data*, that proposes an active attitude in interpreting the data and metadata of a collection. The section ends with a review of topic models, its relation to DH, and its labelling (which is what makes topic models human readable) and evaluation.

Methods for topic labelling go from pure procedural (computer-aided), common in sciences, to pure curated (manually done, usually by experts), common in DH. Procedural labelling is used when the number of topics is high. This variability of labelling nature has produced literature in the evaluation of topic model labelling methods. Reviewed cases show that procedural labelling can be comparable to curated labelling. In any case, the review shows that the use of topic models in text collection interfaces to help exploration and reading is not common.

Thirdly we review the concepts and practice of interfaces to DL and text collections, considering theoretical frameworks and philosophies supporting innovation in interface design over the last twenty years. This review shows design concepts that go beyond the classic search task. Interface design can bring freedom to the user, and improve explorability of collection contents. For the interface standards and trends in text collections, the standards are presented as the features that mainstream DL use. The review shows that interfaces to major institutional DL are conventional text-based web pages. At the same time opportunities to enrich these interfaces are clear; the examples discussed show how data visualisation can be included in collection interfaces, as a whole, or as an element of the interface. This review shows that more innovative interfaces to text collections are found as data visualisation works. In most of times, these works are not integrated into the official sites that host the collections. Finally a short review of interfaces for reading texts on screens and displays shows that the main features for e-reading seems established. Besides that, several cases that experiment with the idea of reading texts in others ways than the usual linear reading, have brought possible solutions for reading parts of collections too big to be read from beginning to end. The following chapter analyses the reviews and identifies some gaps and opportunities for research and innovation in the development text collection interfaces..

3. Gaps, research questions and contribution to Knowledge

This chapter starts by exposing the gaps and opportunities found from the review of background and context of this research project. Then it introduces three research questions that the project intends to answer. Finally we present the contribution to knowledge that this research project has produced. Since this is a practice-led research project it has produced specific digital artefacts in addition to contributions to theory. Detailed information about these creative works, and an analysis of the overall outcomes of the project can be found, respectively, in Section 5.1 and Section 6.1.

3.1. Gaps

After the diverse reviews presented in the previous chapter (2 Context and background) a list of gaps and opportunities for knowledge inquiry and contribution were found. Due to the diversity of the three areas reviewed —digital libraries, text analysis and digital collection interfaces— the findings are diverse and complementary. The process of developing the artefacts involved experimentation in all three fields (see chapter 5).

Following the order of the review, (see chapter 2), the first block reviews digital libraries and text collections: the cases reviewed are major DL that are considered mainstream, and standards as contemporary interfaces to DL in general, and to collections of texts in particular. This first review shows the contrast between the posthuman scale that DL are reaching in their continuous growth, and the poor innovation in the interfaces of the official pages from big public institutions. These institutions host DL that grow everyday, but they use text-based interfaces that

are ineffective at representing collections. This contrast reveals an opportunity in the development of interfaces to large collections of documents that fits with the aim and motivation of this research project.

The second review is dedicated to text analysis. The main gap found here is that the output of text analysis, in particular the topic modelling of text collections, is generally applied at a collection level. As with Moretti's other "distant reading" methods, topic model analysis is commonly used to show overviews of the collection, distribution of topics ("Martha Ballard's diary"), comparison of topics and/or documents and, similarities among documents ("Mining the Dispatch"). In contrast, no cases were found dedicated to representing single documents, or parts of those documents. The generated representations of the collections are presented as diagrams, charts, maps, but rarely are the results of the analysis integrated in the interface. The cases found that integrate the outputs of the analysis of the collection in exploration systems are for image collections, as in "Discover the Queenslander" (Whitelaw, 2014), where the images are analysed to extract a palette of colours, allowing browsing the collection of images by colour. We see an opportunity in integrating the results of text analysis of text collections into interfaces as features to support exploration and reading.

As seen, topic modelling is a powerful tool that it is still new in terms of application. Most of the reviewed approaches do topic model analysis for one or more of these three reasons: to compare topic models themselves, to evaluate them, or to somehow, represent the collection as a whole, for overview, evolution, and comparison. Reviewed interfaces are limited in supporting reading. For example "Mining the Dispatch" shows us "exemplary articles" for specific topics, but the documents are not central, i.e. the documents are exemplars of the topic, rather than the topics being a guide to the documents. "Martha Ballard's Diary" is not an interface, so does not encourage exploration, but again the large scale analysis emphasises topics across the diary, instead of its contents. Obviously, all authors present topic models related to collections of documents, but rarely are the documents central in the approach. In particular topic model analysis is rarely used to support reading documents, and this is a significant gap.

More investigation about the appropriate number of topic labels for each collection and context is necessary. Literature shows that the list of topics that the

algorithm outputs, can be modified by a curation process (Nelson, 2010a). There are tools that help in that process: MetaToMATo (Metadata and Topic Model Analysis Toolkit) (Snyder et al., 2013), is a web based application that allows of collections of texts combining metadata and topic models. It can filter documents by topic, and summarise views with metadata and topic graphs. Another related approach is the prototype TopicExplorer (Hinneburg et al., 2012) that combines topic modelling, keyword search and faceted lists to explore a large collection of Wikipedia documents and other collections. In the process of reviewing topics by curation, some topics can be removed because they are considered too general, or semantically unusable; some can join other topics, in cases where topics overlap or are included semantically in others, or are, in fact, synonyms. The literature and practice reviewed points to opportunities for experimentation in the process of labelling topic models.

In topic model labelling the literature recommends choosing inclusive labels. As an example, in Nelson's "Mining the Dispatch" one of the labels is "fugitive slave ads" (Nelson, 2010b), but some of the texts are not "ads" so Nelson concludes that topic models are better suited to representing the prevalence of a topic in the collection, than to strictly and accurately classifying every text.

According to literature, there is no a clear definition of how to optimise topic model analysis parameters, in particular the number of topics to generate related to the size of the segments of text that make up the collection. In determining the size of the segments, researchers try to find any natural pre-existing segmentation, such as pages of a diary, as in "Martha Ballard's diary (Blevins, 2010). Some research opportunities appear in experimenting with these parameters through the development of real cases. See chapter 6 for a detailed explanation.

In reviewing the cases a generalisation was found: on the one side experts in analysis, usually mathematicians/statisticians, and computer scientists, conduct machine-only processes with no curated methods, other than dataset annotation for training, and, hence, a qualitative evaluation of the machine-based method. On the other side, humanities researchers can find computational methods too opaque due to knowledge barriers. The review shows that from the point of view of the humanities, it looks acceptable, and desirable to mix curated and procedural methods in the same process. This links with the mentioned Drucker's concept of

capta - if data is "taken" for interpretation then curated labelling is a form of interpretation too.

In considering interfaces to digital collections, we find a common interaction that follows the schema of Schneiderman's Mantra. An initial overview allows contextualisation, and/or understanding of the collection, from a faceted lists to a geo-map, from a network visualization to a treemap. After the overview then an interaction shows collection details, and then an item can finally be reached. Some approaches just follow one or two of these three schematic steps. For building interfaces for reading, this approach does not take the shortest path to a text in the collection. Following the three steps strictly, the visitor needs several clicks – decisions — to read a text in the collection.

This suggests an opportunity to revisit the "overview first" vision and see what can be done to put the contents of the collection in the centre of the scene. Even though the overview is necessary, useful, perhaps inevitable, sometimes the overview can act as a substitute for the content. This observation was influential in developing channels to work in through the creation of the artefacts. Lessons learned about this question are analysed in chapter 6.

To build a collection of texts interface, from the data to the last graphic element of design, can be a process with many steps. Some of the possible steps could be based in manual — done by experts — or in procedural — done by computers — methods. In this review a variety of these two methods for similar tasks has been found. Neither of the two methods is *a priori* better than the other. The reasons for choosing one or other method depends on the goals and context of each case.

From the last brief reviewed block about interfaces dedicated to reading text the evidence found is that standards for reading texts on digital devices are already established and offered for most devices and interfaces in contemporary systems. For this reason, when developing interfaces for reading, it is recommended to follow the main features that can be considered standard according to the classic definition of "active reading" by Adler and van Doren (1972).

The following sub-section introduces the two research questions in this research project that cover the two intimate aspects of the sides of the research: the practice led nature of the research, and its theoretical findings.

3.2. Research questions and Contribution to knowledge

One main question is addressed in this research project. The question refers to the practical case of interfaces to text collections, and includes the theoretical methodology and recommendations documented based on the creative work..

The aim of this project is to answer the question: how can we create interfaces to text document collections that let us explore the collection by reading its contents? This question is answered with the creation of real artifacts, as well as through theoretical issues presented in this dissertation, which are documented and described as new methods for creating text collection interfaces focused on reading.

As the contribution to knowledge, a group of digital artefacts are presented in the form of visualisation interfaces to text collections with a focus on exploration-by-reading, and on a rich "back and forth" path between collection inner view and overview. "Inner view" is used here in opposition to overview, in the sense that inner view allows a perspective of the collection from inside, looking at a sampling of the details, instead of a bird's eye view, or overview. Proposed approaches use relationships among collection items as directives to browse the collection, instead of "details on demand" as the end of the turn in the game of an information seeking session. In the approach presented here the reader can move between related detail views, or back to a collection overview and dive again into the contents of the collection.

The collections chosen to work with in the presented artifacts include a range of contents that range from scientific papers to art and historical collections. As a direct consequence of this research project, these collections have gained accessibility, as in readability, and, finally, the collections have gained presence since the created interfaces are accessible online. Therefore, this project contributes to expand public knowledge resources.

The methods used in the analysis of the texts are innovative in the use of state of the art techniques, such as topic models, entity recognition, and text similarity. Beyond one-to-one categorisation of collection objects, topic model analysis is used here to create rich and complex multi-dimensional relationships between

segments of texts. This gives a more complex way to characterise the collection contents, as well as enabling navigable relationships between collection elements. This technique is an essential element of the "deep interfaces" approach developed in detail in chapter 6s.

Since the interfaces developed support and prioritise reading the documents of text collections, it is important to consider how to read collections of texts that are too big to be fully read. The answer to this question is the introduction of the concept of crossreading. Crossreading is the non-linear reading of large collections. To adapt a collection of texts to crossreading, a process is necessary: the texts that belong to a digital collection are divided into small segments, and these segments form the objects of a new collection to be analysed. Once this new collection is analysed with topic modelling, new multidimensional metadata is associated with each object. As a consequence, new ways of exploring and reading the collection appear.

As a contribution to theory the process of generating what we call deep interfaces to text collections is described in detail, so the process can be reproduced with other datasets, and in other contexts. Every step of this process has been validated, and improved through a cumulative development process that, following a reflexive methodology, has produced a set of individual artefacts presented in this dissertation.

All documents and their sources related to this thesis, including the research project, the literature review, source code, source files of images and documents, and publications can be found in the Github repository at <https://github.com/jaumet/myacademydata>

In this chapter, the outputs and opportunities for new research have been presented according to the reviews introduced in Chapter 2. Then the two research questions that this project intends to answer have been introduced. Finally the contributions to knowledge were presented. The following chapter introduces the methodologies, and methods used during the research project.

4. Methodologies and methods

The list of methodologies introduced in this chapter is combined in order to fit the idea of a research based in the production of interfaces to real cases of text collections, and the idea of learning out of the research process in order to compile a list of good practices that could be applied for other projects and by other researchers. Firstly the methodologies will be introduced and, secondly, a list of methods and techniques used in the research project are presented.

4.1. Methodologies

Practice-led and practice-based research: This is a project based on practice, where software creation is the specific practice. According to Linda Candy (Candy and Studios, 2006), there are two types of practice related research: practice-based and practice-led. In her words:

1. If a creative artefact is the basis of the contribution to knowledge, the research is practice-based.
2. If the research leads primarily to new understandings about practice, it is practice-led.

This project is presented as a combination of these two practices. According to the mentioned aims of the project, on the one hand this project has created a number of software artefacts that have contributed to knowledge, making the presented collections more accessible. In that sense this project is practice-based. On the other hand, the cumulative process of creation of artefacts has produced theoretical issues and a methodology to build interfaces to text collections. In that sense the project is practice-led.

This double practice combination of research methodology follows Greame Sullivan's braid metaphor, where creative practice cannot be separated from research theory, in other words, the complexity of that relationship, as in a braid of fibres, reveals all kinds of structures among practical and theoretical research (Sullivan, 2005). Sullivan's metaphor explains the richness, the complexity, and the fruitful relationships between practice and theory in this research project.

Reflexive methodology The symbiosis between practice and theory in this research project is embodied in a process of moving back and forward, learning from experience, developing new understandings and reapplying these to the production process. The reflective practice produces improvements during development due to continuous learning and questioning. I used a reflexive methodology inspired by its tradition in the humanities and creative arts. Sullivan talks about four kinds of reflexive practices: self-reflexive, reflexive, dialogue and questioning. Even though all practices can be related, and used simultaneously, this project is a reflexive dialogue: "the plausibility of an interpretation of research findings will be determined in part by the capacity of the reflexive researcher to openly dialogue with the information" (Sullivan, 2005).

From a more pragmatic point of view, this reflexive idea of continuous inter-practice inquiry could be understood as "learning by doing" as defined by Roger C. Schank (Schank, 1995), that is: since this research project produces artefacts, those act as doing-devices, then it is possible that the learning by doing is a feasible practice.

The main reason, and, in fact, advantage that the reflexive methodology brings to this research project is that it allows experience and insight to be exported from one artefact to the next one in a cumulative process. In fact, the final artefact (Diggers Diaries) is the cumulative embodiment of knowledge and techniques developed during the development of previous artefacts.

Empiricist methodology Empirical methodology is based on sensory experience and the evidence that can be extracted from it. Applied to this research project, sensory experience comes with the creative production of digital artefacts. The evidence comes from the validation of the conducted processes through two forms

of evaluation:s

- Qualitative evaluations: online questionnaires
- Quantitative evaluation: semi-structured interviews, public exhibition at the museum, feedback from users.

4.2. Methods

This section introduces the methods used during the research project. The methods are included in a first subsection, tools and techniques that have allowed the studies and experiments conducted successfully. In a second subsection the methods of evaluation are listed, and how these validate the usability and goals of the generated artefacts is discussed.

4.2.1. Tools and techniques

Curated —conducted by humans— and procedural —conducted by computers— methods are used across the multiple stages of the development of each of the presented artefacts. A list of methods grouped by task and subtasks follows:

Data gathering: Data for the projects was gathered through a combination of procedural methods (APIs, spiders, and scripting for scraping), and curated methods (eventual manual download, annotation, first classification for storing).

Text analysis: choices and evolution:

- Text document segmentation:

Several of the interfaces rely on the segmentation of source texts to enable a cross-reading approach. Thus the process of segmenting texts is a key method. Three methods were tried here: a curated segmentation (Crossreads II, see 5.3), a procedural segmentation (Crossreads I, see 5.2), and a natural segmentation (Diggers Diaries I and II).

In a curated (manual) segmentation process, one or more editors do their personal segmentations of the texts of each document of the collection. In a procedural

process, a software splits documents into pieces based on length, sentences, paragraphs, sections, etc. The reason for these two approaches is that the segmentation task is very subjective. A human expert could add a personal view to the segmentation (Crossreads II). A procedural segmentation (Crossreads I) can accomplish well this task in terms of size of each segment, but it cannot be expected to have the richness of a segmentation curated by an expert.

For the final project, Diggers Diaries (I and II), I found a naturally segmented collection. The collections of diaries are naturally segmented into handwritten pages.

- Topic model analysis

I used manual algorithms and automated tools like MALLET, "a Java-based package for statistical natural language processing, document classification, clustering, topic modelling, information extraction, and other machine learning applications to text." (McCallum, 2002)

- Topic model labelling

I used curated labelling plus labels grouping. This generates a two level hierarchical menu of topics. The topic grouping and labelling process is discussed in depth in Section 6.

Technically, the web application is only-client-side Javascript. Generic Javascript libs used: angularJS, bootstrap, and jquery.

4.2.2. Evaluation and validation

In order to evaluate the generated artefacts, studies were conducted on how users accept and use these interfaces. For the case of Visference the evaluation focused on new features introduced through text visualisation, by comparison with existing or conventional presentations of text. In the case of Crossreads there is no existing interface to compare to; the tool provides a novel exploratory presentation of specific text. For this reason, Crossreads was evaluated using a semi-structured interview only. The results of these two evaluations informed the design of the last artefact, Diggers Diaries, which follows the conclusions of both studies.

The evaluation aims addressed the following questions:

1. Do users detect the new features?
2. Do users still prefer or require access to the existing presentation?
3. Do users understand the new features?
4. Do users feel confident and positive about using the new features?

According to these aims, the evaluation experiment was based on the established technology acceptance model (TAM) (Davis et al., 1989), and task technology fit (TTF) (Goodhue, 1995). TAM attempts to understand why people accept or reject information technologies. TTF says that technologies will be used if, and only if, their available functionalities support the user's activities. Consequently, its focus is on the match between the user's task needs and the available functionalities of a given technology. Aims 1 and 2 are directly related to TAM. Aims 3 and 4 are related to TTF. The questions have been designed following (Taylor-Powell, 1998).

The evaluation was in the form of a mixed method approach, including:

- a)** Online questionnaires for standard users. Here "standard users" means potential visitors of each real scenario, this is, readers of academic journals for Visference.

Online questionnaires are the most popular method to gather quantitative data for statistical analysis from users. Questionnaires allow participation from an unlimited number of people and can collect data on knowledge, beliefs, attitudes, and behaviours (Taylor-Powell, 1998). Online questionnaires also make it easy to protect the privacy of participants.

- b)** Semi-structured interviews with expert users. In order to complement the online questionnaire and provide open-ended evaluation and feedback on each tool from experts a small number of semi-structured interviews with experts in each artefact's scenario or domain was conducted..

A semi-structured interview is a qualitative method. It is open, in the sense of participants expressing their ideas freely following pre-defined questions, allowing new ideas to be brought up during the interview (Kitchin and Tate, 2013).

CHAPTER 4. METHODOLOGIES AND METHODS

The results of the online questionnaire can be found in the Appendix. The analysis of the questionnaire results, as well as the questions and opinions from the semi-structured interviews are discussed in chapter 6 Discussion.

5. Results: the artefacts

This chapter presents the results of the research project. Due to the practice led nature of the research the contents of this chapter have been shaped to introduce the generated output of this research project: the artefacts.. Analysis of results, lessons learned, and general discussion occurs in chapters 6 Discussion and 7 Conclusions and future.

Since all the artifacts are digital, all of them are accessible online. Nevertheless, there is a website that lists and introduces each one as part of this research process. This site is called “Web of Artifacts” and it is accessible at <http://research.nualart.cat/woa>

The structure of the chapter documents all the artefacts created including for each:

- A figure showing a detail of the generated interface
- Artefact name and description: name, version and short description of the artefact
- Narrative: a brief description of the artefact highlighting including motivation and aims, as well as its contribution to the thesis as part of the process evolution.
- Datasets: description of the dataset used for the project, and, in case it exist, description of complementary datasets generated during the research or taken as external sources.
- Collaborators: a list of direct collaborators and their field of expertise, in relation of each artefact.

- Data process: a description of the process related to the data including gathering, reformat, and conversions, analysis, final format, storage and access.
- Data analysis: detailed description of the analysis done, and used tools.
- Interface development: description of challenges and techniques and external libraries used.
- User evaluation studies: presentation of the results of the conducted user studies. The detailed list with numbers, and charts can be found in Appendix A.
- Project outcomes: list of all outputs that the project generated; it can include publications, exhibitions, presentations, etc. Access to all the code is included.

5.1. Visference: Journal of Machine Learning Research



Figure 8. Detail of Visference

Artefact name and description

Visference is a visualization tool for the exploration of conference papers. It uses topic model analysis and sorting table techniques

Datasets

This project has an initial dataset (dataset-1), and a created one (dataset-2)

- Dataset-1: collection of 282 accepted papers from the JMLR Workshop and Conference Proceedings Volume 28: Proceedings of the 30th ICML.
- Dataset-2: we create a golden dataset as a representation of machine learning domain. This dataset was compound by:
 - A list of classic books about statistics, including: Alpaydin Machine Learning 2010, Mackay, Barber, Machine Learning-Tom Mitchell, Data

Mining-Practical Machine Learning Tools and Techniques, Pattern Recognition and Machine Learning CM Bishop, ESLII, RW, Smola Book, Learning Kernel Classifiers.

- Every papers under stats from arxiv since 2010 to 2013

Narrative

Motivation: Most conference proceedings present their content as a one-dimensional, non-interactive list of papers on a web page. However, the reader of this kind of presentation might not know the reason for the paper order, does not get an overview of the contents or relations between the papers, and has very limited search and filtering functionalities available. A list of the identified problems includes: a negative features of a flat and non-interactive list, no sorting options, no overview of the dataset, no relationships among items, only CTRL+f (or COMMAND+f) for searching, no filtering.

Aims: To explore more effective interfaces to represent the contents of conference proceedings. Nevertheless, this prototype is visually conservative in the sense that it tries to not scare usual visitors with unfamiliar visualization but incorporate known interactions, like table sorting. The tool has been evaluated positively by users and experts.

Collaborators

- Dr Wray Buntine, Machine Learning Research Group, NICTA. Specialist in topic models.
- Dr Mark Reid. Machine Learning Research Group, NICTA. Editor of the JMLR.

Data process

- Dataset-1 gathered from a bibtex file
- Dataset-1 transformation: papers in PDF were transformed to text, and then all papers split into pages, producing about three thousand segments of text.

- Dataset-2 transformation: books and papers in PDF were transformed to text, and then all of them split into pages producing about fifty thousand segments of text.
- Final topic model output parsed and transformed to JSON

Data analysis

- From dataset-2, the topics were generated.
- Then topic models for dataset-1 (small) were inferred based on the current dataset-2 (big) trained model.
- Topic model labelling: We asked experts from our group (Machine Learning Research Group at NICTA) about which labels they suggest for each topic. Topic model results and labels can be found online (Buntine and Nualart, 2013)

Interface development

The proposal includes a simple HTML table with the 282 papers as rows and topic models as columns. The table is sortable by columns. Each cell shows a weighted ball proportional to topic model scores.

User evaluation studies

An online questionnaire and semi-structured interviews were conducted to evaluate the acceptance of Visference as a new tool for exploring conference papers compared to the existing static conference page. As average, 80% of the thirty-three participants prefer Visference over the the existing interface for typical tasks. Tasks and percentage of positive answers for Visference are as follows:

- How many papers were accepted in the conference? 91.18%
- Which papers are related to machine learning theory? 70.59%
- How many papers talk about optimisation? 85.29%
- Understanding the topics and themes of the conference. 85.29%

- Finding papers related to your personal interests. 82.35%
- Exploring new topics and discovering new research in this field. 88.24%

See A for detailed results of the questionnaire. Semi-structured interview results are integrated into Chapter 6.

Project outcomes

- A demo site: <http://research.nualart.cat/visference>
- Publication: poster and a presentation at several NICTA retreats and conferences <http://research.nualart.cat/?p=54>
- Ten topics that can represent the machine learning knowledge domain: SVNs and Kernels, Theory, Policies and Games, Images and Neural Network, Experiments, Definitions, Optimisation, Probabilistic Models, Discussion, Topic and Latent Variable Models
- Code is accessible under GPL licence in github (Nualart, 2014)
- Development charts are accessible in github (Nualart, 2014)

5.2. Crossreads I: Eugeni Bonet exhibition



Figure 9. Detail of Crossreads I

Artefact name and description

Crossreads I, a manner to deconstruct linear narrative text in order to read text in multiple orders.

Datasets

Fifty-seven articles by Eugeni Bonet (Barcelona, 1954), video and cinema artist, and writer.

Narrative

This project introduces the concept of Crossreads. The project comes from a proposal by Museum of Contemporary Art of Barcelona (MACBA); they sought ways to represent a collection of texts in a web interface that can be used online and from their exhibition space. The proposal was to build a section of the exhibition “The Listening Eye - EUGENI BONET: SCREENS, PROJECTIONS,

WRITINGS” (May to August 2014). Eugeni Bonet (Barcelona, 1954) is an important figure for the Catalan and European art in video, cinema and digital media in general.

The proposal includes an idea in the domain of information seeking and discovery called Crossreads that proposes an experimental way of reading texts, as an alternative to traditional linear reading. It proposes to break the initial narrative line of a text by segmenting it into smaller parts. Then, the text is reordered according to similarity scores, which finally offer the reader multiple paths to read the text. The aim of this project is to explore and study the effects when a reader processes fragmented information, as well as to analyse user activity and support reader’s exploration with visualization techniques.

Collaborators

Dr Gabriela Ferraro (Machine Learning Research Group, NICTA). Natural language processing specialist.

Data process

- Documents compilation, conversion to text format
- Procedural segmentation of texts into a total of 710 segments. segment length was about seven hundred characters in total which equates to an average of one minute of reading for an adult (Williams, 1998).

Data analysis

The similarity calculus between segments is done with off-the-shelf Natural Language Processing tools and techniques (Nualart et al., 2014). The analysis steps were:

- Tokenization: words in the segments are separated by whitespace and punctuation characters.
- Stop word removal: standard stop word removal.

- Named Entity Recognition: identification and classification of Named Entities (NE) in each segment. We applied the OpenNLP Named Entity recogniser [2], which distilled four types of entities, Person, Location, Organization and Others.
- Similarity Calculus between segments.

The browsing constraints according to the segmented dataset, was that each segment is linked to the most similar segment to which it is not already linked. The drawback of this approach is that links will have a wide range of similarity scores since, in each iteration, the number of segments to compare with is smaller, and the possibility of finding a segment with a high similarity score decreases. For this reason the limitation of this method is that the similarity between segments is smaller on every new step. A very large collection of texts would eventually solve this problem. However, the benefit is that there will not be any orphan segments, i.e. all segments link to other segments, so the reader will always have the possibility of some crossreading.

Interface development

This artefact was intended to be accessed in an exhibition space, in the Museum of Contemporary Art of Barcelona (MACBA). The interface was designed to be easy to read the contents. In the exhibition space several tablets were accessible to browse the collection of texts through Crossreads interface. One of the tablets was projected on a wall.

Technically Crossreads I was a client-side web app, written in javascript, with the libraries Jquery and Bootstrap. The data is stored in online accessible JSON files.

User evaluation studies

The interface development process was evaluated through presentations and reports with a team from the Museum composed of art historians, librarians, and curators. The feedback received by the organisers was positive. No more evaluations were done after the closure of the exhibition in August 2014.

Project outcomes

- A website <http://research.nualart.cat/crossreads/I>
- This project is part of the exhibition "The Listening Eye EUGENI BONET: SCREENS, PROJECTIONS, WRITINGS" at Museum of Contemporary Art of Barcelona (MACBA) 2014 May to August 2014
- The exhibition catalog: EUGENI BONET: ESCRITOS DE VISTA Y ODIO 2014 (In Catalan and Spanish) MACBA ed, 338 pages. ISBN:978-84-92505-69-2 (<http://www.macba.cat/en/publi-eugenibonet>)
- Code is accessible under GPL licence in github
- Development charts are accessible in github

5.3. Crossreads II: "In your computer", by D Quaranta



Figure 10. Detail of Crossreads II

Artefact name and description

Crossreads II, a manner to deconstruct linear narrative text in order to read text in multiple orders.

Datasets

The book "In your computer", by D. Quaranta

Narrative

This project is very similar to its first version, but it was necessary to develop it for several reasons. Version I was a project where the contents and the artefact itself are in Catalan and Spanish. For the project to be shown internationally at the DL2014 conference, it needed a version in English. A second reason for this version was to try a larger number of segments of text, as well as a different criteria when crossreading. The second version also introduced the random button that allows a jump to a random segment of the collection —as when you open a book

randomly and start reading. Finally in this second version, the code from version I was cleaned and improved.

Collaborators

Dr Gabriela Ferraro (Machine Learning Research Group, NICTA). Natural language processing specialist.

Data process

- The book is published under free licenses and is online accessible. Once downloaded, the text was converted to text format.
- Opposite to Crossreads I, in this second version we did a curated segmentation. The text was divided in 1500 segments, twice the segments in version I.
- Segment length was about seven hundred characters in total which equates to an average of one minute of reading for an adult.

Data analysis

The analysis was similar to the one for Crossreads I (see previous Section). version II differs from version I in defining browsing constraints according to the segmented dataset. In version II each segment is linked to its most similar pairing. To avoid repetition of pairs, segments that have already been set as a maximum similarity segment during ten iterations are skipped. After ten iterations, the skipped segments are used again in the similarity calculus.

Interface development

The principles and technical choices are the same for both versions. Version II of the interface has evolved such that it offers link nuances. For instance, with the right link, the reader goes to its most similar segment.

User evaluation studies

No user studies were done for this second version as it follows the same principles as version I that was evaluated and approved by experts during the development.

Project outcomes

- A website <http://research.nualart.cat/crossreads/II>
- A poster and a short paper in the Workshop "The search is over" part of the Conference Digital Libraries 2014 (DL2014), London (University College London).

5.4. Diggers Diaries I: WWI Diaries

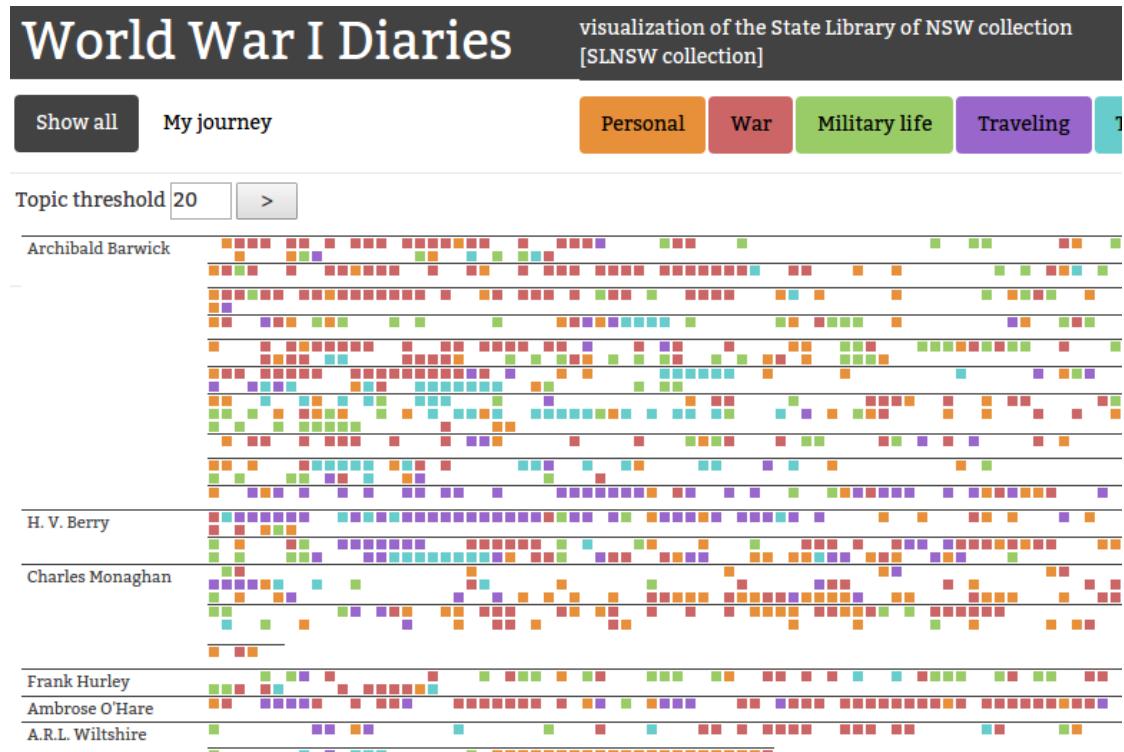


Figure 11. Detail of Diggers Diaries I. This interface has only one way to overview the diaries, this one: by diaries/pages in a vertical list.

Artefact name and description

Diggers Diaries I, an interface to explore a portion of the SLNSW WWI Diaries Collection through general topics

Datasets

126 diaries from WWI as part of the WWI Diaries collection hosted by SLNSW segmented as the actual handwritten pages of the diaries in 11944 pages.

Narrative

This project presents part of the collection "Word War I Diaries" from the State Library of New South Wales (Australia). The artefact is an interface to that collection that learns from its ancestors: Visference and Crossreads. From

Visference it takes topic model analysis and curated labelling. From Crossreads, it takes the narrative multiplicity and fragmented reading. Diggers Diaries is clearly developed in two stages. These two stages match the Versions I and II. Version I uses a smaller corpus of diaries, and the research is focused on the analysis and its consequences for the contents. In this first version, a topic model is used to classify each page; topics are manually grouped and labelled . The process is explained in detail in section 6.1.

Collaborators

No other collaborators.

Data process

- Data was downloaded from the State Library of New South Wales transcriptions site. This task required web spiders and scrapers since the API was not offering to download what was needed.
- The dataset is a collection of handwritten notebooks. For the segmentation of texts, the natural handwritten pages were taken. This time there was no need for a segmentation method. The size of each page is similar to the one already tested in Visference.
- The data was then adapted to a data model where text and metadata of the diaries were stored, as well as the outputs of the text analysis in JSON files freely accessible online.

Data analysis

In this artefact the strategy to analyse the collection goes back to Visference, and topic model analysis is revisited, this time, incorporating the experiences from previously described artefact development. Several different sets of topics were generated and compared through running the curated labelling process on each of the outputs. Test were done with 5, 10, 20, 30, 40, 50 and 100 topics. The number that created better topics and levels was 50. These 50 topics were reduced to 25 final labels grouped into 5 categories: personal, war, military life, travelling, and the accidental tourist. See the two-level labels in Figure 12

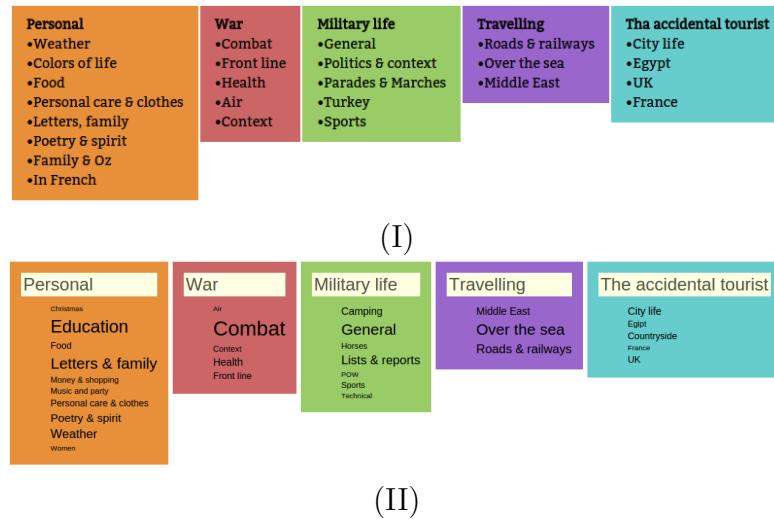


Figure 12. Two-level curated labels for Diggers Diaries versions I and II. Version II incorporates a tag cloud relative to the number of pages under each label.

Topic model analysis was conducted with the open source toolkit MALLET, a Java-based package for statistical natural language processing (McCallum, 2002)

Interface development

This interface is based in two layouts: the overview, and the reader.

The overview offers a list of the 126 diaries by authors' names, and for each diary every page is represented as a square. The squares are coloured in one of the five colour groups when the top scored topic is over the threshold. The threshold is set with a dropdown menu, its default value is 20%. The pages can be filtered through topic groups, coloured menus and submenus.

There are two ways of interaction and navigation: clicking a page takes the user to the reader. Clicking a label from the two-level menu filters the pages related to that label.

The reader shows the text of the current page, the metadata of the diary (author, title, dates), a topic chart for the page, the image of the cover of the diary, and the navigation options (previous/next page) in the diary.

Project outcomes

- The demo site: <http://diggersdiaries.org/1>
- The datasets generated after the collection and transformation of the texts is accessible for reuse as JSON files. All the links are in diggersdiaries.org

5.5. Diggers Diaries II: WWI Diaries

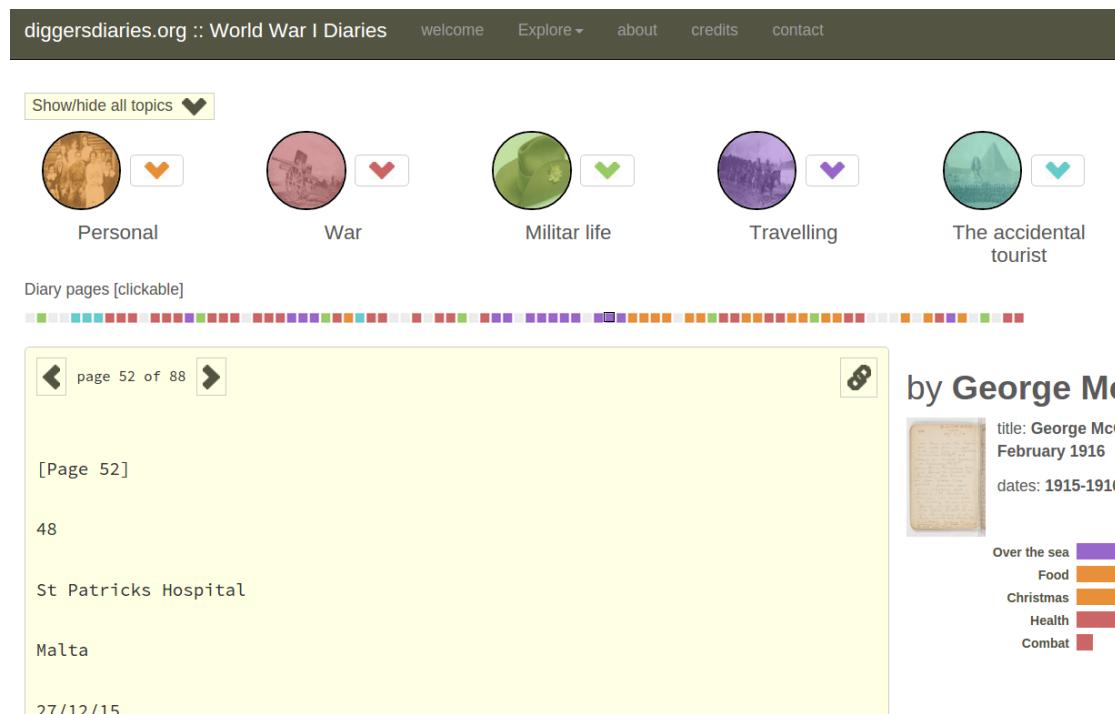


Figure 13. Screenshot of the reader of Diggers Diaries II interface.

Artefact name and description

Diggers Diaries II, an interface to explore crossreading of large collections of documents through topics.

Datasets

685 diaries from WWI as part of the WWI Diaries collection hosted by the State Library of New South Wales (SLNSW), segmented as the actual handwritten pages of the diaries in 81763 pages.

Narrative

As in version I, Diggers Diaries II focuses on the interface development. Since in version II the collection is about seven times bigger, the challenges need to be redefined in relation to version I. The number of diaries (688) and especially the number of pages (81763) brought constraints to the interface design. The pages

are shown as coloured squares that are grouped into images, one image for each diary. When a page is clicked, the reader rolls up, and a bar chart shows the five most relevant topics for that page..

Clicking a new topic, the system jumps to a new page in the chosen topic. This way, Diggers Diaries includes crossreading, through a two-level category menu made with topic model analysis and a labelling process. The demo site is just showing the crossreading idea applied to the collection of WWI Diaries; a production version would include more features to save pages to your account, to build playlists, and so on.

Collaborators

No other collaborators.

Data process

The steps for the data process are similar to the ones described in previous artefact. One problem we encountered was about the API and the URLs of the existing page. During the timebetween version I and II, many URLs changed due to updates, and changes in the SLNSW systeml.

It is remarkable that we downloaded not only all the page transcriptions, but their scanned images too. The 81763 images weighs 45GB. The rest of the data and code weighs 777MB. From the interface only is downloaded the queried data.

Data analysis

The main difference: we used 100 topics. After removing and joining topics initially generated with the software Mallet, we curated them under 30 labels. And again group these under the same five groups created in version I.

Interface development This version II needed a redesign of the interface due to the large dataset and the large number of pages represented.

In this version we keep the “all roads lead to Rome” philosophy from version I, that is, exploration for reading.

There are three overviews:

- Faceted: Diaries/authors (685 diaries)

- Timeline: Dates/duration (685 lines)
- Pixel/matrix: Diaries/pages (81763 pages)

The structure is reader oriented because the reader is accessible from one click in all the three overviews. On click, the reader is embedded into the overview.

The reader is also accessible directly. The reader and the three overviews are offered together as the exploration menu. The reader offers two ways of reading:

- Linear reading: natural reading of diaries, page by page.
- Crossreading: jumping page to page according to the two levels of labels.
When a label is clicked the reader brings a new page that has the label clicked as one of its main topics.

In order to represent all the pages with at least click interactivity, we created static images for each diary, and we map click locations to get which page of which diary is clicked. This image-grids technique allows for efficient interaction with a large amount of items in a standard HTML page.

Project outcomes

- The site in production: <http://diggersdiaries.org>
- Diggers Diaries was presented (at the University of Canberra, November 2015) to experts from the State Library of New South Wales (SLNSW) and the National Archive of Australia (NAA) and with a very positive feedback

6. Discussion

This chapter describes the lessons learned from the creation of the collection interfaces described in the previous chapter and from the literature and practice outlined in Chapter 2. It offers a reflection on the making process and working hands on with the data and in relation to existing concepts, practices and conventions in the field. These lessons and the research outputs have revealed concepts and methodologies to create interfaces to text collections oriented to support reading, and the exploration of the documents within digital collections. In this dissertation, such interfaces are introduced as deep interfaces.

The chapter starts with three sections composed, respectively, of three inter-related components: Text analysis, Text reading, and Text collection Interfaces. These three components have developed following a reflexive methodology in order to create a group of related artefacts. The final section defines the concept of deep interfaces, which integrates the lessons learned through these three components.

The first component discusses how the "miraculous" effectiveness of text analysis has been applied to the creation of the interfaces presented here. Despite the incredible advances of text analysis, the results of applying analysis techniques , to improve the interfaces to a collection of texts are not widely used, as shown in Chapter 2. This section shows the creative process including the use of topic model analysis to classify documents of a collection. The importance of topic model labelling in the integration of analysis output as elements of the interface is also shown.

The second component considers the task of reading at a time when it is difficult to find even a moment for reading. How can an interface to a large text document collection take this into account and offer an interface that is dedicated and oriented to support reading? How can we read, or at least have a taste, of such vast collections of texts? Is there a way to discover and explore a text collection

mainly through reading the texts rather than through viewing representations of the collection as a whole? As seen in Chapter 2, several authors have experimented with fragmented narratives and data multiplicity, that is, the multiple narrative lines that a text can contain. This idea of fragmentation, in combination with the nonlinear reading of crossreading, is presented as a strategy to deal with large text collections.

The third component, "The interface and its codes", discusses the lessons learned from building the interfaces in relation to conceptual interface design. The interface style that has been developed seeks simplicity in the sense that new users of the interface do not need previous explanation in order to start using the interface.

The final section, "Deep interfaces", draws all these findings together. It proposes deep interfaces as a conceptual umbrella that advocates for specific methods and strategies when designing interfaces to large text collections.

In other words, this chapter gives answers to the research question (see Section 3.2):

How might we create interfaces to text document collections that let us explore the collection by reading its contents?

The short answer to this question is creating deep interfaces. A deep interface entails a set of strategies that make the interface-data binary deeper in meaning, and richer in data structures and relationships. The discussed strategies include: simplicity of the interface, the use of text analysis outputs integrated in the interface as visual elements, and a way to read texts that are too long to read in their entirety.

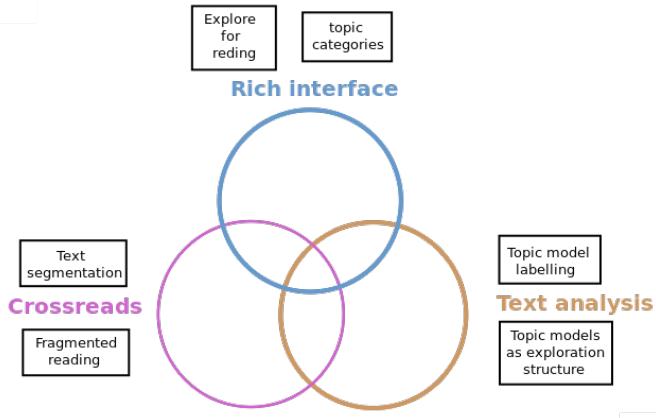


Figure 14. The three components of deep interfaces

6.1. The miraculous analysis

Visference: topic models and topic models labelling.

In the development of Visference — a sortable table for the listing of scientific conference papers — a topic model analysis was applied to a collection of 282 papers published in the JMLR proceedings (Journal of Machine Learning Research). This analysis was done in two phases. In the first phase a kind of golden list of ten topics was created, in order to represent key concepts and fields of enquiry in the field of machine learning. In the second phase, a collection of conference papers was classified according to this golden list.

In more detail, in the first phase a topic model analysis was conducted on a large set of scientific texts related to machine learning and statistics. Over ten thousand papers and about twenty books were included in the corpus; the papers were downloaded from the open access repository, Arxiv.org, and the books were hand-selected by experts (researchers from the Machine Learning Research Group, at NICTA, Australia). This corpus was used as a "golden" or authoritative dataset of texts representing the machine learning discipline. From the topic model analysis of this corpus ten main topics were found. The idea behind this analysis was to find the ten main topics in the field of machine learning research. The ten topics were labelled by experts. The final label election as well as the list of

books in the corpus is presented in Section 5.1 Visference. In the second phase of the process, the scientific papers from the “JMLR Workshop and Conference Proceedings” were segmented by sections and classified according to the ten golden topics. Since these topics aim to define the main categories that represent all content related to machine learning, the goal of classifying papers by topic was to help researchers interested in the conference to easily find papers that could be in their areas of interest. Using the trained model, that is the big corpus mentioned, a topic inference was conducted for the almost three-hundred papers to analyse. Each paper was segmented by section, and the final topic model that represented the paper was an average of the topics contained in their sections. Visference shows two aspects in relation to text analysis: on the one hand it shows that there is a significant opportunity to improve standard lists of papers in academic conferences. It seems surprising that scientific conferences, especially the ones about data analysis, still use very traditional, and non-interactive layouts for long lists of documents. On the other hand, more specific to the kind of analysis conducted in Visference —topic models— it is not very appropriate to represent a whole scientific paper under ten topics. A scientific paper has a length that demands to classify by topics smaller segments of the document.

Crossreads: text segmentation, and text similarity

Crossreads shows a collection of articles in the context of an exhibition (see 5.2 and 5.3). Every exhibition has its particular conditions, but they have at least one aspect in common, which is the limited time that visitors have in front of a piece. Translated to a collection of texts, this time-dependency is even more relevant, since the time needed to read a text is quantifiable in a minimum value of words per minute. In order to solve this first question I opted to segment the texts into small pieces that could be read in around one minute each. After reading one of the pieces, a new related piece of text from the collection appears in front of the reader, and so on. This idea of breaking the unity of the collection texts is discussed in more detail in the following section 6.2 that addresses the criteria for text segmentation, and for determining the similarity of text segments.

The text segmentation is explained in the conference paper published after the

development of Crossreads (Nualart and Ferraro, 2014). Two segmentation approaches were tested, each with different benefits. In both versions of Crossreads each document was divided into segments, such that each segment consisted of one or more paragraphs. A segment length was about seven hundred characters in total, which equates to an average of one minute of reading for an adult (Williams, 1998). In Version I, segmentation was procedural, that is, using computers. In Version II, segmentation was performed manually. While the method used in Version I was fast and capable of processing large collections, the method applied in Version II allowed for greater quality segmentation according to a more complex understanding of context.

The reason for these two approaches is that the segmentation task is very subjective. A human expert could add a personal view to the segmentation (Version II). A machine produced segmentation (Version I) can accomplish this task in terms of the size of each segment, but it cannot be expected to interpret the content of the text the way an expert human reader would. Both methods are compared here as part of the process of practice-led experimentation.

To develop the similarity calculus between segments necessary to create the Crossreads network, the following off-the-shelf Natural Language Processing tools and techniques were used:

- Tokenization: words in the segments separated by whitespace and punctuation characters.
- Stop word removal: standard stop word removal.
- Named Entity Recognition: identification and classification of Named Entities in each segment. The OpenNLP Named Entity recogniser was applied [2], which distilled four types of entities: Person, Location, Organization and Others.
- Similarity Calculus between segments.

The similarity between pairs of segments was calculated as the sum of the following factors,

$$Sim(i, j) = \text{TokSim} + \text{EntitySim} + \text{NESim}/3$$

where TokSim is the token cosine similarity between segments, a common vector based similarity measure. To calculate the similarity, the tokens of each segment are transformed into vectors and then the Euclidean cosine is used to determine the similarity between pairs of vectors. EntitySim is the sum of the Named Entities in each segment, normalised by the number of tokens in both segments, and NESim is the cosine similarity between Named Entities. During this process, the similarity between different NE types (Person, Location, Organization) is calculated separately and its average (/3) is calculated. The path among similar segments was calculated as follows. First, an arbitrary segment i was chosen, and used to calculate the similarity between the segment i and the entire segment collection. Second, the segment with the highest similarity value score was set as the maximum similar segment of $i..$. Since linear reading of a document is enabled in each iteration it was decided to skip links to segments of the same document as segment i . Finally, different constraints to each version were applied as follows:

- Version I: In the Crossreads network, each segment is linked to its most similar segment. The drawback of this approach is that links will have a wide range of similarity scores since, in each iteration, the number of segments to compare with is smaller, and the possibility of finding a segment with a high similarity score decreases. However, the benefit is that there will not be any orphan segments, i.e. all segments link to other segments, so the reader will always have the possibility of some crossreading.
- Version II: Each segment is linked to its most similar pairing. To avoid repetition of pairs, segments that have already been set as a maximum similarity segment during ten iterations are skipped. After ten iterations, the skipped segments are used again in the similarity calculus.

Crossreads shows two important lessons that have been applied to the final artefact. On the one hand a positive aspect is that crossreading is feasible. In the semi-structured interviews the experts saw great possibilities behind the idea of crossreading texts to develop in the future. On the other hand, related to text analysis, the experience shows that syntactic similarity is not as strong as semantic similarity. This is the reason that comparing Visference and Crossreads

experiences, in the final development of Diggers Diaries the chosen analysis was semantic, that is topic model analysis, as shown in the following paragraphs.

Diggers Diaries: natural text segmentation and topic model labelling

Diggers Diaries is the last artefact that, following the reflexive methodology, incorporates the lessons learned in creating the previous artefacts. Diggers Diaries is an interface for a collection of letters and diaries from World War I soldiers from Australia, digitised and held by the State Library of New South Wales (see 5.4 and 5.5). The collection of diaries contains over eighty thousand pages, grouped into over seven hundred volumes and envelopes. The diaries have been analysed using topic models with specific methods for segmentation and topic labelling.

The nature of the collection, that is, handwritten pages and some printed pages too, brought a genuine segmentation that accomplished the conditions pointed out in previous artefacts Visference, and Crossreads, in terms of text length, as well as content fragmentation. A handwritten page, sized in most of the diaries as an A5 page, can be read in about a minute. The reader has the scanned image of the original page, providing additional context, and the controls allow navigation to the previous and next pages easily. The segmentation here is intrinsic to the documents: we are comfortable with reading in pages, and with moving from page to page. Diggers Diaries provides this functionality in the interface, enabling the segments to be joined back together. Thus the page presents a useful fragmentation, the right size but also "readable" and readily joined because it corresponds to pages.

Initially the topic model analysis was tested with various numbers of topics. As explained in Chapter 2, the algorithm used to calculate topic models upon a collection of texts allows setting parameters for the number of topics to be generated. In Visference this parameter was set to ten, here, in Diggers Diaries, the analysis was done to obtain, ten, thirty, fifty, and one hundred. The list of one hundred topics was selected as the base to work with, because it shows more variability of topics. The second step was to classify each topic as either removable, well-defined or a synonym of a well-defined topic. The topic classification involved a subjective process of reading and interpretation to validate and analyse the topics.

The texts (pages of the diaries) were read to validate the topics (following the example of Blevins (2010)). This shows how reading and subjective interpretation is involved in the process of constructing the topics. The removable topics were those representing an irrelevant topic, e.g. "Monday, Tuesday, Wednesday, ...", or non-sense, e.g. "time days leave day good back months home france england week work long weeks great place ago australia camp things", or topics that related to transcribers notes, e.g., "diary letter written pages australian note page letters notes book copy war printed august records Australia transcriber's paper april read". Since the focus of this interface was to show the collection for library visitors, in this version, these notes were not included, although, for other more specific reasons, they could be added as part of the accepted topics. Thirty-eight topics were directly removed. From the rest, some topics were found to be synonymous or almost-synonomous, e.g. these two topics were joined with the label "Over the sea": "port board sydney melbourne ship boat left wharf leave ashore arrived bay harbour p.m troops fremantle town sea cape colombo", and "boat ship boats ashore board harbour water wharf aboard ships shore men port alongside small deck beach side troops coal". Most of the remaining topics were considered synonymous with others, so finally thirty topics were formed.

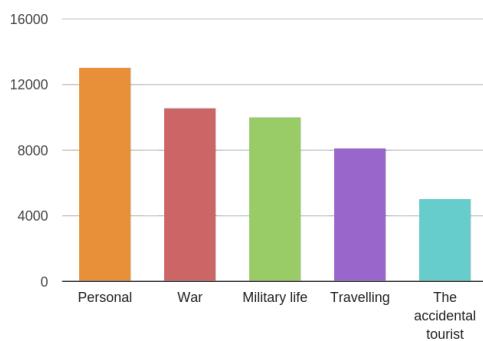


Figure 15. Most common topics by group for the collection of WWI diaries,

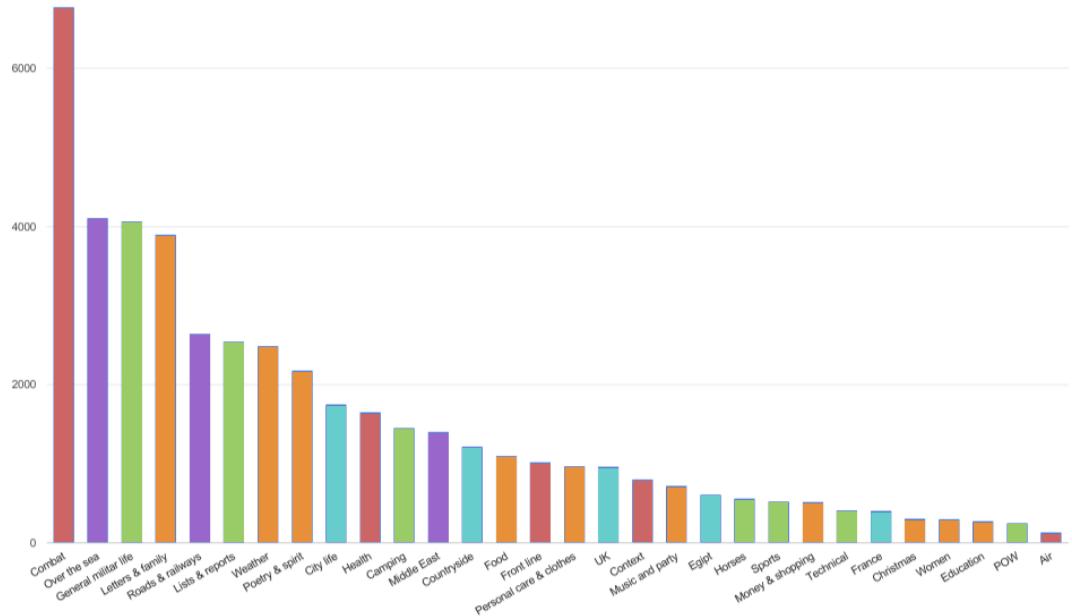


Figure 16. Most common topics for the collection of WWI diaries,

The process of labelling these thirty topics was undertaken in two phases: firstly, the topics were grouped into five general groups. The labels for these groups aimed to be simple and direct: Personal, War, Military life, Travelling, and The accidental tourist. This last label could be simply Tourism, but the chosen label gave a better definition of this topic, which comprises pages where the soldiers talk about impressions they had when visiting European capitals (Paris, London) and Egypt. The soldiers, many of them young men coming from Australian countryside, had significant experiences in visiting these places, and they shared these with their families through their families through the letters and diaries. "The accidental tourist" also refers to the 1985 movie of the same name (Tyle, 1985). This shows that labelling does not happen in a cultural vacuum, labels can include jokes, references, etc. Finally, for each topic some labels were proposed trying to be descriptive. The final groups and topics are shown in Table 3. The most frequent topic groups and topics are represented in two bar charts in Figures 15 and 16.

Groups of topics	Topic labels
Personal	Christmas, Education, Food, Letters & family, Money & shopping, Music and party, Personal care & clothes, Poetry & spirit, Weather Women
War	Air, Combat, Context, Health, Front line
Military Life	Camping, General, Horses, Lists & reports, POW, Sports, Technical
Travelling	Middle East, Over the sea, Roads & railways
The Accidental	City life, Egypt, Countryside, France, UK
Tourist	

Table 3. Groups of topics and topics labels for the analysis of a collection of diaries, as part of the project Diggers Diaries

The process of labelling described above was developed through three artefacts: Visference, Diggers Diaries I, and Diggers Diaries II. In contrast to the described labelling process in Diggers Diaries I and II, in Visference the labelling process was simple: only ten topics were created, and these were labelled through the opinion of several experts. The results listed in table 3 refer to Diggers Diaries II. This is the last artefact, usually referred as Diggers Diaries. Along the three steps the datasets used have grown along with the number of topics calculated.

Entity extraction combined with text similarity used in Crossreads, brought syntactic text similarity. In opposition, topic model analysis does not reveal anything about the syntactic structure of the texts, but about its semantics. in Crossreads, the entity recognition analysis needs an external source (Baldridge, 2005). In contrast, the topic model analysis, conducted in Diggers Diaries, is purely based on statistics. There are no external sources that bring any semantic context; the analysis is based only on counting words in segments of text.

The whole topic model analysis conducted in the final artefact depends on number of topics, segmentation unit, segmentation method, segment length —that impactson the reader, but also on the analysis, because a longer text has more topics— and topic labelling, which includes the curated tasks: validate, join, group, label groups, and label final topic models. See Table 4 for a comparatison of some of the parameters.

<i>Parameter \ artefact</i>	<i>Visference</i>	<i>Crossreads I</i>	<i>Crossreads II</i>	<i>Diggers Diaries I</i>	<i>Diggers Diaries II</i>
<i>Dataset size</i>	<i>papers ~1800 paper sections</i>	<i>700 pieces of text</i>	<i>1170 segments of text</i>	<i>126 diaries 11944 pages</i>	<i>688 diaries 81763 pages</i>
<i>Analysis</i>	<i>Topic model</i>	<i>Entity extraction + text similarity</i>		<i>Topic model</i>	
<i>Number of topics: initial / final</i>	<i>10 / 10</i>	<i>N/A</i>		<i>50 / 25</i>	<i>100 / 30</i>
<i>Segmentation unit</i>	<i>A section of a scientific paper</i>	<i>One or more paragraph of a text</i>		<i>One or more paragraph of a text</i>	<i>A page of a diary or a page of a letter (most handwritten)</i>
<i>Segmentation method</i>	<i>Natural</i>	<i>Procedural</i>	<i>Curated</i>	<i>Natural segmentation (pages)</i>	
<i>Segment length</i>	<i>Variable, but longer than one minute</i>	<i>All segments are close to 700 words</i>	<i>Variable.</i>	<i>Depending on curation</i>	
<i>Labelling method</i>	<i>Curated simple</i>	<i>N/A</i>		<i>Curated advanced</i>	

Table 4. Comparison of some parameters of the text analysis conducted in this research project.

Through the reflexive methodology, the process of text analysis and data enrichment from the analysis was improved across the development of artefacts. The kind of analysis was set to topic analysis that is, a semantic analysis. The preparation of the text to analyse was varied till Diggers Diaries, where the size of the text segments was more appropriate. Finally, the topic model labelling was improved from the very simple process described for Visference to the more elaborate process of topic selection, removal and joining of topics conducted in the two versions of Diggers Diaries.

6.2. The fascination for the data

With the aim of studying ways to deal with the enormous volume of texts available in DL, several strategies have been reviewed: from better and deeper analysis that

can bring more knowledge about the texts, to improvements in the design of the interface to those texts. These two fields — text analysis and interface design — are analysed in this dissertation since they seem to be key to improving our relationship with text collections. However, there is one more field related to this task: reading. Since it is impossible to read all texts available in collections, even the texts of one big collection, the reading task can be revisited and reinterpreted.

Traditionally, reading is a linear task that goes from the beginning to the end of a text. This approach is simply not applicable when the amount of text is too great. One possible solution to this problem is to divide the text into pieces, that is text segmentation. Of course, when a text is segmented, in reality what we break is a narrative that was written to be read altogether, in a row. Nevertheless, in these first decades of the digital age, we find multiple examples of the fragmentation of narrative texts. We find narrative fragmentation in the number of short texts we are exposed every day through short-text messaging (SMS, Whatsapp-like, etc.), microblogging (Twitter, Tumblr, etc.), and social networks (Facebook, Google Plus, etc.).

Reading and fragmented reality

Three aspects of the digital age are relevant here: our growing capacity to communicate in fragments, the scarce time we have to read, and the impossible amount of available texts. The question is: could the task of reading be approached as a fragmented task? Several works in the past have explored the possibilities of breaking the linearity of a text. The philosophers Deleuze and Guattari have described the rhizomatic structure of knowledge: "In a book, as in all things, there are lines of articulation, segmentarity, strata and territories; but also lines of flight, movement, deterritorialization and destratification". In the novel Hopscotch by J. Cortázar (Cortázar, 1966), the author proposes two reading orders for the chapters; the text starts with: "In its own way this book is many books, but mostly it's two books". Project Xanadu from 1960 (Ted Nelson, et al, 1960) is considered the first hypertext project in the digital era, and it was a visionary definition of standards for the WWW that were mostly not included in the standard protocols. One of the Xanadu's features is transclusion, defined by its author, Ted Nelson, as

"the same content knowably in more than one place". One of Xanadu's seventeen rules states: "Every document can consist of any number of parts each of which may be of any data type".

These examples bring confidence and evidence that the traditional idea of reading as a linear task, can be developed to more complex options: multiple narratives. This idea is the base for the development of Crossreads, which is analysed in the following paragraphs.

Crossreading: the project and its process

The first version of Crossreads (see Section 5.2) was developed with the idea of text segmentation in mind as a possible strategy to read long texts. Crossreads was developed to visualize, present, and interface a collection in the context of the museum exhibition ("The Listening Eye" MACBA 2014 (Barcelona, Catalonia) (MACBA, 2014). This is a collection of fifty-seven articles by Eugeni Bonet, an artist working in cinema and video art since the 1970s, mainly in Catalonia and around Europe.

The peculiarities of this collection include:

- One author: all articles have the same author
- The topics are in the range of video, cinema, and art, since most of the articles were originally published in program brochures of video sessions
- The range of published dates is thirty-eight years: the techniques, hence, the language of video technology has changed significantly over this period. Also the point of view and discourse has changed for the past forty years.

Peculiarities of the exhibition context includes:

- Visitors have limited time to engage with the texts
- Use of projections of the texts in the walls
- Use of tablets for visitors with Crossreads interface to read and browse the collection

In Crossreads I, the collection is in Catalan and Spanish. In order to share the work with more people, a second version with texts in English was developed. In this second version a similar text was used. *In Your Computer* (Quaranta, 2011) is a collection of texts written by Domenico Quaranta between 2005 and 2010 for exhibition catalogs, printed magazines and online reviews. The book is published under Creative Commons licences that allow reuse. In this case the author was contacted as a courtesy and agreed to have his text used in this research project.

Crossreads is defined as a manner to deconstruct linear narrative text in order to read text in multiple orders. This project studies data multiplicity, and textual visualization interfaces. The preparation of the texts for crossreading, includes an initial segmentation into small blocks. Then, the textual similarity among the segments is calculated. Finally, a web interface allows exploration of the texts (Nualart and Ferraro, 2014).

From the moment that a long text is segmented into pieces, and these pieces are read in a different order according to, e.g, similarity, there is a different transfer of information for each chosen order of pieces. This fact seems formally logical since the list of segments and their reading order are different. Then formally, the initial linear data, the text creates an explosion of possible combination of elements. Somehow this is a phenomena of data multiplicity. To measure the difference among different sets of segments is beyond the scope of this research. In short, segmentation creates an exponentially increasing set of new permutations of segments.

After Crossreads I and II, a fact related to the length of the segments arose: some segments had no other similar segment in the collection. The reason for that seemed to be that the dataset was not big enough. A much longer corpus, with a much larger number of segments would create more possible similarities among segments. If the number of segments is small, segments with less common content will have few similar segments, or none at all.

Diggers Diaries: when crossreading becomes natural

Diggers Diaries is an interface to a part of the collection "Word War I Diaries" from the State Library of New South Wales (Australia). The collection items

are dated from 1914 to 1920. It contains texts from 337 authors (all ANZAC, Australian soldiers), 688 diaries, letters and military reports. So far, a total of 81763 pages, that is segments of text. All handwritten pages were transcribed to text by professionals and users of SLNSW.

This project proposes a way of exploring and reading the collection based on grouping pages according to the topics they talk about.

One of the aims of Diggers Diaries is to help reading text collections that you are not able to read for one or more of these reasons: time rush, text length, low accessibility and usability of the document (small font size or inappropriate colour palette), and legibility of the documents (low quality in the representation of the original document, and poor conservation of the original).

From the previous artefacts developments, a list of lessons learned have arisen. The first is related to the kind of analysis used to find relations between segments of text that belong to a collection. This question is discussed in 6.1, in short, Diggers Diaries uses a semantic analysis of topics instead of the similarity analysis conducted in Crossreads. To solve the problem of low diversity in text segments due to the small number of segments, a much bigger collection of texts was chosen. Crossreads II had 1115 segments of text, while Diggers Diaries has over eighty thousand segments.

A notable difference with Crosreads version is the segmentation process. In Diggers Diaries the segmentation is natural: since the texts are presented by pages, each segment is a page. There is a wide range of text length in the collection of pages. The text is a page is accompanied by the original image of the page what seems to reinforce the logic of the segmentation in pages. As in Crossreads, in Diggers Diaries the reader can stay in the current diary and go back and forwards, or jump to another page. The similar life situation of the authors (most young men from Australia travelling to the other side of the globe to fight in WWI) generates a common narrative landscape that smooths crossreading.

6.3. The interface and its codes

The interface of a digital collection, is not only the tool we have to interact with a digital collection, it is also all we see of the collection. From the point of view

of an observer using the interface, the interface “is” the collection. In fact, the interface is a simplification of the collection that highlights some aspects of it, and can hide, in some contexts, its rich complexity. In those cases, the interface plays the role of the surface to the collection. To use the interface implies diving into it from the outside, towards the hidden treasures inside.

This section follows a chronological narrative that shows the cumulative process of interface development in the artefacts presented. It starts with Visference (see Section 5.1), which introduces new features in text collection interfaces using standard web interface elements. Then Crossreads introduces crossreading as a practice for reading across a text collection. Finally Diggers Diaries incorporates all the lessons learned from previous artefacts, and includes some new features. It is a text collection interface oriented to support and suggest reading.

Visference: conservative design

Visference is an interface proposed for an academic conference. Academic conferences commonly list accepted papers in a flat, text-based, non-interactive web page or a printed program. The dataset chosen for Visference was a collection of 282 accepted papers from the JMLR Workshop and Conference Proceedings Volume 2 : Proceedings of the 30th ICML (ICML, 2013).

The motivation of Visference is to improve the way we list scientific papers in conference websites. A relevant work mentioned is ”Word storm”, by Castella and Sutto (Castella and Sutton, 2014), that lists accepted papers of a conference adding a word cloud with a modification of the algorithm that makes word clouds easily comparable (see 2.2.1).

Visference had a design compromise from the very beginning: there is much room to innovate in conference paper lists, but at the same time it is necessary to use web conventions that will improve the exploration of the list without the need for an initial tutorial on the interface. This philosophy has been applied to the rest of the interface design for all the created artefacts. The documents are presented in a HTML table with sortable columns. The columns are: paper title, authors, PDF link, abstract/reference link, and the ten topics. The sortable columns allow users to see the most relevant papers for each topic, as well as the related topics.

In terms of interface design, Visference is not innovative because of this low-risk design compromise. The innovation in Visference is the presentation of a list of documents that is traditionally flat, non interactive and not sortable. Visference tries to push for innovation in the conservative websites of scientific conferences, and, in the development of this research project, represents the philosophy that has been applied to the rest of the interfaces: the use of standard visual elements in the interfaces, as well as the avoidance of introductory tutorials to explain the use of an interface. The following paragraphs discuss the role of interfaces when the task of reading is a priority.

Crossreads: fragmented narratives

Crossreads as a proposal for reading text in fragments has been presented and discussed in 6.2. In relation to interface development, Crossreads is a key project of this research project. Crossreads is the first reading oriented interface. Version I was designed for an exhibition space of interaction. This condition had some design considerations, such as a responsive design adapted to the devices used in the exhibition space (10-inch tablets); a length of the exhibited texts adapted to a short time visit; and an interface design that uses standard interface elements as in Visference, because visitors do not have time to learn how to use an unfamiliar interface.

Another central part of Crossreads in relation to its interface is the design of the reader, that is the interface where each piece of text is presented in front of the visitor. The reader is placed in the centre of the screen, and highlighted by a border. The interactions that characterize crossreads are to jump to a similar page of the collection, and to jump to any other random page of the collection. The elements that can be considered "standard" are: next/previous page of the current document, and the timeline that allows browsing documents of the collections by kind of document. These two groups of elements are placed at the same visual level.

Diggers Diaries: a lot of details build an overview

Diggers Diaries is an interface to a part of the collection of diaries, letters and military reports from Australian soldiers —also called diggers— in World War I. The project was developed in two stages: version I and version II (see Sections 5.4 and 5.5). Both versions draw on the same collection of WWI diaries, but version I uses a smaller corpus of diaries. In version I the research focused on the analysis and its consequences in the contents of the collection. The development of Diggers Diaries version I in relation to text analysis is discussed in detail in section 5.4. In version II the dataset increases by over six fold in the number of segments of texts, therefore version II development focused on the interface because it uses the same kind of analysis process developed, for the first time, in version I (topic model analysis, and manual topic grouping and labelling). Version II focuses on elements that eventually will help visitors in reading and exploring the collection. Among these elements there are data visualization elements.

Three data visualization elements are the tools to explore and overview the collections: by pages, by diaries, and by date:

- By-pages overview is a page-grid visualization device that represents the smallest unit of content (page) at the whole collection level. It draws on “show everything” type interfaces, but again the innovation is in using analysed text content to fill in this view. This by-page overview interface offers several display modes, as well as a colouring scheme. The grid of pages is coloured according to each of the five groups of topics. The pages can be filtered by group of topics. Since each topics score is a percentage of that topic for that page, the pages are coloured according to the biggest topic score. Only scores bigger than 20% are coloured.
- By-diaries overview is a faceted view of all diaries that can be sorted by date, and alphabetically by author names and topic names. When a diary is clicked, the diary metadata is shown, as well as a page browser shows the principal group of topics for each page according to five colours. Then, when page is clicked, the reader appears integrated in this diary view. Other diaries can be expanded at the same time, so diaries can be read in parallel.

- By-date overview shows the start and end dates of each diary, and, in consequence, also the duration. This overview allows selection of contents according to the historical moments of the collection (e.g. the reader can easily select texts from the beginning of the war, from the end, etc.). Diggers Diaries follows the model that gives several options to explore the collection from the home page, as in "Explore Australian Prints + Printmaking" (Ennis and Whitelaw, 2014), and the Eugenics Archives (Collective, 2016). See 2.3.3 for more details. In the home page the reader options are available and, at the same visual level, the three mentioned exploration options. The reader takes the visitor to a page of the collection and invites to crossread the collection jumping to pages according to a two level menu of topics (see Section 6.1).

The reader interface in Diggers Diaries introduces a topic to develop in future research projects, that is the inside out exploration of textual document collections, and the idea that the construction of an overview of a collection of texts, hence, of its contents and context, can be created from samples of the collection, and not necessarily from overviews—or distant reading—of the collection.

This section has shown how elements and concepts related to text collection interfaces developed across the creation of artefacts with real collections. From Visference, and its compromise to create a simple interface based on text analysis , to Crossreads and its fragmented narratives as a strategy to read unredably large collections of texts, as well as a way to build multiple narratives upon one text. Finally, Diggers Diaries, presents the first deep interface at the same level standard visual elements used in interfaces with elements that come directly from text analysis outputs, That is, buttons that select pages within the collection that semantically corresponds to topics. The next paragraphs introduce the concept of deep interfaces as the results of lessons learned in the experience of developing the presented artefacts. In the following section this concept is defined and discussed.

6.4. Deep interfaces

This section presents the diversity of findings that this research project has produced and integrated into an umbrella concept: deep interfaces. The construction of the definition of deep interfaces is a cumulative process through the experience of developing the presented artefacts (see Chapter 5).

Deep interfaces is a wide concept used to compile the experiences of developing interfaces to textual document collections. Deep interfaces refers to interfaces that on their surface resemble a standard web page, although they allow deeper interpretation of the meaning, and the structure of the contents of the collection due to their integration of text analysis elements into the collection interface

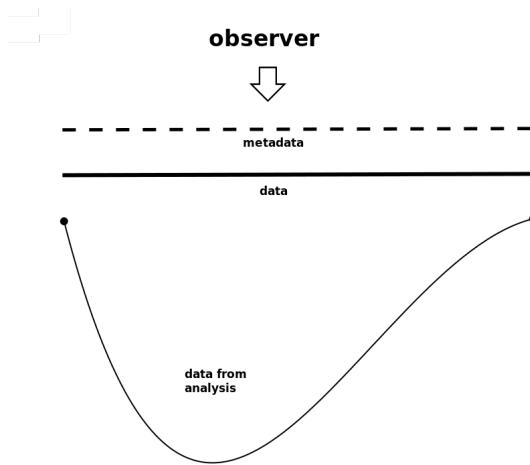


Figure 17. Deep interfaces from the point of view of data:: metadata, data, and data from analysis.

Whilst text analysis is widely applied in more fields every day, there is a gap in its use for digital collection interfaces. In the digital humanities , text analysis is used for so-called “distant reading”; most of the reviewed examples are not dedicated to improve the reading of text collections, but to overview and explore the collectionore the collection (see Section 2.3). Deep interfaces fill this gap, incorporating text analysis into the interface to text collections as visual elements.

Text analysis of every text of the collection adds a new layer of metadata that enriches the information structure of the collection, usually represented by standard metadata fields. New structures within the documents of collections offer

new relationships among collection items, and, therefore, new meanings and ways to interpret the collection. In this context, "deep" is about the volume of data exposed at the interface. Text analysis enables document content to be represented at the interface.

In order to get to know big textual collections —those too big to be fully read—a strategy has been practised during this research project, that is, crossreading. In short, to cut the documents of the collection into pieces according to some criteria, and then read part of the collection with the help of structures among those pieces. Crossreading is presented as a fragmented narrative reader. According to this idea, collections of texts can be overviewed reading small parts of them. That is, a reader can construct her/his own very subjective idea of the collection upon details of the collection. Since to generate an overview of the collection is not the goal of crossreading, in a sense, to crossread the collection will eventually generate a view —or opinion— of the collection.

Crossreading is a concept that complements the idea of building interfaces for reading, or oriented to read, a text collection. In a reading interface it is not necessary to explore the collection, nor to overview it, before starting to read texts of the collection. Deep interfaces are oriented to take the visitor of the collection directly to an item of the collection, no matter how big the collection. Diving into the collection this way, and then jumping to other items of the collection according to topics or similarities among the items (or parts of the items); in this way deep interfaces can be explored. In that sense, the texts of the collection are accessible from the home page, or at the distance of a single click. At the same time this proposed system or strategy to build interfaces to document collections is not incompatible with other standard techniques for exploration, like the use of data visualization techniques to offer the items of the collection according to some standard metadata fields, such as: timelines, geo-maps, faceted lists by authors, etc.

This practice led research has developed several interfaces to textual document collections with the technical aim of bringing together textual collections with the power of text analysis. Its broader goal is to generate ideas to bring to life knowledge that lies stored in institutions in the form of textual document collections. Since this research project methodology follows a reflexive—cumulative—process,

the final compilation of lessons learned and implemented is presented under the label of deep interfaces. The idea of deep interfaces is a concept that has three legs: the way we analyse textual documents and use its results, the way we read long texts, and the way we build interfaces for textual collections.

The following last chapter recaps the whole research project and points out ideas to develop in the future related to interfaces to digital collections.

7. Conclusions and future

This practice-led research project has produced accessible interfaces to real textual document collections that offer new ways to interact with these collections. Drawing on the experience of creating those interfaces, it introduces the concept of deep interfaces. This approach compiles the findings of this research into a set of recommendations for building new interfaces to text collections, where features are derived from state-of-the-art text analysis, and dedicated to support the reading of texts within the collections.

This dissertation has explained the rationale, the results, and the overall experience of creating interfaces to text collections. As the introduction showed the context, motivation, and aims for this research are linked to the growing amount of text available —some coming from a digitization process, some digital born—and the urgency for better means to interact with that quantity of information and its posthuman scale.

A review of current practices in major public digital libraries revealed consistent features, such as the aggregation of contents from other institutions, the dominance of traditional text-based interfaces, and sometimes limited access to the objects of the collections. This review also considered cases that offer innovative ways to reuse their collections, like New York Public Library that offers a snapshot of the public contents to download and reuse. As a general view, while institutions are beginning to innovate through these created “labs” —departments of the institutions that experiment with its contents— the official sites remain very conventional.

A brief review of the field of text analysis introduced its standard methods for finding similarities among texts, such as document clustering, topic model analysis, and related work in topic labelling and evaluation. While text analysis techniques like topic models have been applied to text collections in DH, they are not used to support interfaces and reading, rather to do “distant reading” analysis. During

the research project most of these techniques were used to manipulate the original collections of texts, and to create interfaces to interact with them. This project ultimately uses semantic similarities calculated with topic model analysis, and supported by a curated process of topic clustering and labelling which makes the statistical analysis human-readable and usable.

The third literature review considered the development of interfaces to digital libraries in general, with a specific focus on text collections. It reviewed the theoretical framework as well as new paradigms and metaphors proposed by various authors that have influenced interface development for the last twenty years, starting with the popular “Overview first, zoom and filter, then details-on-demand” and “Previews and Overview” by Schneiderman and their collaborators, moving on to “Show everything” by Stamen, and, finally, more recent proposals of “Information Flaneur” by Dörk, and “Generous interfaces” by Whitelaw. These metaphors brought fresh air to the design of new interfaces, and new features such as the opportunities beyond the search box, and support for serendipity in front of task oriented design. These sort of features contribute to the freedom of the user when interacting with the data represented in the interface.

With this idea of innovation with simplicity, a review of standard practices in interfaces to text collections and the use of data visualization as interfaces to these collections showed that most innovative proposals come from practitioners and researchers that work outside of major institutions, and from work in data visualisation rather than digital libraries. Finally a brief review of the concepts of reading texts on screens, that is e-reading, found that the standards for e-readers established three decades ago, today have become standards and, therefore, are included as common features in all kinds of digital readers.

From these multiple literature reviews several gaps and opportunities for new research were identified and introduced in Chapter 3. One key gap is the lack of the use of interfaces not based on a traditional search box in official sites of major public institutional DL. Innovation in these interfaces is found in projects from external practitioners and researchers. In a promising tendency institutions are beginning to include some of these external projects in their official sites, see for example the works of Whitelaw (<http://mtch1.net/>). Following this tendency, Crossreads I was developed outside of the institution, but included in its official

site too (see Section 5.2).

While computational text analysis has been applied in fields such as information retrieval for decades, and, more recently in literature, history and the digital humanities, text analysis is not used in current interfaces to text collections. As a contribution to knowledge, this research project uses text analysis outputs as visual elements integrated into large-scale collection interfaces. Finally the study of reading as a task shows that it is impossible to read text collections at the scale considered here. In response this project proposes a method to get a sample of an un-readably large collection of texts, that is, the concept of crossreading. The idea is to divide long texts into segments and through computer aided recommendation, jump from segment to segment according to topics or other relevant features.

In order to explain the process followed in these experiments and make them repeatable with other collections, Chapter 4 explained the project’s methodologies and methods. The main methodologies include a reflexive research mode, that makes the whole research a cumulative process building on the lessons learned and exported from artefact to artefact. In that sense all the project’s key advances have been applied to the last developed artefact, that is Diggers Diaries. The practical nature of this research, is where theory can arise from practice, as a practice-led and a practice-based research. About the methods and tools used during the research project, a list of software is described and links to all source code and data are provided in order to encourage other researchers and practitioners to copy and improve the methods presented here.

Chapter 5 described the results of the project in detail. Five practical projects were introduced: Visference, Crossreads (I and II), and Diggers Diaries (I and II). This chapter also describes the process of creation, including data gathering, transformation, and analysis, interface development —libraries and languages used— evaluation studies, and project outcomes. A brief narrative explains the experience and highlights issues related to each artefact.

Drawing on these artefacts, and the lessons learned from the creative experience the discussion chapter outlines a narrated timeline from three key points of view: text analysis, reading, and the interface. These components are the three ”legs” that support the project’s central proposal for deep interfaces. Firstly, the miraculous text analysis offers rich opportunities for experimenting with text analysis

outputs to be integrated into user interfaces to text collections. The second leg responds to the question of how we can encompass the posthuman scale of digital text collections now accessible online. This question is answered through the creation of Crossreads (I and II) artefacts and the concept of crossreading that the final project, Diggers Diaries, presents. It can be said that the melting pot of experiences during the practical creation of the artefacts has been used in Diggers Diaries (the final artefact) as all the lessons learned in a theoretical level, are used to define the concept of deep interfaces.

This chapter continues with a description, with examples, of the limitations, potential applications and future work prompted by this research. The chapter ends by briefly drawing together some thoughts to conclude this dissertation.

This project presents significant contributions to the creation of new interfaces to digital text collections. At the same time its limits should be acknowledged. Each interface works with a real collection of digital texts. Using real collections demonstrates the validity of these techniques, but it also means that the design of each of the interfaces responded to specific constraints and contexts of development. These constraints also introduce specific limitations related to the design and features of the interfaces.

While they offer important new approaches to text collection interfaces, the artefacts created in this research project do not innovate in interface design. In order to avoid the need for introductory tutorials to explain new features and interactions, the works presented are conservative in terms of design.

The features offered in the interfaces are quite limited. For example none of the interfaces offers a full text search box for the collection. In part this is done for pragmatic reasons, to avoid server-side queries, and propose a simple, client-side application. Also the sites that host the collections already offer full-text search and the focus of this research project is to propose new features, instead of reproducing already existing ones.

The concept of crossreading and the design of the reading experience are based on the personal subjective experience of working with the collections. To support further work on reading collections it will be necessary to know more about the reading process – to test and validate approaches such as crossreading.

Another limitation, especially for the case of Diggers Diaries, is one of data

synchronization —that the collection of diaries hosted by the State Library of New South Wales grows as new donated diaries are digitized. To add them to the Diggers Diaries interface would require the data analysis to be rerun and an updated dataset built. At the moment this is a process that one person could do in a matter of days. Depending on the volume of new data, the topic analysis, clustering and labelling would also need to be redone, which is a more significant task.

There are some general preconditions that future applications of the methods defined in this research project should take into consideration: data access and licensing. In projects where text analysis is done at a collection item level, full text access to every object of the collection is required. In these cases, accessing the data requires some technical expertise in order to use APIs or scrape the contents. The data needs to be under a license that allows reuse, modification, and publication.

In relation to the nature of the collections there is a question about the content constraints for the application of crossreading to a collection of texts. In this research project all the collections had some homogeneity. In the texts used in Crossreads I and II, each collection had a single author, although the texts were diverse in format and structure, including interviews, opinion pieces, presentations, catalog notes, etc. These formats impact on the segmentation process as part of the preparation of the collection for crossreading. In contrast to this, in Diggers Diaries, the collection of diaries is highly homogeneous. The personal situation of the young men far from their families on the front line, or at the opera in London, leads the authors to share similar feelings, experiences and situations that they express often in a similar style, in line with their common cultural background. The experiences with these collections, very different in both structure and content, is not enough to make conclusions. Certainly more research would be needed about the nature of texts suitable for crossreading, and the way we read digital texts.

Deep interfaces could be applied to a range of other collections. Inspired by “Mining the Dispatch” (Nelson, 2010a), the deep interface concept could be applicable to crossread press news, that is, jumping to other articles with similar subjects. Another field that could be tested is poetry, jumping across authors, styles and topics, whilst reading the poems.

One more subject to work in the near future in interfaces to text collections, is to revisit and expand the concept of digital reading. Digital reading can be understood as reading, listening, and/or viewing. Text documents can be complemented with images, audio and video from similarly accessible collections, so that these collections can be not only crossread, but cross-listened, and cross-viewed.

The software produced in the project has been designed to be applied easily to new collections. Technically, this project is completely described in order to assure reproducibility, and even more important, repeatability with other collections of texts. All the source code, data and other related documents are accessible in the respective repositories (Nualart, 2016). The techniques, and software necessary to repeat the work is free software, and the libraries and computer languages used are generic.

As a possible real application of the methods, at the moment there is a private proposal in discussion from a start-up about building a deep interface to legal texts for experts and professionals. This collection would include all the laws, and all the court sentences of a specific country. One of the challenges of this project is to work with a dynamic collection of texts, as is the case with court sentences. Probably the system would need to rebuild the model regularly including updated content.

This work shows the feasibility of creating deep interfaces to text collections, exposing the content of collections to improve their readability. A future challenge will be to implement such interfaces. Institutional collections could lead the way here, however there are several factors that might slow this development. Challenges for institutional collections include limited resources, a reliance on proprietary (vendor provided) DL software, and limited in-house technical capacity. Alternatively, deep interfaces might be implemented outside the institutions, like some of the examples already discussed. This option has many advantages, but also introduces difficulties associated with data synchronisation, sustainability, and data access: servers, databases and APIs are not alive forever.

At the moment of writing this paragraph Google has announced the release as free software of their "SyntaxNet" (Petrov, 2016), an open-source neural network framework that offers state-of-the-art syntactic analysis of natural language. This will empower everyone with access to a standard computer to analyse large

amounts of text with a state-of-the-art tool. This can be seen as a further step in the democratization of text analysis. Besides that, it is worth mentioning once more Mallet (McCallum, 2002), Machine Learning for Language Toolkit, free software that implements techniques for text classification, sequence tagging, and topic modelling.

Another avenue for future work relates to the advances expected in data analysis, and in so-called Natural Language Understanding. Deep learning techniques promise significant improvement in performance and in features (Unit, 2012). This will help the process of integrating data analysis and visualisation in education, and support the development of new ways of using data, new interactions, and, thus, new interfaces.

In short, the rationale of this research is as follows: we live in a sea of digital text —more than we can ever imagine. In fact, libraries, museums, archives are digitising and preserving all kinds of materials. The value of all these processes of preservation is dependent on accessibility: what use is all this if it is not read? Therefore, there is an urgency for new forms of access and interaction with these digital texts. This project faces this urgency, demonstrating the feasibility of deep interfaces for text collections.

Bibliography

Mortimer J Adler and Charles Van Doren. *How to read a book: The classic guide to intelligent reading.* Simon and Schuster, 1972.

Nikolaos Aletras, Timothy Baldwin, Jey Han Lau, and Mark Stevenson. Representing topics labels for exploring digital libraries. In *Proceedings of the 14th ACM/IEEE-CS Joint Conference on Digital Libraries*, pages 239–248. IEEE Press, 2014.

AAA American Anthropological Association. Cultural anthropology collections. http://www.culanth.org/curated_collections, 2016. (Visited on 01/31/2016).

Keith Andrews, Wolfgang Kienreich, Vedran Sabol, Jutta Becker, Georg Droschl, Frank Kappe, Michael Granitzer, Peter Auer, and Klaus Tochtermann. The infosky visual explorer: exploiting hierarchical structure and document similarities. *Information Visualization*, 1(3-4):166–181, 2002.

archive.org. Internet archive: Digital library of free books, movies, music & way-back machine. <https://archive.org/>. (Visited on 01/31/2016).

Jason Baldridge. The opennlp project. URL: <http://opennlp.apache.org/index.html>, (accessed 2 February 2012), 2005.

Laurent Baleydier. Wikipedia kartoo, 2001. URL <http://en.wikipedia.org/wiki/Kartoo>. [Online; accessed 12-November-2015].

et al. Bernhardt. Deutsche digitale bibliothek visualized. <http://infovis.fh-potsdam.de/ddb/>, 2015. (Visited on 02/01/2016).

- Howard Besser. *The past, present, and future of digital libraries*. Wiley Online Library, 2004.
- David Blei. Topic modeling and digital humanities. *Journal of Digital Humanities*, 2(1):8–11, 2012.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- Cameron Blevins. Topic modeling martha ballard s diary. *Pers. Blog*, 2010. (Visited on 12/12/2015).
- Bono. Data never sleeps 3.0. <https://www.domo.com/blog/2015/08/data-never-sleeps-3-0/>, 2015. URL <https://www.domo.com/blog/2015/08/data-never-sleeps-3-0/>. [Online; accessed 10-May-2016].
- Christine L. Borgman. What are digital libraries? competing visions. *Inf. Process. Manage.*, 35(3):227–243, 1999.
- Wray Buntine and Jaume Nualart. Icm 2013 simple topic model. <http://research.nualart.cat/visference/visference-topicmodels.html>, 2013. (Visited on 11/28/2015).
- Vctnnevar Bush. As we mov think. *Perspectives on the computer revolution*, page 49, 1989.
- L. Candy, S. Amitani, and Z. Bilda. Practice-led strategies for interactive art research. *CoDesign*, 2(4):209–223, December 2006. ISSN 1571-0882. doi: 10.1080/15710880601007994. URL <http://www.tandfonline.com/doi/abs/10.1080/15710880601007994>.
- Linda Candy and Cognition Studios. Practice based research : A guide practice and research. 2006.
- David Carter and Luiz Fernando Capretz. 3d user interface for a file management system. *IEEE Can. Rev*, 44:13–15, 2003.

- Quim Castella and Charles Sutton. Word storms: Multiples of word clouds for visual comparison of documents. In *Proceedings of the 23rd international conference on World wide web*, pages 665–676. ACM, 2014.
- Tim Causer and Valerie Wallace. Building a volunteer community: results and findings from transcribe bentham. *Digital Humanities Quarterly*, 6, 2012.
- Matthew Chalmers and Paul Chitson. Bead: Explorations in information visualization. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 330–337. ACM, 1992.
- Jason Chuang, Christopher D Manning, and Jeffrey Heer. Termite: Visualization techniques for assessing textual topic models. In *Proceedings of the International Working Conference on Advanced Visual Interfaces*, pages 74–77. ACM, 2012.
- Cisco. The cisco® visual networking index, 2016. URL <http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/mobile-white-paper-c11-520862.html>. [Online; accessed 10-May-2016].
- Jeff Clark. Spot, 2009. URL <http://www.neoformix.com/2012/IntroducingSpot.html>.
- Jeff Clark. Grimm s fairy tale metrics, 2013. URL <http://www.neoformix.com/2012/IntroducingSpot.html>.
- Collective. The eugenics archives. <http://eugenicsarchive.ca/>, 2016. (Visited on 08/04/2016).
- Julio Cortázar. Hopscotch (rayuela). *New York: Pantheon*, 1966.
- Fred D Davis, Richard P Bagozzi, and Paul R Warshaw. User acceptance of computer technology: a comparison of two theoretical models. *Management science*, 35(8):982–1003, 1989.
- DCMI. Home: Dublin core® metadata initiative (dcmi). <http://dublincore.org/>, 2001. (Visited on 12/22/2015).

- DDB. Startseite - deutsche digitale bibliothek. <https://www.deutsche-digitale-bibliothek.de/>, 2016. (Visited on 08/04/2016).
- Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391, 1990.
- G. Deleuze and F. Guattari. Introduction: rhizome. *A thousand plateaus: Capitalism and schizophrenia*, pages 3–25, 1987.
- Giorgio Maria Di Nunzio. Visualization and classification of documents: a new probabilistic model to automated text classification. *Bulletin of the IEEE Technical Committee on Digital Libraries (IEEE-TCDL)*, 2(2), 2006.
- Marian. Carpendale Doerk, Sheelagh, and Carey Williamson. The information flaneur: a fresh look at information seeking. In *Proceedings of the 2011 annual conference on Human factors in computing systems*, pages 1215–1224, 2011. URL <http://dl.acm.org/citation.cfm?id=1979124>.
- Johanna Drucker. Humanities approaches to graphical display. *Digital Humanities Quarterly*, 5(1):1–21, 2011.
- Erik Duval, Wayne Hodgins, Stuart Sutton, and Stuart L Weibel. Metadata principles and practicalities. *D-lib Magazine*, 8(4):16, 2002.
- Butler. Ennis and Mitchell Whitelaw. Australian prints + printmaking. <http://printsandprintmaking.gov.au/>, 2014. (Visited on 02/01/2016).
- Europeana. Europeana - homepage. <http://www.europeana.eu/portal/>, 2008. (Visited on 12/07/2015).
- Christiane Fellbaum. *WordNet*. Wiley Online Library, 1998.
- Gabriela Ferraro, Hanna Suominen, and Jaume Nualart. Segmentation of patent claims for improving their readability. *Proceedings of the 3rd Workshop on Predicting and Improving Text Readability for Target Reader Populations (PITR)@EACL 2014*, pages 66–73, 2014.

- Paula Findlen, Dan Edelstein, and Nicole Coleman. Mapping the republic of letters, 2011.
- Bryant Foo. Navigating the green book | nypl labs. <http://publicdomain.nypl.org/greenbook-map/>, 2013. (Visited on 01/31/2016).
- Trevor Fountain and Mirella Lapata. Taxonomy induction using hierarchical random graphs. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 466–476. Association for Computational Linguistics, 2012.
- Edward A Fox. The digital libraries initiative: update and discussion. *Bulletin of the American Society for Information Science*, 26(1):7–11, 1999.
- Edward A. Fox and Ohm Sornil. Digital libraries. 2003.
- Michael Friendly and Daniel J Denis. Milestones in the history of thematic cartography, statistical graphics, and data visualization. *URL* <http://www.datavis.ca/milestones>, 2001.
- Edward B Fry. Readability. reading hall of fame book, 2006.
- A. Goldst. dfr-browser. <https://github.com/agoldst/dfr-browser>, 2014.
- Andrew Goldstone and Ted Underwood. What can topic models of pmla teach us about the history of literary scholarship. *Journal of Digital Humanities*, 2(1):39–48, 2012.
- Yihong Gong and Xin Liu. Generic text summarization using relevance measure and latent semantic analysis. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 19–25. ACM, 2001.
- Dale L Goodhue. Understanding user evaluations of information systems. *Management science*, 41(12):1827–1844, 1995.
- Stephan Greene, Gary Marchionini, Catherine Plaisant, and Ben Shneiderman. Previews and overviews in digital libraries: Designing surrogates to support

- visual information seeking. *Journal of the American Society for Information Science*, 51(4):380–393, 2000.
- Seth Grimes. Unstructured data and the 80 percent rule. *Carabridge Bridgepoints*, 2008.
- Ralph Grishman. *Computational linguistics: an introduction*. Cambridge University Press, 1986.
- Jiawei Han, Micheline Kamber, and Jian Pei. *Data mining: concepts and techniques: concepts and techniques*. Elsevier, 2011.
- Roland Hausser and R Hausser. *Foundations of computational linguistics*. Springer, 1999.
- Jef Heer. A conversation with jeff heer, martin wattenberg, and fernanda viegas, 2010.
- Alexander Hinneburg, Rico Preiss, and René Schröder. Topicexplorer: Exploring document collections with topic models. In *Machine Learning and Knowledge Discovery in Databases*, pages 838–841. Springer, 2012.
- Peter Hirtle. A new generation of digital library research. *D-Lib Magazine*, 5:7–8, 1999.
- Thomas Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 50–57. ACM, 1999.
- ICML. Proceedings of the 30th international conference on machine learning (icml-13). In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 2013.
- Matthew L Jockers. *Macroanalysis: Digital methods and literary history*. University of Illinois Press, 2013.
- Steve Jones, Stephen Lundy, and Gordon W Paynter. Interactive document summarisation using automatically extracted keyphrases. In *System Sciences, 2002*.

HICSS. Proceedings of the 35th Annual Hawaii International Conference on, pages 1160–1169. IEEE, 2002.

Mehmed Kantardzic. *Data mining: concepts, models, methods, and algorithms*. John Wiley & Sons, 2011.

Rob Kitchin and Nick Tate. *Conducting research in human geography: theory, methodology and practice*. Routledge, 2013.

Krista Lagus, Samuel Kaski, and Teuvo Kohonen. Mining massive document collections by the websom method. *Information Sciences*, 163(1):135–156, 2004.

Jey Han Lau, David Newman, Sarvnaz Karimi, and Timothy Baldwin. Best topic word selection for topic labelling. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 605–613. Association for Computational Linguistics, 2010.

Jey Han Lau, Karl Grieser, David Newman, and Timothy Baldwin. Automatic labelling of topic models. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies- Volume 1*, pages 1536–1545. Association for Computational Linguistics, 2011.

M Lee, Brandon Pincombe, and Matthew Welsh. An empirical evaluation of models of text document similarity. *Cognitive Science*, 2005.

Kalev Leetaru, Shaowen Wang, Guofeng Cao, Anand Padmanabhan, and Eric Shook. Mapping the global twitter heartbeat: The geography of twitter. *First Monday*, 18(5), 2013.

Jure Leskovec, Anand Rajaraman, and Jeffrey David Ullman. *Mining of massive datasets*. Cambridge University Press, 2014.

library of congress. Library of congress. <https://www.loc.gov/>, 1800. (Visited on 01/31/2016).

Erika C Linke. Million book project. *Encyclopedia of Library and Information Science: Lib-Pub*, page 1889, 2003.

- Tie-Yan Liu. Learning to rank for information retrieval. *Foundations and Trends in Information Retrieval*, 3(3):225–331, 2009.
- MACBA. The listening eye. <http://www.macba.cat/en/exhibition-eugeni-bonet>, 2014. (Visited on 11/28/2015).
- Pattie Maes et al. Agents that reduce work and information overload. *Communications of the ACM*, 37(7):30–40, 1994.
- Christopher D Manning and Hinrich Schütze. *Foundations of statistical natural language processing*. MIT press, 1999.
- Xian-Ling Mao, Zhao-Yan Ming, Zheng-Jun Zha, Tat-Seng Chua, Hongfei Yan, and Xiaoming Li. Automatic labeling hierarchical topics. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 2383–2386. ACM, 2012.
- Gary Marchionini, Catherine Plaisant, and Anita Komlodi. Interfaces and tools for the library of congress national digital library program. *Information Processing and Management*, 34(5):535 – 555, 1998. ISSN 0306-4573. doi: [http://dx.doi.org/10.1016/S0306-4573\(98\)00020-X](http://dx.doi.org/10.1016/S0306-4573(98)00020-X). URL <http://www.sciencedirect.com/science/article/pii/S030645739800020X>.
- Mario Perez-Montoro and Jaume Nualart. Visual articulation of navigation and search systems for digital libraries. *International Journal of Information Management*, 35(5):572 – 579, 2015. ISSN 0268-4012. doi: <http://dx.doi.org/10.1016/j.ijinfomgt.2015.06.005>. URL <http://www.sciencedirect.com/science/article/pii/S0268401215000614>.
- D. Masad and S. Nayar. Sec document clustering - david masad and sanjay nayar, css 739. <http://www.davidmasad.com/sandbox/FirmClusters.html>, 1011. (Visited on 12/08/2015).
- Viktor Mayer-Schönberger and Kenneth Cukier. *Big data: A revolution that will transform how we live, work, and think*. Houghton Mifflin Harcourt, 2013.
- Andrew Kachites McCallum. Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu>, 2002. (Visited on 11/28/2015).

- Meeks, E. Documents | digital humanities specialist. <https://dhs.stanford.edu/comprehending-the-digital-humanities/documents/>, 2011. URL <https://dhs.stanford.edu/comprehending-the-digital-humanities/documents/>.
- Qiaozhu Mei, Xuehua Shen, and ChengXiang Zhai. Automatic labeling of multinomial topic models. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 490–499. ACM, 2007.
- Jessica Milstead and Susan Feldman. Metadata: Cataloging by any other name... *ONLINE-WESTON THEN WILTON-*, 23:24–31, 1999.
- Franco Moretti. *Graphs, maps, trees: abstract models for a literary history*. Verso, 2005.
- Tamara Munzner. Data types, 23, 2015.
- Robert K Nelson. Mining the dispatch. <http://dsl.richmond.edu/dispatch/pages/home>, 2010a.
- Robert K Nelson. Mining the dispatch. <http://dsl.richmond.edu/dispatch/Topics>, 2010b. (Visited on 29/04/2016).
- David Newman, Youn Noh, Edmund Talley, Sarvnaz Karimi, and Timothy Baldwin. Evaluating topic models for digital libraries. In *Proceedings of the 10th annual joint conference on Digital libraries*, pages 215–224. ACM, 2010.
- Jaume Nualart. Contributors to jaumet/visference Â· github. <https://github.com/jaumet/visference/graphs/contributors>, 2014. (Visited on 11/28/2015).
- Jaume Nualart. Jaume nualart repos at github. <https://github.com/jaumet>, 2016. (Visited on 11/02/2016).
- Jaume Nualart and Gabriela Ferraro. Towards a rhizomatic narrative. 2014.
- NYPL. Collections. nypl digital collections. <http://digitalcollections.nypl.org/collections?sort=recent#/?scroll=0>, 2016a. (Visited on 07/04/2016).

- NYPL. Nypl releases hi-res images, metadata for 180,000 public domain items in its digital collections. <http://www.nypl.org/press/press-release/january-6-2016/nypl-releases-hi-res-images-metadata-180000-public-domain-items>, 2016b.
- NYPLlabs. Gutenberg authors. <http://tools.nypl-labs.biz/gutenberg/>, 2015. (Visited on 08/04/2016).
- University of Oxford. [ota] the university of oxford text archive. <http://ota.ox.ac.uk/>, 1976. (Visited on 12/02/2015).
- Kenton O'Hara and Abigail Sellen. A comparison of reading paper and on-line documents. In *Proceedings of the ACM SIGCHI Conference on Human factors in computing systems*, pages 335–342. ACM, 1997.
- OpenDOAR. Opendoar - home page - directory of open access repositories. <http://www.opendoar.org/index.html>, 2005. (Visited on 12/07/2015).
- Alexander Pak and Patrick Paroubek. Twitter as a corpus for sentiment analysis and opinion mining. In *LREc*, volume 10, pages 1320–1326, 2010.
- Stephanie Pappas. How big is the internet, really? <http://www.livescience.com/54094-how-big-is-the-internet.html>, 2016. [Online; accessed 10-May-2016].
- F.V. Paulovich and R. Minghim. Hipp: A novel hierarchical point placement strategy and its application to the exploration of document collections. *Visualization and Computer Graphics, IEEE Transactions on*, 14(6):1229–1236, Nov 2008. ISSN 1077-2626. doi: 10.1109/TVCG.2008.138.
- Slav Petrov. Announcing syntaxnet: The world's most accurate parser goes open source, 2016. URL <http://googleresearch.blogspot.be/2016/05/announcing-syntaxnet-worlds-most.html>. [Online; accessed 13-May-2016].
- Simone Paolo Ponzetto and Michael Strube. Taxonomy induction based on a collaboratively built knowledge repository. *Artificial Intelligence*, 175(9-10):1737–1756, 2011.

- D. Quaranta. *In Your Computer*. Lulu. com, 2011.
- Daniel Ramage and Jason Chuang. Dissertation browser | information. <http://nlp.stanford.edu/projects/dissertations/>, 2012. (Visited on 01/27/2016).
- Allen H Rehear. Text encoding. *bIGITAL IIUMA\ ITIES*, page 218, 2004.
- Lisa Rhody. Topic modeling and figurative language. *Journal of Digital Humanities*, 2(1):19–35, 2012.
- Gerard Salton and Michael J McGill. Introduction to modern information retrieval. 1986.
- Roger C Schank. What we learn when we learn by doing. 1995.
- Benjamin M Schmidt. Words alone: Dismantling topic models in the humanities. *Journal of Digital Humanities*, 2(1):49–65, 2012.
- Michael Seadle and William Y Arms. The 1990s: the formative years of digital libraries. *Library Hi Tech*, 30(4):579–591, 2012.
- Michael I Shamos. Machines as readers: A solution to the copyright problem. *Journal of Zhejiang University Science A*, 6(11):1179–1187, 2005.
- Ben Shneiderman. The eyes have it: A task by data type taxonomy for information visualizations. In *Visual Languages, 1996. Proceedings., IEEE Symposium on*, pages 336–343, 1996.
- Carson Sievert. Ldavis. <https://github.com/cpsievert/LDAvis>, 2015.
- SLNSW. State library of nsw - transcripts. <http://transcripts.sl.nsw.gov.au/>, 2010. (Visited on 12/22/2015).
- SLNSW. World war 1 diaries | transcripts. <http://transcripts.sl.nsw.gov.au/project/World%20War%201%20Diaries>, 2014. (Visited on 11/28/2015).
- Justin Snyder, Rebecca Knowles, Mark Dredze, Matthew R Gormley, and Travis Wolfe. Topic models and metadata for visualizing text corpora. In *HLT-NAACL*, pages 5–9, 2013.

- M Stefaner. X by y project data visualizations. *Moritz Stefaner Information Aesthetics and Ludwig Boltzmann Institute for media. art. research, nd* [Retrieved on Dec. 17, 2010] from the Internet: <http://moritz.stefaner.eu/projects/x-by-y>, 2010.
- G. Sullivan. *Art Practice as Research: Inquiry in the visual arts*. SAGE Publications, 2005.
- Hanna Suominen, Tobias Schreck, Gondy Leroy, Harry Hochheiser, Lorraine Goeuriot, Liadh Kelly, Danielle L Mowery, Jaume Nualart, Gabriela Ferraro, and Daniel Keim. Task 1 of the clef ehealth evaluation lab 2014 visual-interactive search and exploration of ehealth data. In *Proceedings of CLEF 2014*, 2014.
- Ah-Hwee Tan et al. Text mining: The state of the art and the challenges. In *Proceedings of the PAKDD 1999 Workshop on Knowledge Discovery from Advanced Databases*, volume 8, page 65, 1999.
- Ellen Taylor-Powell. Questionnaire design: Asking questions with a purpose. *University of Wisconsin Extension*, 1998.
- Ted Nelson, et al. Project Xanadu, 1960. URL <http://xanadu.com/>. [Accessed: 2014-07-26. (Archived by WebCite at <http://www.webcitation.org/6RLo0HzFo>)].
- Lucy A Tedd and J Andrew Large. *Digital libraries: principles and practice in a global environment*. Walter de Gruyter, 2004.
- Miles A Tinker. *Legibility of print*, volume 1. Iowa State University Press Ames, 1963.
- LLC TouchGraph. Touchgraph, 2001. URL <http://www.touchgraph.com/seo>. [Online; accessed 12-November-2015].
- Transcriptorium. transcriptorium. <http://transcriptorium.eu/>, 2013. (Visited on 01/31/2016).
- Trove. About trove. <http://trove.nla.gov.au/general/about>, 2009. (Visited on 01/03/2016).

- Trove. Trove press collections. <http://trove.nla.gov.au/newspaper/>, 2010. (Visited on 01/03/2016).
- Yuen-Hsien Tseng. Automatic thesaurus generation for chinese documents. *J. Am. Soc. Inf. Sci. Technol.*, 53(13):1130–1138, November 2002. ISSN 1532-2882. doi: 10.1002/asi.10146. URL <http://dx.doi.org/10.1002/asi.10146>.
- Anne Tyle. The accidental tourist, 1985. URL https://en.wikipedia.org/wiki/The_Accidental_Tourist. [Online; accessed 09-May-2016].
- UCL. Transcribe bentham. <http://blogs.ucl.ac.uk/transcribe-bentham/>, 2000. (Visited on 01/31/2016).
- Economist Intelligence Unit. Rise of the machines. moving from hype to reality in the burgeoning market for machine-to-machine communication, 2012.
- U.Washington. University of washington digital collections. <http://digitalcollections.lib.washington.edu/>. (Visited on 12/22/2015).
- Marcos Weskamp. Newsmap. *Webdesigning Magazine, June*, page 86, 2004.
- Mitchell Whitelaw. Discover the queenslander. <http://mtchl.net/discover-the-queenslander/>, 2014.
- Mitchell Whitelaw. Generous interfaces for digital cultural collections. *Digital Humanities Quarterly*, 9(1), 2015a.
- Mitchell Whitelaw. Representing digital collections. *Performing Digital: Multiple Perspectives on a Living Archive*, 2015b.
- J. R. Williams. Guidelines for the use of multimedia in instruction. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, volume 42, pages 1447–1451. SAGE Publications, 1998.
- James Wise, James J Thomas, Kelly Pennock, David Lantrip, Marc Pottier, Anne Schur, Vern Crow, et al. Visualizing the non-visual: spatial analysis and interaction with information from text documents. In *Information Visualization, 1995. Proceedings.*, pages 51–58. IEEE, 1995.

Bibliography

Wen Zhang, Taketoshi Yoshida, and Xijin Tang. A comparative study of tf-idf, lsi and multi-words for text classification. *Expert Systems with Applications*, 38(3): 2758–2765, 2011.

A. appendix

A.1 Visference: Results of the online questionnaire:

Number of participants: 43

- Gender:

Answer	Count	Percentage
Female	12	30.77%
Male	24	61.54%
Other	1	2.56%
No answer	2	5.13%

- How old are you?

Answer	Count	Percentage
less than 20	0	0.00%
20 to 30	4	10.26%
30 to 40	12	30.77%
40 to 50	16	41.03%
50 to 60	2	5.13%
more than 60	3	7.69%
No answer	2	5.13%

- How would you rate your technical knowledge of computers and the web?

Answer	Count	Percentage
None	0	0.00%
Some	0	0.00%
Good	13	33.33%
Very good	12	30.77%
Expert	12	30.77%
No answer	2	5.13%

- This section asks about your attitudes and opinions regarding web interfaces:
How often do you use web browsers?

Answer	Count	Percentage
Once a month or less	0	0.00%
Once per week	0	0.00%
Several times a week	1	2.56%
Every day	5	12.82%
Several times a day	31	79.49%
No answer	2	5.13%

- Do you like to encounter new features in the pages you visit the most?

Answer	Count	Percentage
Not at all	0	0.00%
Rarely	2	5.13%
No opinion	7	17.95%
Sometimes	23	58.97%
Often	5	12.82%
No answer	2	5.13%

- Are you happy with the information tools and interfaces that you use?

Answer	Count	Percentage
Very happy	1	2.56%
Happy	19	48.72%
Indifferent	5	12.82%
Unhappy	11	28.21%
Very unhappy	1	2.56%
No answer	2	5.13%

- Have you ever published a paper in an academic conference?

Answer	Count	Percentage
Yes (Y)	19	55.88%
No (N)	15	44.12%
No answer	0	0.00%

This section will ask you to answer a set of questions using the Visference tool as well as the conventional presentation of the conference papers. You'll need to visit the two interfaces: JMLR (volume 28) with existing interface: [<http://jmlr.org/>] JMLR (volume 28) with Visference interface: [<http://research.nualart.cat/visference>]

- Which tool makes it easier to answer each of the following three questions:
[How many papers were accepted in the conference?]

Answer	Count	Percentage
JMLR existing interface	0	0.00%
JMLR Visference interface	31	91.18%
no difference	2	5.88%
No answer	1	2.94%

- Which tool makes it easier to answer each of the following three questions:
[Which papers are related to machine learning theory?]

Answer	Count	Percentage
JMLR existing interface	2	5.88%
JMLR Visference interface	24	70.59%
no difference	7	20.59%
No answer	1	2.94%

- Which tool makes it easier to answer each of the following three questions:
[How many papers talk about optimization?]

Answer	Count	Percentage
JMLR existing interface	0	0.00%
JMLR Visference interface	29	85.29%
no difference	4	11.76%
No answer	1	2.94%

- Which of the two interfaces would you prefer for each of these three tasks?
[Understanding the topics and themes of the conference]

Answer	Count	Percentage
JMLR existing interface	0	0.00%
JMLR Visference interface	29	85.29%
no difference	4	11.76%
No answer	1	2.94%

- Which of the two interfaces would you prefer for each of these three tasks?
[Finding papers related to your personal interests]

Answer	Count	Percentage
JMLR existing interface	1	2.94%
JMLR Visference interface	28	82.35%
no difference	4	11.76%
No answer	1	2.94%

- Which of the two interfaces would you prefer for each of these three tasks?
[Exploring new topics and discovering new research in this field]

Answer	Count	Percentage
JMLR existing interface (A1)	0	0.00%
JMLR Visference interface (A2)	30	88.24%
no difference (A3)	3	8.82%
No answer	1	2.94%

A2. A.2 Diggers diaries II: List of the one hundred topics with terms and first labelling draft not

Note 1	Notes 1	ID	Terms
military life		62	0.20607 men time work great made troops number present days part fact large make place officers australian war general order military
		54	0.18175 good time lot back bit boys pretty big things don't chaps night morning round chap day put thing place bad
		10	0.16237 time told back found made man men thought asked officer place put called gave find knew make long decided looked
		6	0.09093 day sunday tuesday morning monday today saturday thursday wednesday tuesday afternoon night june good august usual april september quiet
Personal?	thoughts of life	8	0.08711 man like men people war years english good world french young great country women australian war read soldier mind things
?		76	0.08259 time light great air lights long minutes guns smoke sight sound hear men noise round half hour side heard dark
Personal?	homesick	56	0.07774 day leave days time home good spent england year great months xmas dec australia week left weather france return christmas
Personal	Weather	3	0.07364 water day hot sand feet men night dust mud cold good wet heat sleep camp rain bad hard bath dry
?	landscape	26	0.0708 miles side hills country river town view place road high top sea fine distance water long small great large
Personal	Weather	15	0.07068 cold night day morning rain snow weather wind heavy day wet fine raining yesterday pay frost warm ground blowing sunday
??		97	0.06925 time war days long back day france things don't home life months thing hope great feel good place hard ago
EDITOR NOTES		4	0.06825 transcribed page previous blank preceding transcript n/a error photograph html duplicate printed transcription reverse preceding pages repeat certificate version multi-page
Personal	Letters Family	42	0.06768 letter letters dear received home mother time mail love hope write loving news send father good week writing wrote days
Personal	Morning / breakfast	93	0.06223 night sleep bed good room morning time tea breakfast slept mess till day put camp comfortable hit blankets men tent
Military life	Camping	90	0.06215 camp move orders ready arrived morning a.m left transport p.m horses men night pack march tents back day road moved
War	Front line	30	0.05909 line trenches front men night back party wire dug post work company yards firing battalion fatigue digging shell day
Personal	Food	60	0.05713 tea bread dinner biscuits tin food jam good day breakfast beef rations meat butter bully milk eggs fruit eat coffee
Personal / Tourist?	free time Entertainment	61	0.05636 good concert evening show night tea afternoon home dinner party half music spent y.m.c.a band time room give splendid enjoyed
Tourist	city life	9	0.05606 town people streets place city women soldiers fine native french large street natives shops round english houses places small buildings
Traveling	Railway	5	0.05508 train station arrived left camp p.m a.m journey leaves night morning marched railway hours good trip miles caught reached town
Personal	Letters Family	20	0.05448 hope letter time dear good write don't i'm letters love heart things back mrs give news long glad i've writing
War	Combat	48	0.05434 wounded men killed shell dead man hit head poor shot died back leg left buried bullet badly blown wounds blood
War	Health	21	0.05377 hospital ward sick doctor bed patients bad feeling day cases medical leg put days wounded sister wound morning arm feel
Tourist	city tourist spots	55	0.05035 church building place built fine large stone walls room beautiful wall inside house high years cathedral side round tower small
Military Life	Parades & Marches	66	0.04931 parade march drill morning afternoon camp day route order inspection men marched guard full marching church evening company usual leave
War	Combat	34	0.04759 shells guns friz gun shell night fire artillery heavy bombardment line shelling trenches firing enemy big close quiet fired day
War	Air force	95	0.0472 planes bombs dropped guns plane air machine night brought aeroplanes friz enemy german aeroplane flying taube lines bomb flew raid
Traveling	Over the sea	45	0.04578 sea deck day ship passed boat weather night morning land call board miles sight rough wind port today side hot
Military life	general	74	0.04528 general brigade officers major battalion division officer col capt staff leut command colonel australian gen men corps birdwood army div
War	News	13	0.04474 never french war german british germans front great traps prisoners france today english fighting papers big australians army germany captured
Tourist	Countryside	32	0.04105 trees green country fields beautiful lovely pretty flowers grass crops tree village fruit garden place road small growing gardens leaves
War	Combat	46	0.04035 turks men enemy trenches wounded attack position left guns fire killed machine heavy casualties night artillery fighting infantry captured
Personal	Personal care & clothes	80	0.03995 boots pair white socks black clothes wear red issued post hat shirt wearing picture card clothing cap clean trousers uniform
EDITOR NOTES		78	0.03836 diary letter written pages australian note page letters notes book copy war printed august records australia transcriber's paper april read
Tourist	London	73	0.03767 london train back tea met hotel home park place dinner afternoon walked hours arrived bus left caught lunch city walk
War	Front line	88	0.03721 line front guns enemy night attack barrage artillery friz road position forward stunt heavy left prisoners moved advance gun morning
Military life	General	38	0.0368 men office officers man major court orderly sergeant charge corporal guard company told colonel mess military room put martial private
Military life	Turkey	14	0.0358 horses miles camp turks water night left camels turkish horse back moved camped regt wadi camel arish brigade desert sand
		39	0.03545 road shell wood village mud left shells roads german ground ypres bapaume place holes runs town blown traffic friz dead
Military life?		40	0.03498 miles marched march battalion camp night line village albert moved back left days day place move time billets morning big
Personal?	Killing time	99	0.03444 day home tea till bed afternoon evening morning dinner back fine read wrote letters breakfast good usual spent rain sunday
Military life?	General	57	0.03342 round good general men lunch afternoon mess found today evening morning rode tonight colonel returned officers h.q. back day bde
Military life	navy	29	0.03286 ships ship troops submarine passed boat boats port british destroyers cruiser submarine board sea speed miles sunk transports night a.m.
Personal	Colors & beauty	0	0.03279 sun blue sky white beautiful light bright green red colour black beauty grey night shining world clouds moon gold full
		98	0.03274 day morning afternoon tea jan wrote parade dinner night feb evening letters sunday march letter monday home good tuesday saturday
War?	Health	86	0.03262 wounded station dressing ambulance bearers day stretcher night field post amb cases friz section back number busy patients line work
Traveling	French roads	27	0.03044 village town left miles road motor amiens billets marched omer french baileul march arrived back somme place hazebrouck kilos lorry
??	Health	75	0.03021 capt sgt killed pte l.h. wounded margin major indecipherable cpt batt jack sick met smith leut col note bill left
Traveling	Over the sea	69	0.02931 boat ship boats ashore board harbour water wharf aboard ships shore men port alongside small deck beach side troops coal
War	Combat	58	0.02917 turks trenches beach fire night shells firing guns sharpen turks gun day morning hill heavy artillery quiet men rifle gully
		59	0.0289 work good day time to-day working hard job days hours things put week fair section-to-night making usual duties dead
Military life	Sports	36	0.02886 played won football sports match game cricket afternoon team day playing good race held boxing games play cards officers parade
Military life	chaps & comrades	44	0.02742 band boys men great people passed crowd played gave flags playing soldiers king troops cheers marched past cheering french singing
		53	0.0274 head long side feel end sketch top small horse front piece black rope drawing left cut water wheel hand hair
Travelling	Over the sea	18	0.02731 port board sydney melbourne ship boat left wharf leave ashore arrived bay harbour p.m troops fremantle town sea cape columbo

100-labelled-crossreads_keys

Personal	sadness	71	0.02712 dear great son death sympathy god boys loss home heart feel soldier died mother sad life proud kind mrs men
Personal	Money	91	0.0271 pay money paid bank book dear office comforts leave london soldiers sydney fund mrs received send made funds australia case
		24	0.02664 morning o'clock back today afternoon tonight night camp time good regt put horses duty received dinner troop arrived horse ill
		23	0.02545 indecipherable today home crossed horses usual bed january monday cairo friday tuesday thursday april mail february letters wednesday tomorrow day
Personal	Family & memories	65	0.025 pm sister sisters miss nice mrs duty i'm love home a.m time happy hope lovely day tea poem matron afternoon
Tourist	Drinking	41	0.02435 beer canteen bought money bottle wine drink french sold photos buy cost price frances bottles pay paid films glass coffee
Travelling	UK & France	28	0.02217 camp salisbury leave france weymouth left england london hut back draft arrived hill boat train miles training sutton plymouth days
Military life?	Egypt & Turquie	47	0.02133 hospital alexandria cairo left lemmos egypt arrived troops island camp april wounded heliopolis to-day aust gallipoli ship back anzac base
Personal	Religion	82	0.02129 church service sunday parade morning held chaplain padre sermon attended afternoon communion services y.m.c.a evening address rev holy present gave
Military life	General?	35	0.02073 sydney australia battalion served pte mrs embarked nsw private australian lieutenant gallipoli field hmat returned service enlisted france october father
		72	0.01979 generally make feel things long interesting find makes pass talk return thing takes mind times case word past trouble matters
Tourist	Egypt	11	0.01974 cairo mile pyramids egypt egyptian heliopolis desert mosque camp gardens menu native sphinx tram natives pyramid citadel hotel city donkeys
War	Context	50	0.01952 war german germany australia government military british terms minister sir peace conscription vote general president country britain forces governor states
Military life?	Camp & horses?	19	0.01947 canal camp cairo suz left col fuller desert arrived major kanbara brigade march men train february officers sand port returned
Personal	Letters & family	67	0.01819 html fine morning afternoon day cleaning harness round evening horses stables friday back paddock saturday night thursday warm tuesday december
Military life?	Camp & horses?	49	0.01811 home mum dad wrote letter ellis write letters mrs george recd meet day indecipherable lovely play auntie walk bed night
Military life	Training & weapons	96	0.01789 horses bty battery horse lines wagon night l.p.m fine a.m day men sgt today guns gun weather morning evening
Tourist	France	12	0.01778 school gun training rifle drill work practice range today bayonet instruction coy musketry shooting afternoon lecture class lewis machine parade
Travelling	Railway & roads	22	0.01727 paris hotel french place rue dinner cafe visit opera indecipherable walked find club people english round city met train girls
War	Cemetery	7	0.01653 train engine line railway large trucks arrived left back run station boys depot number engines made yard started road trains
		37	0.01641 graves dead war grave buried crosses battle cemetery land soldiers html cross death fought lie god broken peace bones soldier
		43	0.01488 house room back engine home html park load farm made french broken fire friz apl supplies left wood oil
		17	0.0146 envelope mrs active justice service reverse post south field walesaustralia censor card addressed australia n.s.w signed postcard image fergusonsupreme shows
		52	0.0137 good letter p.m letters girls day hope night morning nash hospital dear kitty a.m sydney australia joseph mrs colonel egypt
		83	0.01359 sydney german island rabaul ship emden men wireless board officers naval native germans islands captain guinea left sova sept fleet
War	Combat	92	0.01257 gas helmets helmet trenches alarm billets June bombardment masks shells attack tear heavy steel ammunitions france issued trench box mask
		2	0.01193 aug wet sat mon tues sun fri jan letter nov mar sept dec jun feb mother thur july thurs oct
Military life	Communications	25	0.01138 turned p.m a.m breakfast till tea dinner fed rested fell came cleaned camp returned guard day stand saddled hot watered
War	Prisoners	16	0.01123 stop division anzac landing turks telegraph position troops bay army ashmead attack positions suvia turkish august baba gallipoli report enemy's
		94	0.01057 german prisoners camp today received food english germans officers dated germany parcels room cross red escape prisoner french tonight bridge
Military life	navy	84	0.01005 signal time naval flag wireless message berrima ship brigadier station received herbertshohe despatch messages receiving beresford telephone transmitting troops
Personal	Poetry?? CHECK	89	0.00886 november mrs december rup met fri thursday october tuesday hotel london monday wednesday saturday lunch afternoon evening home miss club
		89	0.00886 november mrs december rup met fri thursday october tuesday hotel london monday wednesday saturday lunch afternoon evening home miss club
War	Air force	81	0.0088 sea ship ships fleet arrived anchored proceeded german p.a. aug squadron admiral harbour convoy a.m anchor weather cruiser left firing
		85	0.00712 book tonight today books read mail night yesterday reading back morning full poems boche early day half indecipherable long frank
		64	0.00692 oct nov gen officer left indecipherable afternoon col tidworth troops morning visited wednesday called sunday thursday tuesday pencil monday saturday
		70	0.00653 wher camp interests whch intemeee haewe prif military day soldiers commandant compound owing issued made australia guard police charge received
		87	0.00645 day night quiet fine letter rec sun wed tues mon posted html sat indecipherable writing letters fri any heavy arr
		79	0.00634 south wales transcribed library state spell misspell transcriber's possibly judy gibmert john peter smith ditto mayo betty lynne adrian bicknell
		33	0.00536 lieut enemy squadron bde machines air report pilot machine aerodrome flying capt observer reconnaissance reported a.f.c area pilots m.c wadi
Personal	Food	1	0.00508 hun food hun men parcels russian russians misery awful received whilst english prisoners control reliable promptly british caused bread french
CHECK	songs?	77	0.00439 men rations baked n.c.o's unit hospital n.c.o's man temp leave bakery n.c.o.'s strength rejoined or flour ovens bakeries dough personnel
Personal	In French	63	0.00438 yer we're song puff voice we'll sing there's love don't tommy we've ill it's you're blokes html snake you'll that's
		51	0.00402 les des pour bon vous nous c'est tres pas qui guerre est monsieure agrave par une dans france tout sont
War	Context	31	0.00312 miss mrs indecipherable hill rita coy t.n.r.a.n.c/o king grafton j.t batt john health nurse tindale baby trained testimonials
		68	0.00107 page kms miles france east south west north belgium called village resting harvey flanders german centre battle coast tel australian

