



IMPERIAL COLLEGE LONDON

C401

MENG COMPUTING FINAL REPORT

NewSumm
NEWS AGGREGATOR AND SUMMARISER

Author:
Kunal M L Wagle

Supervisor:
Anandha Gopalan

June 4, 2017

Abstract

Abstract to be written

Acknowledgements

Acknowledgements to be written

Contents

1	Introduction	11
1.1	The Problem	11
1.1.1	The Media	11
1.1.2	The Consumer	12
1.2	News Aggregator and Summariser: A potential solution	13
2	Market Research	15
2.1	News Reading Habits	15
2.1.1	Digital News Report 2016	15
2.1.2	Potential User Survey	17
2.2	Related Products	25
2.2.1	Google News	25
2.2.2	Flipboard	26
2.2.3	Yahoo News Digest	27
2.2.4	Apple News	27
2.2.5	Comparing the existing products	28
3	Background Research	30
3.1	Machine Learning Techniques	30
3.1.1	Topic Modelling Techniques	30
3.1.2	Topic Labelling Techniques	31
3.1.3	Clustering Techniques	33
3.2	Summarisation Techniques	35
3.2.1	Extractive Summarisation Techniques	35

3.2.2	Abstractive Summarisation Techniques	36
3.3	Natural Language Processing Libraries	40
3.3.1	Aspects of Natural Language Processing	41
4	Design	42
4.1	Front End Architecture Diagram	42
4.2	User Story	42
4.2.1	Home Page	42
4.2.2	Topic Dashboard	43
4.2.3	Search Results Page	43
4.2.4	Topic Viewer	44
4.2.5	Topic Settings	45
4.2.6	Article Viewer	46
4.2.7	Topic Lists	48
4.2.8	Profile Settings	48
4.3	Back End Flow Diagram	49
4.4	Language and Platform Choices	50
4.4.1	Front End	50
4.4.2	Back End	51
4.5	Infrastructure	51
4.5.1	Hardware	51
4.5.2	Database	52
4.5.3	Infrastructure Decisions	53
5	Implementation	56
5.1	Database Schema	56

5.1.1	Schema Diagram	56
5.1.2	Articles	56
5.1.3	Topics	57
5.1.4	Clusters	58
5.1.5	Users	59
5.1.6	Digests	60
5.2	Potentially Useful APIs	61
5.2.1	News Outlets	61
5.3	Potentially useful libraries	66
5.3.1	Mallet	66
5.3.2	Natural Language Processing	66
5.4	Machine Learning	70
5.4.1	Topic Modelling	70
5.4.2	Topic Labelling	72
5.4.3	Clustering	74
5.5	Summarisation	74
5.5.1	Extractive Summarisation	74
5.5.2	An attempt at Abstractive Summarisation	74
5.6	Restlet	74
5.7	Server Tasks	74
5.8	Front End	74
5.9	Key Classes	74
6	Optimisation	75
6.1	Speed Optimisations	75

6.1.1	Topic Modelling	75
6.1.2	Topic Labelling	75
6.1.3	Clustering	75
6.2	Memory Optimisations	75
7	Evaluation	76
7.1	Machine Learning and Summarisation	76
7.2	Summary Analysis	76
7.3	User Interface Evaluation	76
8	Conclusion	77
9	Future Work	78
9.1	Foreign Languages	78
9.2	Further Optimisations	78
9.3	Other Apps	78
References		79
Appendices		82
A	Source Code	82
B	API	83
C	User Guide	84
Index		85

List of Figures

1.1	The Digital News Report 2016 on trust in the news	11
1.2	A Google News Search on 24th May 2017	13
2.1	Why people use News Aggregators	15
2.2	Concerns with personalisation of News Feeds	16
2.3	What is a good way to get to the news?	16
2.4	Survey Graph surrounding how often news is read	18
2.5	Survey Graph asking how many sources are used	19
2.6	Survey Graph about consistency of sources	19
2.7	Survey Graph regarding how much of an article is read	20
2.8	Survey Graph regarding summarisations of articles	20
2.9	Survey Graph about combining articles	21
2.10	Survey Graph concerning News Aggregators and summarisation . . .	22
2.11	Survey Graph regarding searching for topics in the news	22
2.12	Survey Graph regarding news aggregators and searching	23
2.13	Survey Graph about News Digests	23
2.14	Survey Graph about News Digests' frequency	24
2.15	Survey Graph concerning potential platforms	24
2.16	Screenshots from Google News	25
2.17	Screenshots from the Flipboard app	26
2.18	Screenshots from the Yahoo News Digest app	27
2.19	Screenshots from the Apple News app	28
3.1	An example of a Dependency Tree	37
3.2	A concept created during multimodal semantic summarisation	39

3.3	A flowchart showing the process of semantic graph summarisation	40
4.1	A basic diagram of the expected architecture of the front end	42
4.2	A wireframe of the Home Screen	43
4.3	A wireframe of the Topic Dashboard	44
4.4	A wireframe of the Search Results page	45
4.5	A wireframe of the Topic Viewer	46
4.6	A wireframe of the Topic Settings page	47
4.7	A wireframe of the Article Viewer	47
4.8	A wireframe of the Topic List page	48
4.9	A wireframe of the Profile Settings page	49
4.10	A guide to the expected flow of the Back End	50
5.1	A diagram of the MongoDB schema	56

List of Tables

2.1	Related Products in the market and their benefits and drawbacks . . .	29
4.1	Advantages and Disadvantages of different Infrastructure combinations	54
5.1	Readership statistics for selected UK outlets	62
5.2	Space used for OpenNLP models	68
5.3	Parameters for training topic models	71
5.4	Parameters for estimating topic models	72

Listings

1	A sample document in the Article table	56
2	A sample document in the Topics table	57
3	A sample document in the Summaries table	58
4	A sample document in the Users table	60
5	A sample document in the Digests table	60
6	A sample response to an API call to The Guardian	62
7	A sample response to an API call to News API	63
8	A sample response to an API call to Wikipedia's API	65
9	Analysing an annotated document using Stanford's CoreNLP	68

1 Introduction

1.1 The Problem

1.1.1 The Media

The Fourth Estate is currently enduring one of its toughest passages of time to date. Some of the challenges it faces are familiar, whilst some are new ones that have been created by the onset of the digital age.

The term ‘fake news’ is now at the forefront of our minds in an atmosphere that is highly charged. The Digital News Report 2016 found that in the United Kingdom only 50% of readers trust what they say in the news, whilst only 29% say that they trust journalists (see Figure 1.1). This statistics appear to be even more dire for the media in the Untied States of America, when the percentages are 33% for news channels, and 27% for print journalists.



Figure 1.1: Digital News Report findings showed that when it came to trust in the news, in only a few countries was trust as high as 50% in the accuracy of a piece of news. In countries such as the United States it is far lower.

In addition to this, the media have been submerged in a deluge of websites that also claim to disseminate news to readers. In an increasingly digital world, anyone with access to the internet could write a blog about a news story. The popularity of social networks such as Facebook and Twitter also mean that these amateur ‘journalists’ have a platform on which they can compete with professionals from the industry.

In 2016, the concept of fake news dominated the general election in the United States. Shortly after the country went to the polls, Hillary Clinton’s campaign became dogged in rumour after more than a million tweets were sent out about the ‘pizzagate’ scandal. The fake news went both ways. In July 2016, it was reported

widely on blogs that Donald Trump had said in an interview with *People* in 1998 that:

'If I were to run, I'd run as a Republican. They're the dumbest group of voters in the country. They believe anything on Fox News. I could lie and they'd still eat it up. I bet my numbers would be terrific.'

The claim was later found to have been fabricated.

Fake news and the people's trust is not the only major issue facing the media at the moment. Another key is sheer volume. Given the abundance of media outlets in the world at the moment, and perhaps more importantly (due to the internet) the accessibility around the globe, there is a lot of competition between outlets.

1.1.2 The Consumer

As a result of a lot of the aforementioned accessibility, consumers also have a challenge in that they can't efficiently read all the news that is out there. If I were to search Google News on the 24th May 2017 for news about Donald Trump, I'd find a story about the President's meeting with the Pope at the Vatican (see Figure 1.2). On further inspection, I can see that there are several pages of results about the meeting itself, from dozens of different news sources. A consumer has little to no hope of reading all of these.

So which news sources should a consumer prioritise? Often, that comes down to political opinion. In fact, the Digital News Report 2016, conducted by the Reuters Institute for the Study of Journalism, says that only 34% of consumers in the United Kingdom believe that the media is 'free from undue political influence'. In the United States, which doesn't have a state broadcaster like the United Kingdom does in the BBC, the corresponding statistic is significantly lower at 21%.

The immediate consequence of statistics like this is that if someone hasn't read all the sources in a piece of news, they are reasonably likely to develop some form of political bias on the issue (perhaps without even realising it). For example, someone who reads about a story in a left-leaning publication such as *The Guardian* may feel differently when compared to someone who reads the corresponding story in a right-leaning publication such as *The Daily Mail*.

The simple solution to this would be reading the story in as many sources as possible before forming a final opinion. However, as mentioned before, this is just simply not feasible. In July 2016 a study was conducted by the Statistic Brain Research Institute that said that the average attention span of a person is now as low as 8.25 seconds. In addition to this, the Institute's findings also said that the average

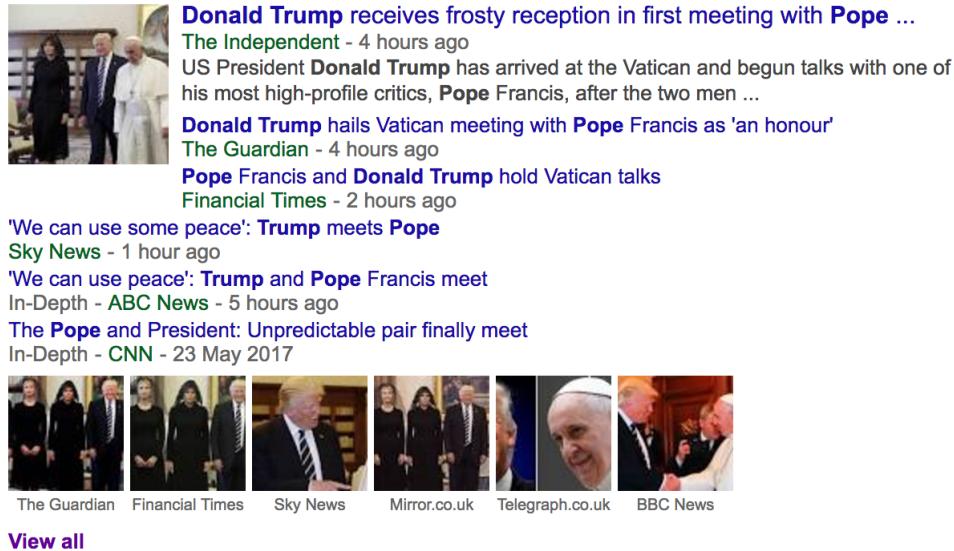


Figure 1.2: A Google News search for stories about President Trump's visit to the Vatican on 24th May 2017

person only reads 28% of a webpage of average length (593 words). Another key finding was that users spend only 4.4 seconds to cover each additional 100 words on a webpage. It's simply not possible for every news source to capture a user's attention, especially if they are all broadly covering the same story.

1.2 News Aggregator and Summariser: A potential solution

The fundamental aim of my project would be to create a project that directly addresses the consumer's concerns mentioned above, whilst indirectly aiding the concerns of the media.

The project would be a News Aggregator, that would allow users to see the news from multiple mainstream (and therefore 'trusted') sources in one place. In addition to this though, the aggregator would also act as a summariser, so that a user doesn't need to read as much as before to gain the same opinion.

However, a key aim has to be to address the issue of there just simply being too much news for someone (who is reading from multiple sources) to take in. To try and achieve this, the News Aggregator will need to identify articles across sources that are about the same topic (for example, Donald Trump's visit to the Vatican), and summarise them. This will mean that a user only has to read a single article to

get the full information from multiple sources about a news piece.

In an indirect attempt to allay the concerns of the media and the public about the phenomena that is fake news the News Aggregator will need to find a way to counteract it. One obvious way to do this would be through the limitation of the number of outlets to only those that are ‘trusted’. In addition to this, a clear indication when reading the summary of what facts different outlets all agree on, and what facts they either differ on or omit, will allow the user to easily corroborate parts of the story.

A further secondary aim to this could be to allow a user to customise their summary, taking in information from the outlets of their choice (from the original ‘trusted’ list). The advantage to this would potentially be that if a user decides that too many issues from one source are uncorroborated, they can remove them, and see only information from the other sources that have written about a topic. The disadvantage to this could be that if a user dismisses an entire section of the political spectrum as fake continuously, they’ll still pick up a political bias on the facts of the topic. Although, it could be said that the Digital News Report 2016 findings suggest that this is a feature that won’t be used much. One of the more interesting findings in the report was that more than 60% (in the United Kingdom) were concerned about the potential impact of personalising their news feeds. They felt that it could lead to both ‘missing information’ and ‘missing challenging viewpoints’.

2 Market Research

2.1 News Reading Habits

2.1.1 Digital News Report 2016

Conducted by the Reuters Institute for the Study of Journalism[24] in Oxford University[31], the Digital News Report[23] is an annual study that serves to investigate how people access and find out about news in various countries across the globe, including the United Kingdom.

The study attempts to cover, amongst other things how people get the news, how they use social media, what types of news they trust, and both the reasons to use and concerns about news aggregators. Findings of the report relevant to my concept are listed below:

Reasons to use News Aggregators

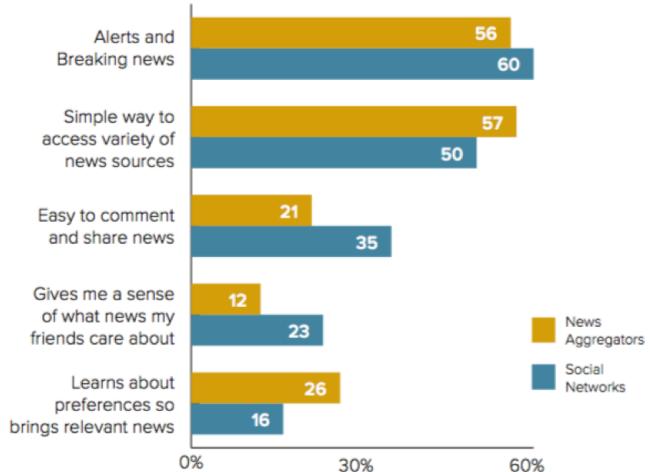


Figure 2.1: A graph showing why people use News Aggregators and Social Networks as news sources

As shown in the graph in Figure 2.1 the key to the popularity of a News Aggregator is that it's simplicity and speed. It's important to bring news as soon as it happens, and it needs to provide access to a variety of different news sources. Of lesser importance is the social aspect - the need to comment and share the news, although the fact that 26% of people want their preferences to influence the news that they're interested in should not be discounted.

Concerns about Personalised News

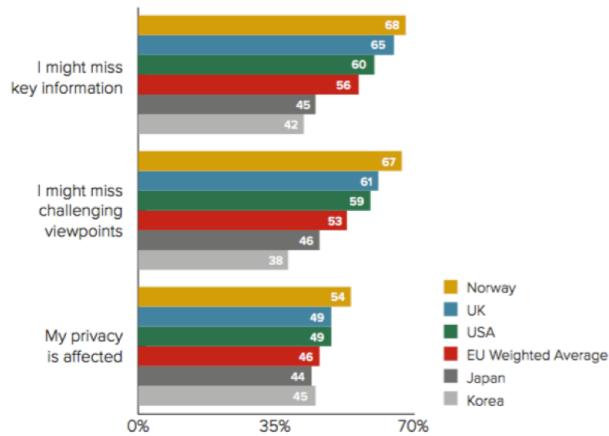


Figure 2.2: A graph showing the main three concerns that people have with Personalised News Feeds determining what they read.

Key concerns, especially in the United Kingdom, with Personalised News is that information might be incomplete. It is not difficult to see why this may be a concern - if you limit your news to a subset of the sources that might be available then it's possible to miss parts of the story - the parts that challenge aspects of the articles you're reading. It's therefore important that the project reflects this concern, which is why it is listed as a secondary aim (See section ??).

What should determine what someone reads next

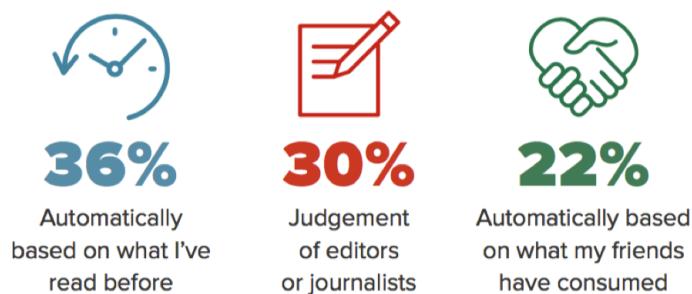


Figure 2.3: Key statistics in response to the question: Is this a good way to get the news?

Least popular are algorithms that present news based on what people's friends have been reading. In general, according to the Reuters Institute[24], 'People think *they*

are the best judge of what's important to them.' More interesting, is that people are now less inclined to allow journalists or editors to push articles to the top of reading lists than having an algorithm suggest similar stories to what they've just consumed. This could perhaps indicate that they aren't as trusting in journalists and editors as before - something that the Reuters Institute touches upon in an essay near the end of the report.

Conclusions

The Reuters Institute for the Study of Journalism is a trusted and respected organisation, and this study has been supported by institutions varying from media outlets such as the BBC[3], news search engines such as Google, and other Universities, such as University of Canberra[30]. Importantly, it's also supported by regulators, including Ofcom[22], which is the regulator in the United Kingdom.

The survey was conducted on their behalf by YouGov[39], which is a highly respected polling company in the United Kingdom, and it surveyed more than 50,000 people, including over 2000 in the United Kingdom. In addition, the report is coupled with numerous essays on various key points that arose from the study. The writers of these essays varied from the Director of Research at the Reuters Institute to the Chief Executive Officer of The New York Times[29].

As a result of this, I deemed the Digital News Report 2016[23] to be a very reliable source of information with regards to how people read the news, and their preferences in the field.

2.1.2 Potential User Survey

Following this initial research I decided to conduct a survey of potential users to delve into more specific aspects of News Reading Habits. I was more focused on the number of sources people use to ascertain the facts about a piece of news. I also asked questions regarding the summarisation of news and news digests.

I was interested in the questions about summarisation and digests after reading some data from a study conducted by the Statistic Brain Research Institute[25]. Statistic Brain conducted a study in July 2016 that suggested that the average attention span of a human is now only 8.25 seconds[16]. They coupled the presentation of this data with statistics about people browsing the internet taken from the paper *Not Quite the Average: An Empirical Study of Web Use*[33]. The statistics suggest that the average user only reads 28% of the words on the average webpage, and that for each additional 100 words beyond that, the user would only spend an additional 4.4

seconds on the page.

Admittedly, the paper was written in 2008, and so might not reflect average web browsing activity, but it raised a question worth answering in my opinion. Do users actually read full articles on the news, and if not, would summarising articles help make sure they didn't miss out on key points, thus counteracting an issue discovered in the Digital News Report 2016[23] (Section 2.1.1)?

The survey was presented on Google Forms[10] and was answered by 96 respondents. Further findings are shown below:

1. How often do you read the news?

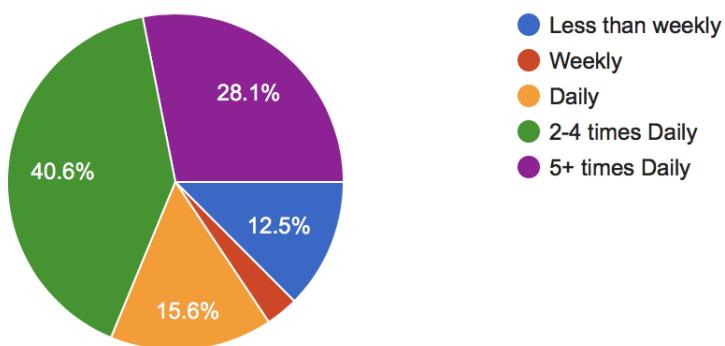


Figure 2.4: A pie chart showing how often people read the news.

The first question (Figure 2.4) on the survey was designed to ascertain how often people read the news. Overall, well over 80% of respondents said that they read the news at least daily, and nearly 60% more than once per day.

2. How many sources do you read?

Next I asked about how many sources a user read. It was quite rare for someone to use only one source for a piece of news. In fact, over 90% said that they used two or more sources, with three and five sources being the most prominent selections.

3. Are these the same sources every time?

The respondents were a lot more split on the question of consistency of sources, although more than half said that they used the same sources every time, with

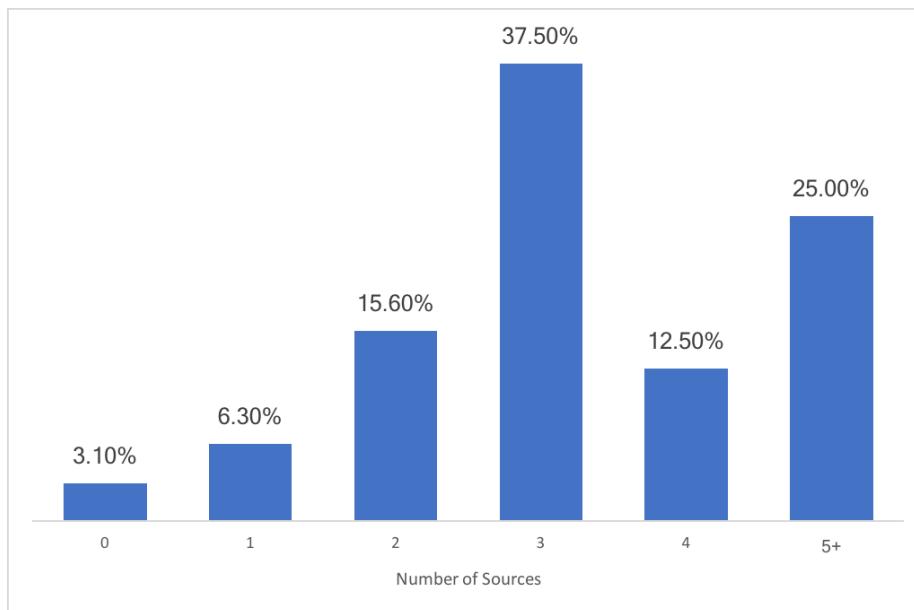


Figure 2.5: A bar graph showing how many sources people normally use for a source of news.

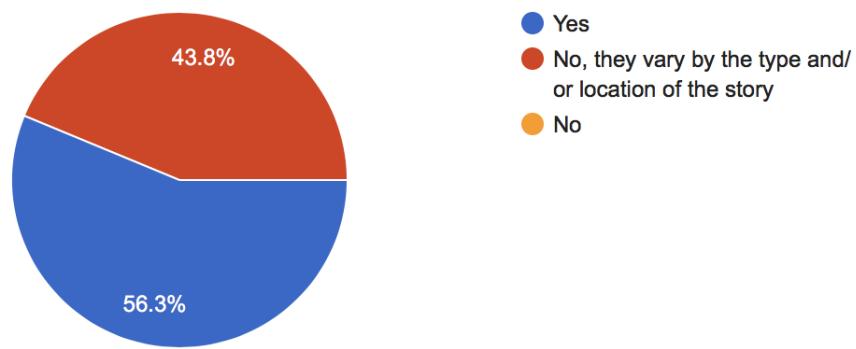


Figure 2.6: A pie chart showing whether people use the same sources every time.

43.8% saying that they vary their sources based on the type and/or location of the news story.

4. How much of an article do you normally read?

I further expanded upon the initial reading I had done about attention spans by asking potential users how much of articles they normally read. Only 18.8%

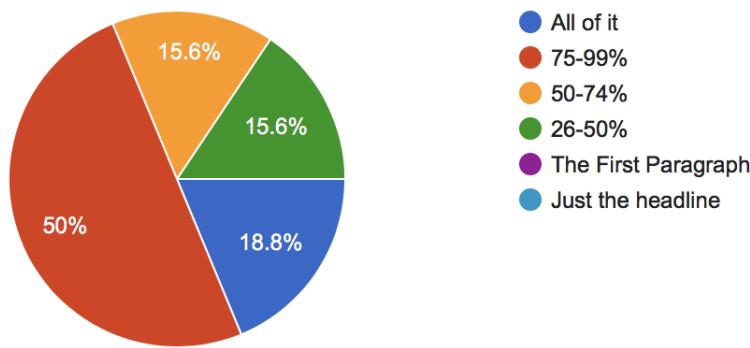


Figure 2.7: A pie chart showing how much of articles people normally read.

of respondents said that they read the entirety of articles (Figure 2.8), and as much as 15.6% confess that they read less than 50% of an article on average. It's plausible that the users could be missing key facts through not reading the entire article. On reflection, I could have also asked if this was a concern to those users.

5. If an article was presented as a shorter summary, would you read more/all of it?

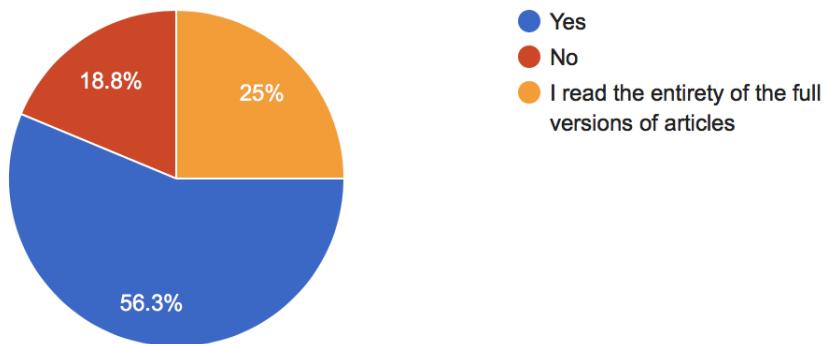


Figure 2.8: A pie chart showing what percentage of people would read a shorter summary of an article.

Another question asked was whether reading a shorter summary would result in users reading more or all of the article. Encouragingly, more than half

(56.3%) of respondents said that they would read a shorter summary.

6. Would you read a summary of different articles on the same topic combined?

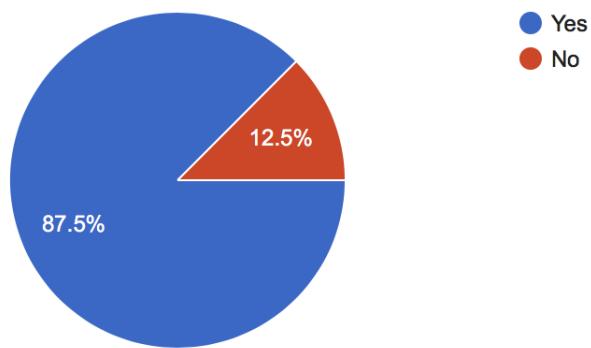


Figure 2.9: A pie chart showing whether people would read a summary of different articles on the same topic combined.

With this question (Figure 2.9) I aimed to obtain an idea from potential users as to whether combining articles about the same topic from different sources and then summarising it would be considered useful. The question got a resounding response, with as much as 87.5% saying that they would read a summary like this.

7. Would you use a News Aggregator that summarised articles?

Consolidating the findings of the previous question was the fact that 81.3% said they would use a News Aggregator that summarised articles (Figure 2.10). Amongst those who said they wouldn't, there were comments along the lines of 'I'm not great with technology' given as explanation. An interesting comment however, said that 'nuance could be lost in the summarisation'. This is a valid point, and so a key point of the evaluation process will have to be focused on the summarisation of articles itself, to ensure it's not losing important information at any stage.

8. How often do you search for news on a specific topic?

For the questions in Figure 2.11 I tried to get an idea of how much people search for news on a specific subject of interest to them. In general, this came

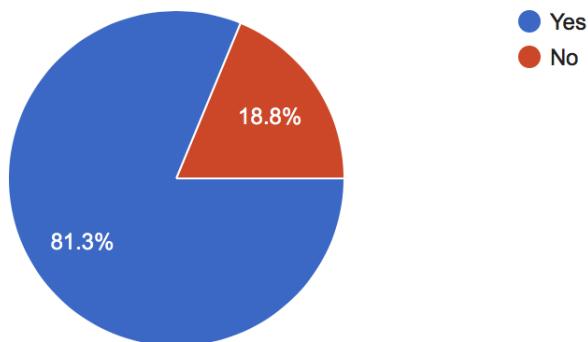


Figure 2.10: A pie chart showing answers to the question: ‘Would you use a News Aggregator that summarised articles’

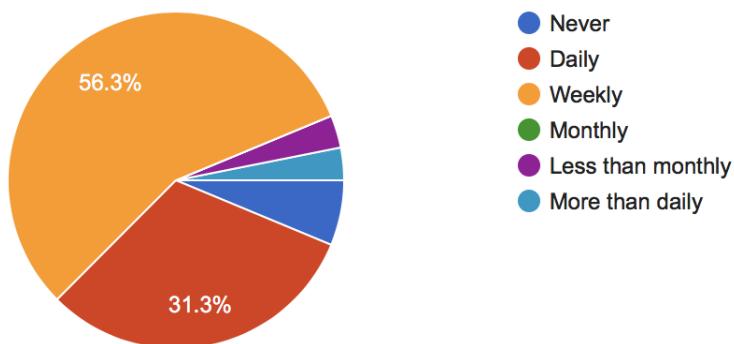


Figure 2.11: A pie chart showing how frequently people search for news on a specific topic.

out to be less frequent than reading the news itself, with just over half the respondents saying that they search for a specific topic weekly. However, a healthy proportion (31.3%) search daily for specific topics, and some search even more often than this.

9. Should the news aggregator have the ability to search for a specific topic?

As much as 81.3% of respondents agreed that the News Aggregator should have a function to search for specific topics.

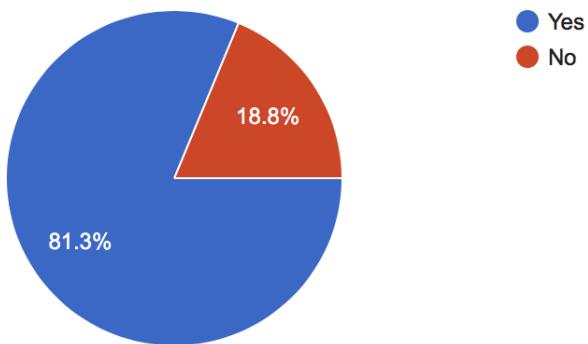


Figure 2.12: The answers to a question asking if the News Aggregator should have the ability to search for a specific topic.

10. Do you like services such as news digests that prepare lists of articles each day that apply to a topic you may be interested in?

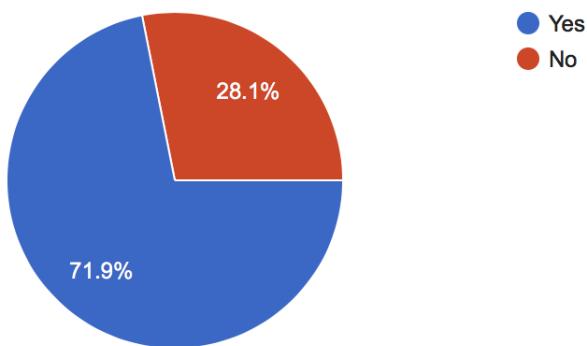


Figure 2.13: A pie chart showing people's attitudes towards News Digests.

The next question, shown in Figure 2.13 centred around News Digests. Nearly 72% of respondents said that they liked services that provide news digests.

11. How often should these news digests be updated?

People who were in favour of news digests in general felt that these should be updated at least daily (47.8%), with a further 34.8% on top of that saying it

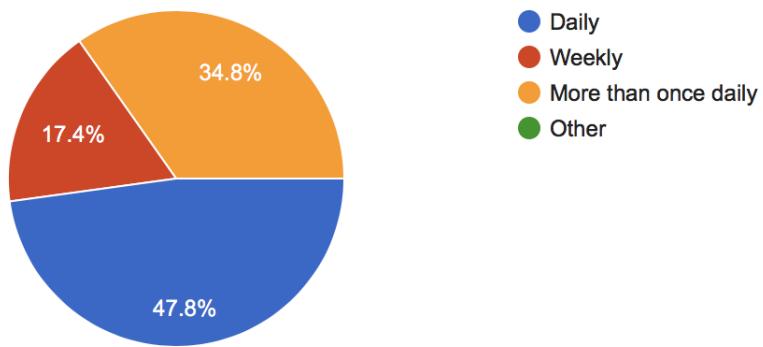


Figure 2.14: Responses regarding how frequently news digests should be updated.

should be updated more than once.

12. What platform should the aggregator be developed for?

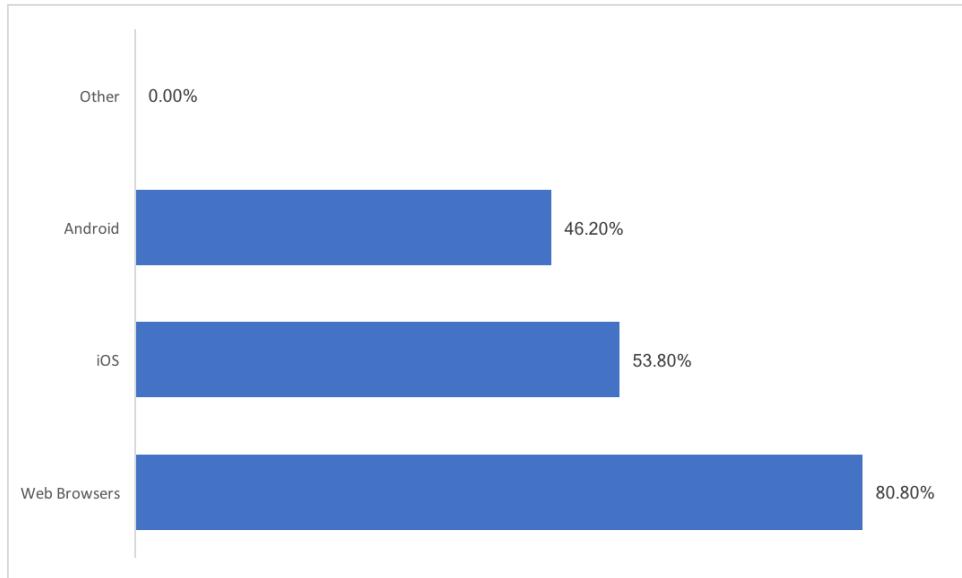


Figure 2.15: A bar chart showing responses when asked which platform they'd rather see the news aggregator developed for.

Finally, respondents were asked which platform they would prefer to see the news aggregator developed for. 80.8% said they'd want a version for web browsers, and 53.8% for iOS and 46.2% for Android. This would suggest that

I should develop a web solution and then move to a mobile platform afterwards if there is time.

2.2 Related Products

2.2.1 Google News

The screenshot shows the 'News' section of the Google News homepage. The 'Top Stories' header is visible, followed by a large image of Martin McGuinness. Below the image, the headline reads 'Martin McGuinness resigns as deputy first minister of Northern Ireland'. The story is from The Guardian and includes a link to 'See realtime coverage'. To the right, there's a sidebar for 'Recent' news items like 'Jeremy Hunt says four hour A&E target only applies to 'urgent' cases' and 'Brexit Briefing: Angela vs Theresa. Sign up for your new newsleine newsletter'. At the bottom right, there's a weather forecast for Edgware, England.

(a) A section of the homepage for Google News.[11]

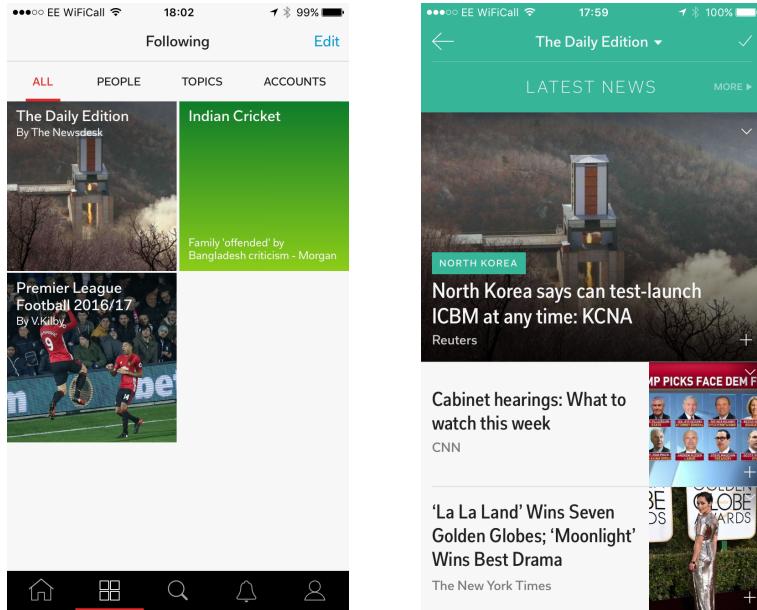
The screenshot shows a Google search results page for the query 'uber'. The 'News' tab is selected. The results are dated '21 Dec 2016' and sorted by relevance. Several news articles are listed, all related to Uber's self-driving car pilot being revoked in San Francisco. The articles include titles from TechCrunch, The Sun, and Financial Times, along with links to Patch.com and the San Francisco Examiner. A 'View all' button is at the bottom.

(b) Google News uses clustering techniques to group articles about the same topic together.

Figure 2.16: Screenshots from Google News

Initially developed early in the century and released in 2006, Google News[11] is a free-to-use news aggregator. Google News operates in a similar manner to the traditional Google[9] search engine, thus making it a go to aggregator when searching for a specific topic. Google News also groups ‘similar’ articles together using Clustering techniques[15]. Google News operates purely as a go-between - when clicking on an article the user is taken straight to the media source itself, rather than being able to read the article on Google News itself.

2.2.2 Flipboard



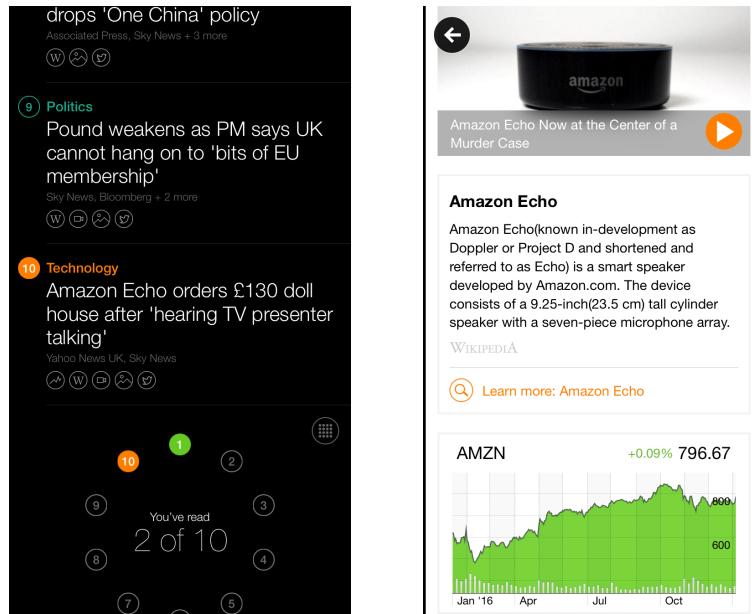
(a) Users can subscribe to ‘magazines’ that they may be interested in.

(b) All users are subscribed to The Daily Edition, which puts together the latest most popular news.

Figure 2.17: Screenshots from the Flipboard app

Flipboard[6] is a much more recent attempt at a news aggregator (developed in 2010), and relies on the concept of users subscribing to ‘magazines’ on different topics. There’s a central ‘cover page’ on the home page that shows the most recent stories from across all a user’s subscriptions. Like Google News, links from the desktop website send the user to the original media source, while links from within the mobile applications open a browser within the app itself.

2.2.3 Yahoo News Digest



(a) A section of the home screen, which displays the top ten stories for the day.

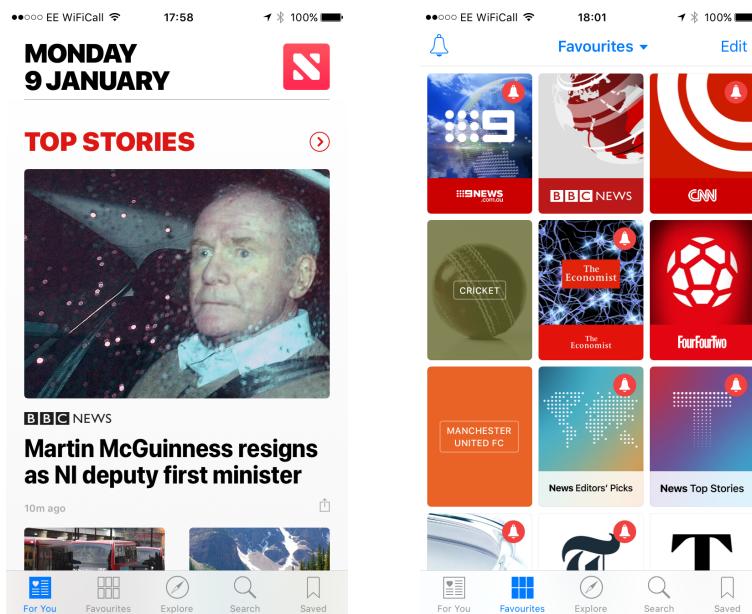
(b) The Yahoo News Digest app provides links and infographics that are relevant to the article.

Figure 2.18: Screenshots from the Yahoo News Digest app[38]

Yahoo News Digest[38] is a direct evolution from Summly[26], which was an app that summarised news. Yahoo News Digest is a phone application that creates two digests a day: one in the morning and one in the evening. Each digest contains the ten leading articles from the previous 12 hours. Each article provides a summary of the story and links to relevant other pages - such as articles from Wikipedia[36].

2.2.4 Apple News

Apple News[2] is an application that is installed by default on all recent iOS devices. Apple's default attempt at a news aggregator allows users to select their preferred news sources, and from a selection of topics. Apple[1] then presents on a home screen news from those sources and topics. Clicking on each article keeps it in the native application, rather than sending the user to the media source itself.



(a) The home screen for the Apple News app

(b) Apple News allows users to select sources and topics to be their favourites.

Figure 2.19: Screenshots from the Apple News app[2].

2.2.5 Comparing the existing products

The existing products are compared in Table 2.1 on page 29.

Product	What it does well	What it doesn't do well
Google News	<ul style="list-style-type: none"> Groups similar articles together 	<ul style="list-style-type: none"> Doesn't host articles itself - therefore making it harder to navigate than perhaps could be possible. Although, this could be to avoid any copyright issues.
Flipboard	<ul style="list-style-type: none"> Similar to a traditional search engine and so is easy to use Obtains articles instantly, thus providing most up-to-date information when searching 	<ul style="list-style-type: none"> Provides a home page that allows a user to see the most popular stories at the time related to the user's subscriptions. Available as a mobile application
Yahoo News Digest	<ul style="list-style-type: none"> Yahoo News Digest won the 2014 Apple Design award[37]. The articles also provide key infographics, quotes, and other information potentially relevant to the article. 	<ul style="list-style-type: none"> Topics are much broader than Google News, thus meaning that users can't necessarily search for topics specific-enough for them. Flipboard requires registration before being able to read articles
Apple News	<ul style="list-style-type: none"> Allows selection based off both topics and the news sources themselves. 	<ul style="list-style-type: none"> There are only ten articles available per digest, and no capability for searching by topic. Digests are only produced in the morning and the evening, so the news articles presented could be out of date. Navigation on the application is not simple. If a user has accessed many articles from push notifications, then the user could have to press the back button several times to get back to the home screen. The topics that a user can subscribe to can be quite limited, and aren't reactive to current affairs - for example, if a natural disaster occurs, you couldn't then subscribe to that natural disaster as a topic.

Table 2.1: Comparing the relative benefits and drawbacks of Google News, Flipboard, Yahoo News Digest and Apple News

3 Background Research

3.1 Machine Learning Techniques

3.1.1 Topic Modelling Techniques

Topic modelling is a subsection of Machine Learning that aims to determine what topic a given document is about. The topics wouldn't be named at this stage, they would simply be given generic names such as Topic A and Topic B. Assigning names to topics will be done at a later stage (see Section 3.1.2).

Latent Semantic Indexing

Also known as Latent Semantic Analysis[19], Latent Semantic Indexing (LSI) was one of the initial forerunners in the field of topic modelling. It uses singular value decomposition to locate patterns in the text of a document and thus form a basis on which to categorise the document. A major benefit of LSI is that it is fast to train, but in general it has lower accuracy when compared to models that are probabilistic, such as Latent Dirichlet Allocation[35].

Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a generative probabilistic model that assumes that the topic distribution has a Dirichlet prior. The theory behind LDA is that a document of words contains a mixture of different topics, and this is reflected in the final answers given for the algorithm.

A rough algorithm for performing LDA[34]:

1. **Set n to be the number of topics there are in the document.** We can do this by trial and error.
2. **Assign every word w in the document d to a random topic.** These topics are temporary. At this stage we can remove function words, such as 'the'. However, we keep duplicates - in fact, at this stage they could be in different topics.
3. **Check and update topic assignments.** To do this, we loop through each word in the document, taking note of how prevalent the word is across documents, and prevalent those topics are in the document. These two probabilities are then passed to a sampling algorithm to generate a new topic for the word. This step is usually completed using the statistical model *Gibbs Sampling*.

4. Repeat step 3 until there are no more topic-reassignments.

3.1.2 Topic Labelling Techniques

Topic Labelling techniques in the project for labelling the topics that are generated from the LDA analysis of the document in Section 3.1.1. When conducting research into this topic specifically, I found a paper (*Automatic Labelling of Topic Models* by Lau, Grieser, Newman and Baldwin in 2011[20]) that documents the creation of an algorithm that labels topics using Wikipedia[36] titles. This could work very well in my project, as Wikipedia titles as title headings would allow users to search for topics that are both broad and specific.

A rough version of Lau, Grieser, Newman and Baldwin's algorithm:

1. **Calculate the top 10 topic terms.** This is done by finding the marginal probabilities of each word from the original topic models, and taking the top ten. The marginal probability of a term is the probability of that term being randomly selected given a topic t .
2. **Search Wikipedia using these terms.** We also search Google[9] using a site restricted search (to Wikipedia) and take the top eight results from each. These are called the *primary labels*.
3. **Isolate all ‘noun chunks’ from the terms.** In this case noun chunks are combinations of words from the terms that appear next to each other. For example, with the term ‘Summer Olympic Games’ the noun chunks would be ‘Summer’, ‘Olympic’, ‘Games’, ‘Summer Olympic’, ‘Olympic Games’ and ‘Summer Olympic Games’. Note that ‘Summer Games’ is not a noun chunk as the words don’t appear juxtaposed. These noun chunks are added to the primary labels from step 2 and are deemed *secondary labels*.
4. **For each noun chunk:**
 - Check to see if the noun chunk is the title for a Wikipedia article
 - Remove the noun chunk if it doesn’t correspond to a Wikipedia article
5. **Calculate the *Related Article Conceptual Overlap* scores.** Related Article Conceptual Overlap (RACO), developed by Grieser et al in 2011[12], is a calculation designed to identify the strength of relationship between terms by

inspecting the category overlap between the terms' corresponding articles. We do this for each secondary label still remaining - details on how to calculate the RACO scores are explained in further detail below.

6. **Discard all secondary labels with RACO score of less than 0.1.**
7. **Add five highest topic terms to the list.** Now we return to the original list of topic terms from step 1 and add the five highest to the remaining candidates.
8. **Perform candidate ranking.** There are multiple ways to do this, but the original paper recommends using a variety of statistical methods, based around a T-test, the Chi-squared test and a log-likelihood test. The aim is to estimate how closely related the candidate is to all the terms in the topic. We then take the top candidate as our final answer.

Calculating the *Related Article Conceptual Overlap*

Related Article Conceptual Overlap (RACO) was first introduced as a concept in the paper *Using Ontological and Document Similarity to Estimate Museum Exhibit Relatedness*[12] by Grieser, Baldwin, Bohnert and Sonenberg in 2011. It compares two terms and estimates their similarity by comparing the two terms' similarity on Wikipedia. The core calculation for RACO goes as follows:

$$\text{Category} - \text{Overlap}(a, b) = \left| \left(\bigcup_{p \in O(a)} C(p) \right) \cap \left(\bigcup_{p \in O(b)} C(p) \right) \right|$$

In this equation, $O(a)$ represents the outlinks (the links to other articles from the Wikipedia article) of an article a and $C(p)$ represents the set of categories that article p is a part of.

An issue with the Category-Overlap calculation as it is, is that there's a bias for articles that are larger than others as they will have more outlinks but won't necessarily be in more categories. As a result it's normalised using Dice's coefficient to produce the final RACO equation:

$$\text{sim}_{\text{RACO}}(a, b) = \frac{2 \times \left| \left(\bigcup_{p \in O(a)} C(p) \right) \cap \left(\bigcup_{p \in O(b)} C(p) \right) \right|}{\left(\bigcup_{p \in O(a)} C(p) \right) + \left(\bigcup_{p \in O(b)} C(p) \right)}$$

3.1.3 Clustering Techniques

Cluster analysis is a machine learning technique that is used to put items that are similar to each other into groups. In practice it's used by Google News[11] to put articles about the same topic together for a user[15]. There are multiple commonly used types of cluster analysis:

Centroid Clustering

In Centroid Clustering[4], also known as k-means clustering, there are k clusters. A vector is calculated for each article in the list. The article is then assigned to the cluster that is closest to its vector score. A major downside to this method however is that it requires k to be defined in advance. In the context of this project, that's not applicable as we don't know how many different articles we are clustering based on.

Density Clustering

Density Clustering[4] also involves calculating a vector score for each article. Once these have been calculated, they can be graphed, and the areas of the graph that have highest density are chosen as the clusters. An advantage of this over the centroid clustering methods is that we don't need to know the number of clusters beforehand. However, density clustering can become less accurate as it requires areas of sparse density on the graph to precisely separate the different groups, which isn't always possible.

Hierarchical Clustering

Hierarchical Clustering[14], which is also known as Connectivity Clustering, is based on the idea of using distance measures between articles to identify which ones are most similar. There are two approaches to Hierarchical Clustering:

- *Agglomerative*, which is a bottom-up approach, assigns each item to its own cluster and then merges pairs that are closer together.
- *Divisive*, a top-down approach, that begins with all items in a single cluster, and then proceeds to split it into multiple clusters.

Creating a vector score for each item

The first step in each of the three clustering techniques is to create a vector score

for each item. As this will play a big part in the final results, it's important to get this stage right. This can be split into two steps:

1. **Find definitive terms within the article.** This can be done using techniques such as the popular *Term Frequency-Inverse Document Frequency* (TF-IDF), which is explained in further detail below. Proper nouns would be useful in this step, as they are more likely to be relevant to what the article is specifically about.
2. **Create a vector.** This vector, based from the definitive terms from the first step, would be a set of keywords and corresponding weights.

Term Frequency-Inverse Document Frequency

TF-IDF[27] is designed to identify terms that appear frequently in one article that don't occur a lot over the entire set of articles (also known as the *corpus*). Variations of it are commonly used by search engines to identify search results that are most relevant to a query. The calculation for TF-IDF is as follows:

Term Frequency

Term Frequency can be most simply calculated as the frequency of a term in a document. However, this could result in a bias towards terms that appear in longer articles. A more accepted way to calculate the Term Frequency therefore is to use a normalising function, called augmented frequency:

$$tf(t, d) = 0.5 + 0.5 \times \frac{f_{t,d}}{\max\{f_{t',d} : t' \in d\}}$$

Inverse Document Frequency

Inverse Document Frequency is used to check in how many documents of a corpus D a given term t occurs. It is an inverse function, so as to minimise the value for the term when it is common amongst various different documents. It is also logarithmically scaled. It is calculated using the following formula:

$$idf(t, D) = \log \frac{|D|}{|d \in D : t \in d|}$$

Term Frequency-Inverse Document Frequency

$$tfidf(t, d, D) = tf(t, d) \times idf(t, D)$$

3.2 Summarisation Techniques

Summarisation techniques will be (predictably) used for summarising the merged articles that were identified by the processes of Topic Modelling , Topic Labelling and Clustering. There are two types of summarisation:

- **Extractive Summarisation**, which consists of taking sentences that are important from the original text, and discarding the rest.[13]
- **Abstractive Summarisation**, which aims to generate a piece of text using natural language techniques. A key to this is that some words in the final summary may not have been the original piece of text.

3.2.1 Extractive Summarisation Techniques

LexRank and TextRank

There are two well known examples of extractive summarisation: LexRank and TextRank.

LexRank and TextRank have very similar methods for extracting a summary[5]. Initially a graph is constructed, that consists of one node for each sentence in the corpus. Then a clustering algorithm is applied, using a TF-IDF calculation to determine similarities.

There are a couple of key differences between LexRank and TextRank. The first arises in the calculation. Both use a TF-IDF calculation, but they are varied slightly. LexRank uses a cosine similarity function in order to weight the final calculation, whereas TextRank uses a more simple logarithmic weighting to perform the calculation.

Both use Google's famous PageRank algorithm to then rank the importance of each sentence based on the calculations in the previous step. With d being a damping factor, the PageRank of a node u is given as:

$$p(u) = \frac{d}{N} + (1 - d) \sum_{v \in adj[u]} \frac{p(v)}{\deg(v)}$$

After this has been calculated, the top ranked sentences are taken by the TextRank algorithm to form the summary. However, in the LexRank algorithm sentence position and length are also taken into consideration. Also, when adding a sentence to the summary, the LexRank checks the sentence against the summary to ensure that it won't be redundant. As a result of this extra step, LexRank is considered more suitable than TextRank when summarising multiple documents, whereas TextRank is only normally used for summarising a single document.

3.2.2 Abstractive Summarisation Techniques

There are six common methods for abstractive summarisation, which can be split evenly into two distinct categories[17]. These two categories centre around the creation of a representation of the given document:

- **Structure based methods** consist of techniques that involve determining the important information in a document by considering its structure[18]. Ways of doing this include fitting the given document to a template, or converting the text in a tree-like structure.
- **Semantic based methods** involve building a semantic representation of the given document and then feeding that into an algorithm that generates natural language that forms the final summary.

Structure based summarisation

Tree based summarisation

With tree based structuring, the first step is to create a dependency tree to represent the document. A dependency tree is a tree with a node for each word in a sentence, the links between the trees show which words depend on each other. For example, given the sentence *This is an example*, we can construct a tree as in figure 3.1:

In the tree (figure 3.1), the word *This* is dependent on *is* and so is linked. *an* is linked to *example* in a similar way, and *an example* is dependent on the verb *is* and so also has a link to it.

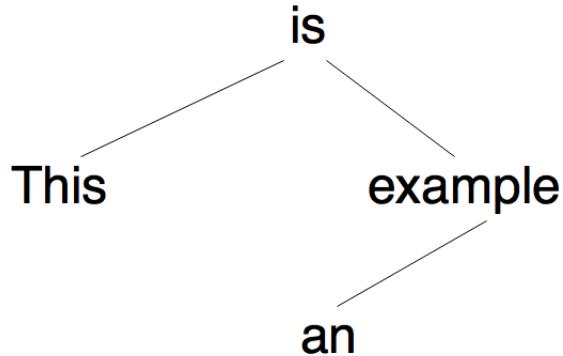


Figure 3.1: A possible dependency tree for the sentence This is an example

Once dependency trees have been created for each of the sentences a similarity algorithm is used to find sentences that are similar. The common phrases between these similar sentences are then taken and form the basis of the final summary. A language generator then combines the common phrases and arranges it to create a final set of summary sentences.

An obvious disadvantage to this method is that if only common phrases are taken from the sentences, context could be lost from some of the sentences. However, on the other hand, the use of a language generator means that a more coherent, more grammatically correct summary is formed.

Rule based summarisation

In rule based summarisation[8], the process is centred around a list of pre-determined categories. With each of these categories, there is a pre-determined set of questions to be answered.

For example, with a theoretical category *Product Launch* we could have the following questions:

1. What is the name of the product?
2. What company is launching it?
3. What type of product is it?
4. What's new about it?
5. What price is it retailing at?

This represents only a subset of the possible questions we could have for this category.

To perform rule based summarisation, the first step is to analyse the document and determine which category it fits best. Once that's been determined, the next stage is to analyse the text to find answers to the questions corresponding to that category. These answers are then fed in to a natural language generator that forms sentences, and thus the summary.

Results for this method of abstractive summarisation have been promising, but the method has a major disadvantage in that the list of categories and questions needs to be pre-determined. As a result, it might not be an optimal algorithm to use for the ever-changing world of news reporting.

Ontology based summarisation

Methods of ontology based summarisation have been developed, primarily using domain based ontology. In this method a 'domain expert' defines a domain ontology for a news event. Each new document is then classified into a topic using these domain ontologies. Important phrases are determined by how close they are to items in the ontology. These phrases are then passed into a natural language generator to form the final summary sentences.

A key disadvantage to this method is that a lot of the domain ontologies has to be manually determined by the 'domain experts' and so can be very time consuming. As a result this may also not be particularly optimal for summarisation of news events.

Semantic based summarisation

Multimodal Semantic summarisation

Multimodal semantic summarisation can work on documents that contains both text and images. First, a semantic model is built to represent the document. This model is made up of different concepts that are surmised from the text. For example, the sentence *Multimodal Semantic summarisation is an example of an algorithm that performs abstractive summarisation* could form the concept shown in figure 3.2:

Concepts are gradually filled with more information as the entire text is analysed. Links are also added between concepts that share some relationship. For example in figure 3.2 if there was another sentence that surmised a concept called *Abstractive Summarisation* then there would be a link from *Algorithm1* to that new concept.

The next step is to rank the concepts. This is done by taking into account the completeness of the concept, and the number of links that the concept has to others. This

Algorithm1
Name: Multimodal Semantic Summarisation Subset: Abstractive Summarisation Complexity: Unknown
Sentence: Multimodal Semantic summarisation is an example of an algorithm that performs abstractive summarisation

Figure 3.2: A possible concept created from the analysis of the sentence
 Multimodal Semantic summarisation is an example of an algorithm that performs abstractive summarisation

way, the concepts are ranked by which is most important to the original document. Once the key concepts have been identified summary sentences can be generated featuring these concepts.

Information Item based method

Information Item based summarisation[7] relies on the content of the summary being determined from an abstract representation of the original document, rather than the sentences from the document themselves. To do this, the document is first scanned so that Information Items (InIt) can be generated. An information item is defined as being ‘the smallest element of coherent information in a text or sentence’.

Once the information items have been created, they are then ranked using frequency analysis to find the most important predicates and entities from the original document. This step is near identical to the term frequency stage in extractive summarisation (Section 3.2.1). These information items that are ranked highly are then combined and fed into a natural language generator to form the final summary sentences.

Semantic Graph summarisation

Semantic Graph summarisation centres around a rich semantic graph (RSG). The document is first converted into a RSG. Each node in the RSG represents a noun or verb in the document, and the links between the nodes represent the semantic and

topological relations between these nouns and verbs. In the second stage, heuristic rules are used to reduce the semantic graph to a more minimalistic version. This will form the basis of the final summary. In the final step, the minimised RSG is passed into a generator that creates the final summary sentences.

The steps are outlined in the flowchart provided in figure 3.3.

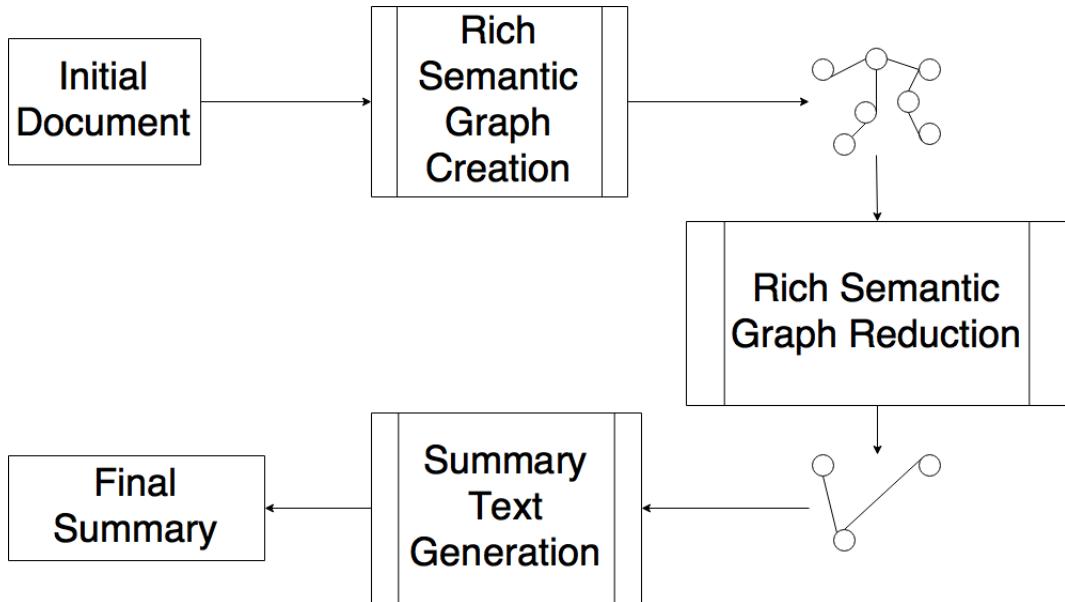


Figure 3.3: This shows the process of semantic graph summarisation in a flowchart based on those provided in [17, 18]

Semantic Graph summarisation has had success with producing an abstractive summary that has fewer redundant sentences, and is also good at producing grammatically correct sentences. However, it's only designed for use with a single document as input. As a result it might not be suitable as a method for summarising the multiple documents needed for the aggregator, unless it's combined with another method.

3.3 Natural Language Processing Libraries

A lot of the summarisation techniques specified earlier, especially for abstractive summarisation, would require a full analysis of the semantics of a body of text. As a result, it's clear that a Natural Language Processing library will be needed for this task.

3.3.1 Aspects of Natural Language Processing

There are several different tasks that a Natural Language Processor. A few key ones that are most likely to be required for the project are explained briefly below:

- **Tokenisation** is the splitting of words in a given document.
- **Sentence Segmentation** is the splitting of sentences in a given document.
- **Part-of-Speech (POS) tagging** takes a list of tokens and analyses them, returning tags for each. A tag represents the grammatical function of the token in the sentence (for example if it is a noun, verb or adjective).
- **Named entity extraction** identifies the proper nouns within a document, such as people, dates or locations, as well as other categories of noun.
- **Chunking** takes a list of tags from a POS tagger and groups sets of tokens by their function in a sentence. Examples can include noun phrases and verb phrases.
- **Parsing** takes a sentence and develops a tree showing the functions of each section of the sentence.
- **Coreference Resolution** identifies noun phrases within a document that are the same. This can commonly be used for replacing pronouns with the original noun phrase.

4 Design

4.1 Front End Architecture Diagram

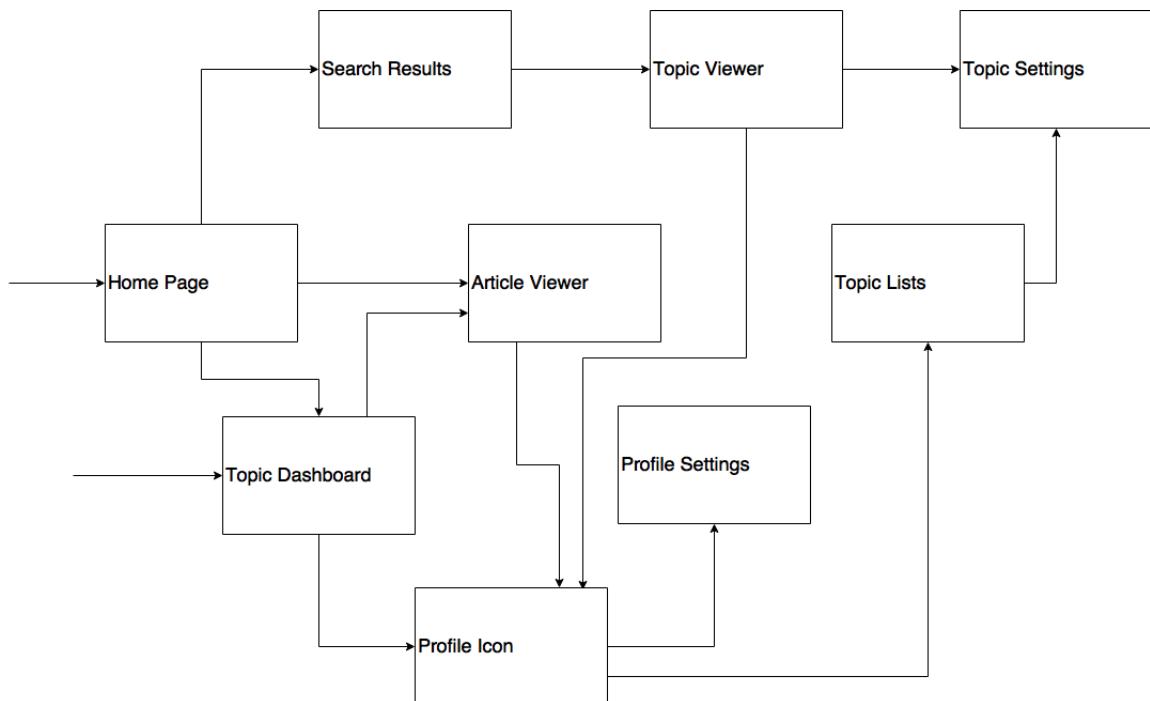


Figure 4.1: This diagram doesn't show all anticipated connections between each facet of the application's front end. For example, the profile icon is in theory accessible anywhere when the user is logged in.

4.2 User Story

4.2.1 Home Page

If a user has not logged in, the first page they will see on opening to the website will be the Home Page. The key to the design of the home page is that it's easy to get started for a user. They are presented with a large search box that they can use without having to sign in. There'll also be a row of trending articles displayed below the search bar. This row will consist of icons consisting of images and titles below, as shown in the wireframe. There's a navigation bar at the top, that is present on all screens, with a button to register and a button to sign in on the top right, and a link to information about the application itself on the top left.

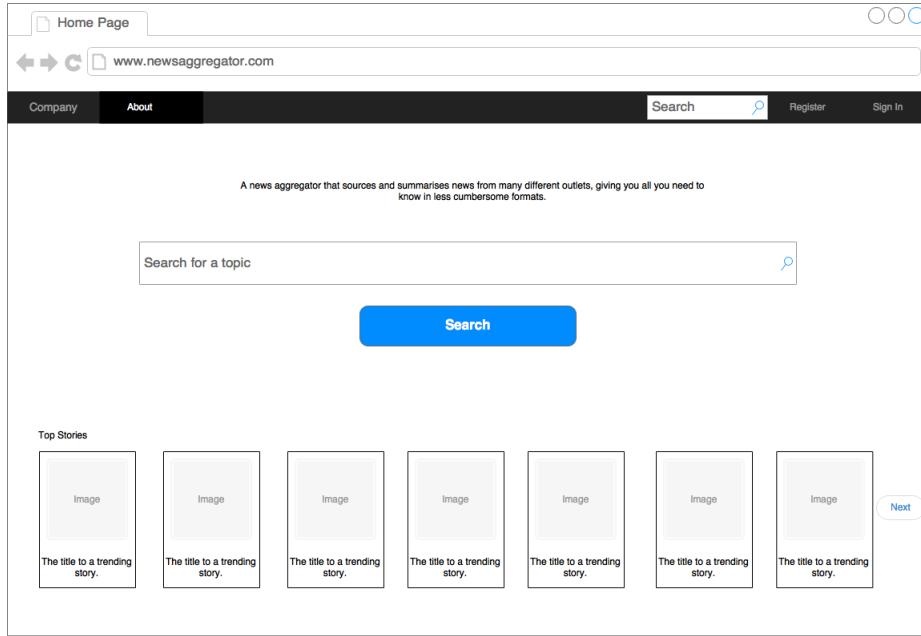


Figure 4.2: A key focus of the home page is that the user can easily search for articles without having to sign in or register.

4.2.2 Topic Dashboard

For a user that is logged in, or alternatively a user that has logged in (or signed up) from the home page will first see the topic dashboard screen that shows sets of articles corresponding to each topic they are subscribed to. For users who don't have any topics that they are subscribed to, they will instead go back to the home page. The topic dashboard initially shows leading articles from across all the topics they are subscribed to. They can drag to one direction (akin to a sideways swipe on a phone or tablet) to see articles pertaining to the next topic.

Articles themselves are presented as panels. On a panel is an image for the article, with the title overlaid. In the top right hand corner are icons corresponding to the sources that the article has been summarised from. Leading articles (those out of the most recent that have the most sources attached) are placed in the larger panels on the top left and bottom right of the screen, as shown in the diagram.

4.2.3 Search Results Page

The search results page is designed to allow users to see as much as possible and do as much as possible without having to leave the page. The results are presented in

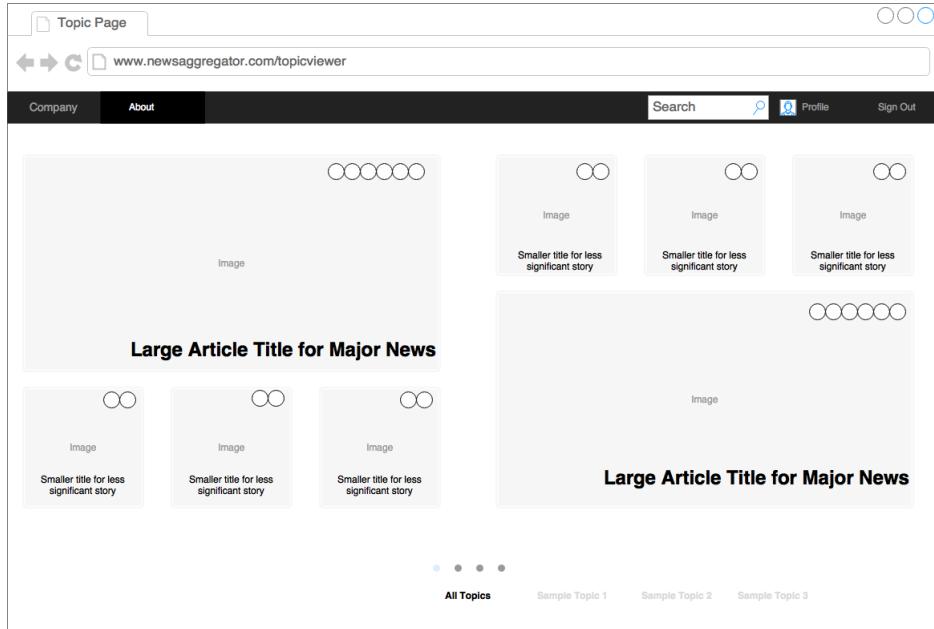


Figure 4.3: The topic dashboard has the ability to show articles across all topics, or pertaining to a specific topic.

the application’s panel design, with an image corresponding to each panel. Below the panel, is a description of the topic, that is likely to be created from the first paragraph of the Wikipedia[36] article about that particular topic. Below this are two buttons, one to view the topic’s articles, which will go to a topic viewer page (shown later), and a subscribe button. If logged in, the user is taken to the topic settings page, where they can choose any particular preferences they have for this specific topic. If the user isn’t logged in, they are taken to a page allowing them to log in or register, before being taken to the preference page.

It’s necessary at this stage to have a user sign in if they want to subscribe to a topic. This way, they can access their topics from anywhere. In addition to this, the search bar is still prominent on the search results page, so that users can edit their queries easily in case of minor errors.

4.2.4 Topic Viewer

The topic viewer is accessed if a user has clicked on the view button from the search results page. This screen is similar to that of the topic dashboard. Minor changes include the page control at the bottom of the screen being replaced by page statistics, including figures about how often articles are created, the number of subscribers,

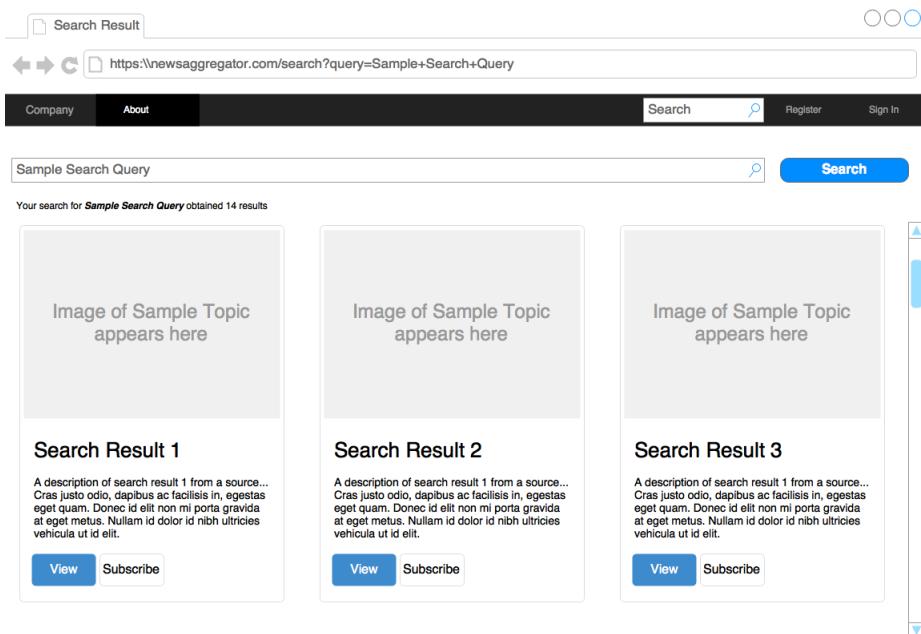


Figure 4.4: The Search Results page has a focus on showing users as much information as possible about each topic.

and the total number of hits that the page has had. The other major difference is the presence of a subscribe button in the top right hand side of the screen. As with the search results page, this leads to the topic preferences screen if logged in, and to a page prompting a user to sign in or register when not already logged in.

The remaining aspects of the screen are the same - namely the panelled articles, and the icons in the top right of each panel that indicate which sources the article originated from.

4.2.5 Topic Settings

The aforementioned Topic Settings page sets out two preferences from the user for a specific topic. These are namely:

- **Should the topic's articles be included in the user's email digest?**
This is answered by a simple on/off switch.
- **Which sources should be used to construct summaries**

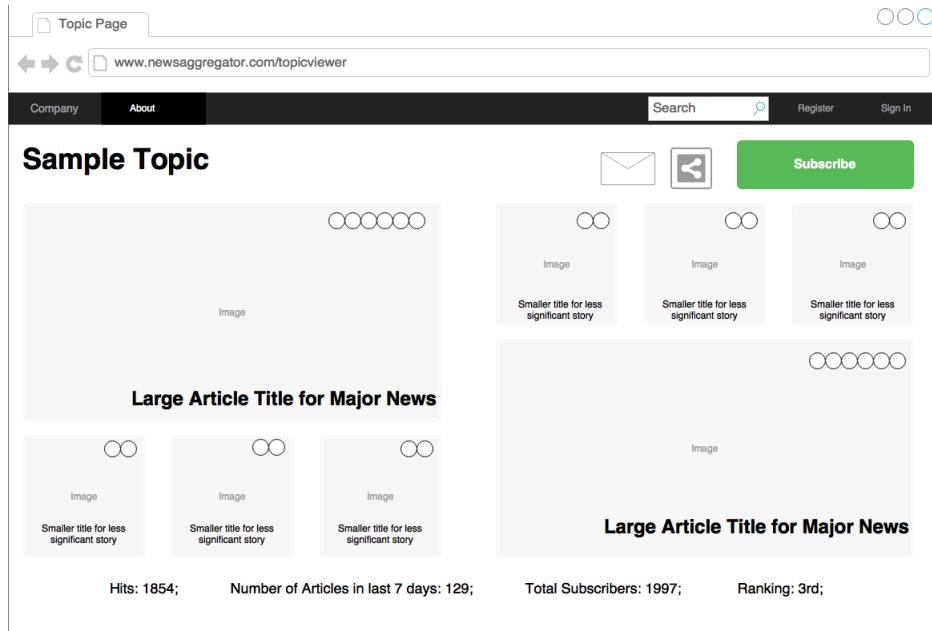


Figure 4.5: The topic viewer is similar in structure to the topic dashboard, but with the addition of statistics and a subscribe button.

This aspect is done using the panelling effect used in other screens. Each news outlet's logo forms the basis of the panel, with the title of the outlet below. These panels form buttons that the user can click on to either select or deselect an outlet. The panel will either have a coloured outer rim, or an embedded effect to indicate selection.

4.2.6 Article Viewer

When the user clicks on an article panel in the topic viewer or topic dashboard they are presented with the summarised article itself. This screen is modelled on standard news interfaces, and produces the headline, image and article body on the left hand side of the page. On the right hand side are links to the original articles that the summary was generated from, and links to other articles that are in the same topic.

Like most other screens in the application, there is also a share button on the top right, and an email button.

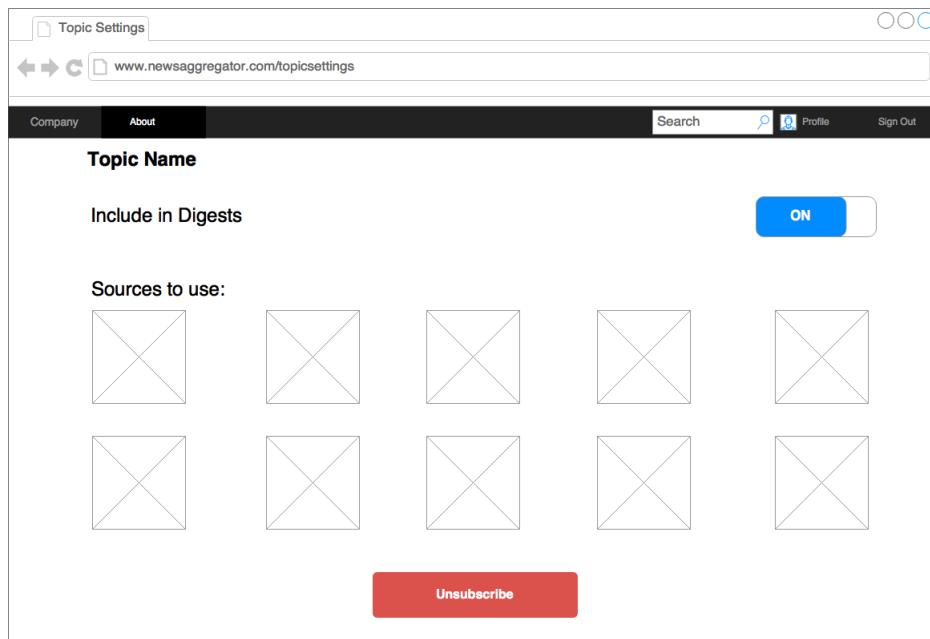


Figure 4.6: The Topic Settings page allows users to specify key preferences about the topic.

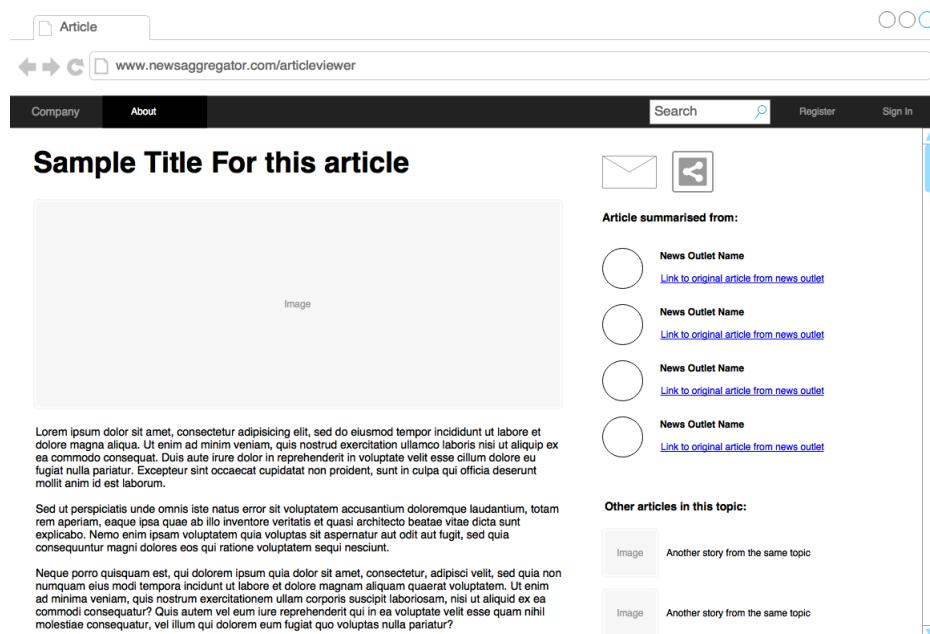


Figure 4.7: The Article Viewer is modelled on standard news outlets' interfaces.

4.2.7 Topic Lists

The Topic List page is accessed by clicking on the profile button when logged in to the application. A panel appears, listing two choices - My Topics, and My Settings. This page is presented on the selection of the My Topics button.

The aim of the topic list is self-explanatory. It lists all topics that a user is subscribed to. Upon clicking on a topic in the list, a user would be taken to the Topic Settings page for that topic.

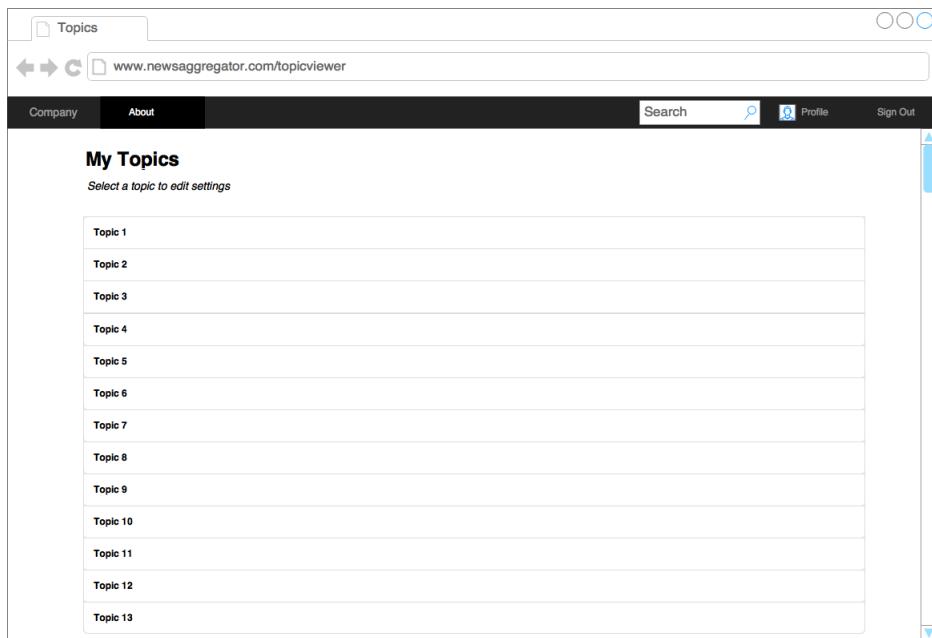


Figure 4.8: The Topic List page allows users to select and adjust the settings of a topic.

4.2.8 Profile Settings

The other option upon clicking on the profile button in the navigation bar of the application is to access the Profile Settings page. Here, the user is able to change their profile picture using the button in the top half of the page. Also present on this page are the user's digest settings. Here the user can set whether they want to receive a news digest in the morning, the afternoon, at both times, or not at all. They can also specify the email address that they wish the news digest to be sent to.

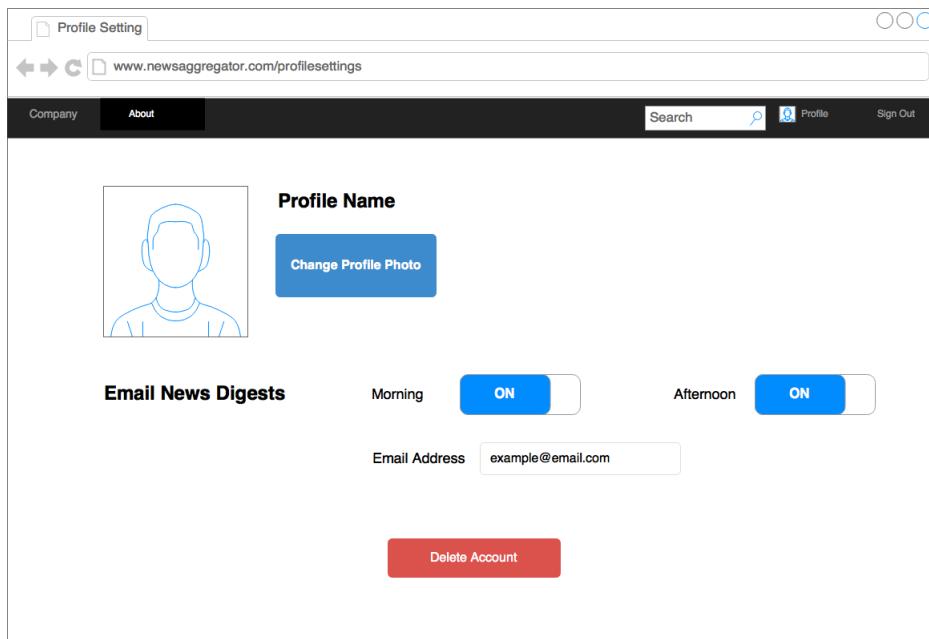


Figure 4.9: The profile settings page allows the user to adjust their profile picture and digest settings.

4.3 Back End Flow Diagram

In Figure 4.10 the expected flow of the machine learning aspects of the back end is presented. The process of an article being summarised is as follows:

1. The **Article Fetcher** queries the various News APIs used for new articles.
2. The **Article Curator** takes the response from the News APIs and performs any scraping that may be necessary, before passing the article to the Machine Learning phases.
3. The article is then passed to the **Topic Modelling** phase, which identifies which topics it consists of.
4. The **Topic Labelling** section takes the results of the Topic Modelling phase, and proceeds to label its topic
5. **Clustering** then occurs. The article, along with its topic label is passed to

this phase, and the program pulls articles from the database that have already been assigned the same topic. These articles are then clustered, and the cluster containing this article is passed to the next phase.

6. Now that the cluster has been passed, **Summarisation** occurs on the article, which is then put in the Summarised Articles database, ready to be called for a user to read.

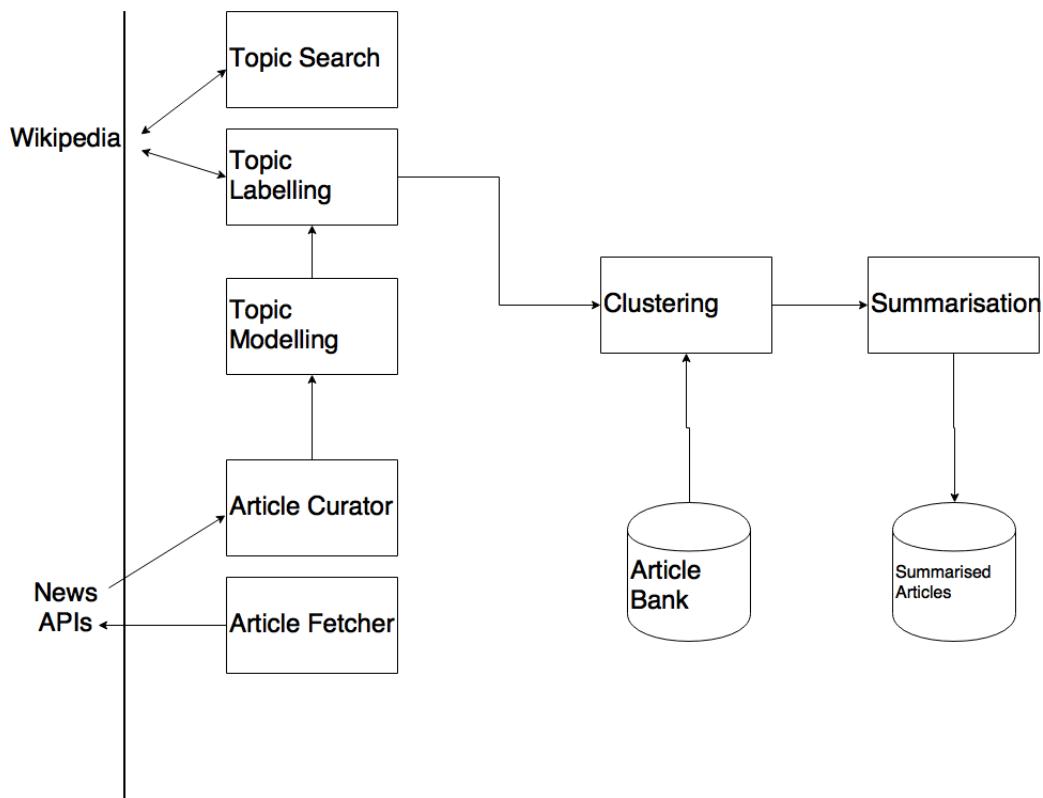


Figure 4.10: The expected flow of the Back End.

4.4 Language and Platform Choices

4.4.1 Front End

The survey that was conducted in Section 2.1.2 clearly presented the result that a website would be the preferred option for the application, followed by a mobile

application in either iOS or Android. Therefore I've decided to follow this and create a WebApp as a primary platform. If time permits, I'll then develop the application for a second platform, which will be a mobile platform.

For the WebApp I will primarily use Bootstrap for the Interface. One design framework I used to great effect during my Industrial Placement in Summer 2016 was React and Redux. React is a JavaScript library developed by Facebook that is designed for building scalable and reusable user interfaces. Redux is a design pattern for JavaScript that is designed around a central predictable state container that is the single source of truth for a User Interface.

The benefit of this framework is that it allows for very clean code in the user interface. This is because the concept of a single source of truth allows for data flow that is purely in only one direction, making it easier to understand.

4.4.2 Back End

The back end will be constructed using Java. The reasoning behind this is that I have experience from my Industrial Placement using Java as the basis of a server and a back end.

4.5 Infrastructure

4.5.1 Hardware

One of the first key decisions I made about the infrastructure surrounded the hardware that the server would be deployed on.

Amazon Web Services

Amazon Web Services (AWS) is a set of cloud computing platforms provided by Amazon. A key benefit to it is that it has a free tier, which allows for a subset of their services to be used without charge.

Some of the key platforms offered by AWS include:

- **EC2** allows a user to create a single virtual machine for free on which to run software. This would be with a limit of 750 hours of use per month. In my case this would be used for running the back end of the News Aggregator.

- **S3** is a storage system provided by Amazon that allows a user 5GB of space for free.

Digital Ocean

Digital Ocean is an infrastructure provider that provides ‘cloud computing for developers’. Digital Ocean offers ‘droplets’ that can act as virtual machines, on which I would deploy my backend. There are no free tiers, but the student pack provided by GitHub would cover the first two months of charges.

In addition to this, one of the advantages of Digital Ocean is that if resizing of the virtual machine is necessary, it can be done in a matter of minutes, thus making the final solution very highly scalable.

In contrast though, unlike Amazon, Digital Ocean don’t provide alternative microservices such as S3. Having said that, it could be argued that the droplet itself would offer that with its storage capabilities. In fact, the cheapest droplet offered by Digital Ocean (priced at \$5 a month), entails a virtual machine with a 20GB SSD disk, which is four times the amount provided by the free tier in AWS.

4.5.2 Database

MongoDB

The first database solution that I researched was MongoDB, as I have had prior experience with it.

MongoDB is a document database that is designed to be both flexible and highly scalable. Documents are written in a format that is highly similar to JSON. This format is what results in a Mongo database being flexible, as the schema of documents in the collection can be changed quickly.

MongoDB allows for the connection of items between tables through linking. Auto-generated ‘ObjectId’s (that are generated based on the timestamp) allow for easy referencing of an item in another.

A MongoDB database would have to remain in a virtual machine on either Amazon or Digital Ocean, and will therefore have to count against the storage provided by that virtual machine.

DynamoDB

DynamoDB is a No-SQL database provided by Amazon Web Services as a separate micro-service to EC2 and S3.

The structure of a document in DynamoDB is more rigid than in MongoDB. It is primarily designed as a key-value system, with a primary key. Searching in the database is driven by this primary key, or by ‘secondary indexes’. These secondary indexes need to be declared when the table is first created. Therefore this means that if there needs to be a schema change, the table needs to be recreated.

A DynamoDB database doesn’t need to be stored in the virtual machine, but there are trade offs. Whilst Amazon Web Services provide 25GB of storage for the database as part of the free tier, there are throttles on the throughput of the database. As part of the free tier, a user is allowed to allocate 25 ‘units’ of throughput to both reading and writing from a database, where one unit corresponds to transfer of 1KB per second. Beyond these 25 units, Amazon would start charging.

4.5.3 Infrastructure Decisions

Comparison table

Table 4.1 shows the key points of comparison between the four main options for hardware and database combinations.

Decision Process

I tried all possible combinations listed in table 4.1 before finally arriving at the decision of using Digital Ocean and MongoDB.

I initially started with the configuration with AWS and MongoDB. However, it became clear quite quickly that there was a distinct possibility that the memory provided by AWS would not be sufficient to hold a large number of articles (about six weeks or greater) in the database. There certainly would not be scope for adding extra outlets near the end of the project.

In an attempt to counter this, and trying to avoid solutions that would require the spending of money, I then switched the database option to DynamoDB. This worked well for a time, but as the database grew, it became clearer that the throttling of the read and write operations to my database was going to cause major issues down the line for my background jobs on the server, as well as for simple API methods. In order to counteract this, I resolved that I would use vouchers for AWS from the GitHub student pack in order to fund some extra read-write operations per second.

Infrastructure	Advantages	Disadvantages
AWS and MongoDB	<ul style="list-style-type: none"> • AWS provide free services • Mongo provides a flexible schema, meaning that the table structure can be changed further down the line. 	<ul style="list-style-type: none"> • Mongo would need to be hosted on the EC2 instance, or on S3, which would limit the free storage to 5GB.
AWS and DynamoDB	<ul style="list-style-type: none"> • DynamoDB would provide 25GB of storage on its own, meaning that I wouldn't have to use the S3 storage. 	<ul style="list-style-type: none"> • DynamoDB's key-value schema is quite rigid, meaning that there isn't much flexibility on the model. • The throttling of reads and writes by DynamoDB could severely impact the performance of key back end tasks.
Digital Ocean and DynanoDB	<ul style="list-style-type: none"> • Digital Ocean offers the ability to resize the server as needs be within a matter of minutes. 	<ul style="list-style-type: none"> • As before, there are some key issues regarding the speed of DynamoDB whilst on the free tier. • This solution would involve maintenance of a database solution provided by one company, and a server on another, which is potentially quite wasteful.
Digital Ocean and MongooDB	<ul style="list-style-type: none"> • MongoDB database can be stored compactly on the Digital Ocean droplet 	<ul style="list-style-type: none"> • Whilst covered by vouchers for the first two months of use at least, the Digital Ocean instance would require expenditure if used for longer than that period, versus AWS EC2 which would be free for at least twelve months. • The flexibility of the MongoDB model would mean that fields can be removed and added easily from the schema, allowing for further changes down the line.

Table 4.1: Comparing the advantages and disadvantages of potential combinations of hardware and database options.

Having resolved to do this, I then decided to make use of Digital Ocean. This was because I wanted to experiment to see the effects of having an instance with extra RAM would do for the running of my program, and knew that the vouchers on the GitHub student pack would easily cover a Digital Ocean droplet until the end of the project. What I found was that my background jobs (such as pulling in articles, labelling, clustering and summarisation) were being performed faster in general, and so I decided to permanently move to Digital Ocean.

In the end, a final nail in the coffin for my use of Amazon Web Services came when I wanted to change the way my articles table worked in the database. I had wanted the ability to search by date published (as I was considering at the time removing articles published before a certain date), but found that I would have to recreate my entire table in order to add this functionality to my table. As a result, and combined with the fact I now had a much larger amount of storage space on my new Digital Ocean droplet, I took the decision to make a new MongoDB database, knowing that it was still early enough in the development process that losing all my previous database work would not affect the finished product, or the development timeline.

5 Implementation

5.1 Database Schema

5.1.1 Schema Diagram

Figure 5.1 shows the basic components of the various collections in the Mongo Database, along with some of the connections involved. Further details on each are given in each of the following sections.

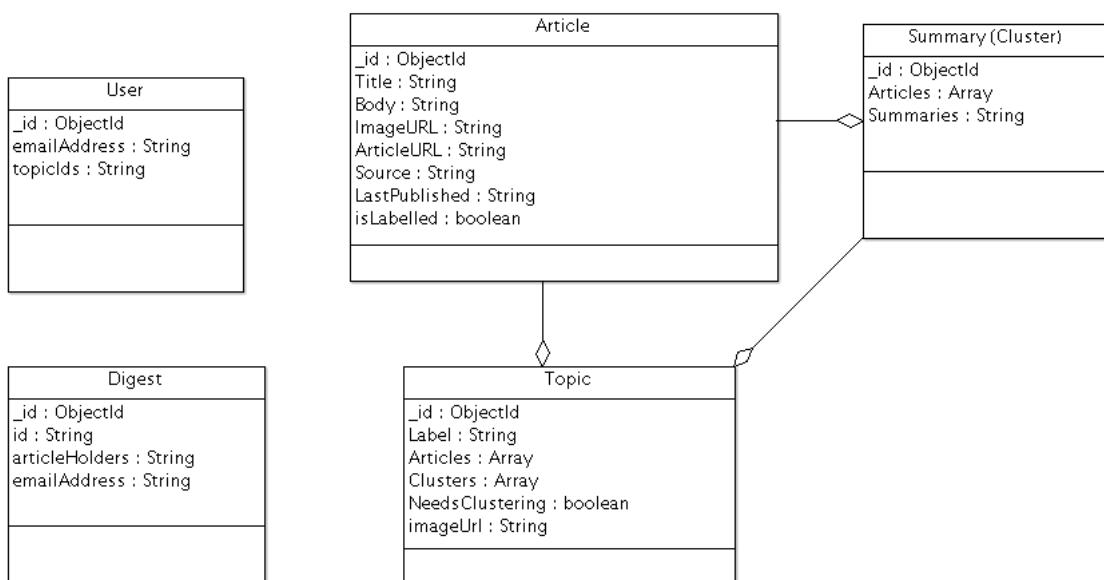


Figure 5.1: A diagram of the core MongoDB schema.

5.1.2 Articles

A sample document for an article can be seen below:

Listing 1: A sample document in the Article table

```

1  {
2      "_id" : ObjectId("591df7cfacea820abe29709f"),
3      "Title" : "Kirsten, Simons on five-man panel in hunt
       for new SA coach",
  
```

```

4     "Body" : "Cricket South Africa has nominated a five-man
      panel including two former national coaches, Gary
      Kirsten and Eric Simons, to recommend a suitable
      candidate for the position of head coach, which it
      aims to fill by the beginning of September
5     ...
6     ",
7     "ImageURL" : "http://www.espnccricinfo.com/db/PICTURES/
      CMS/157600/157626.5.jpg",
8     "ArticleURL" : "http://www.espnccricinfo.com/southafrica
      /content/story/1098337.html",
9     "Source" : "espn-cric-info",
10    "LastPublished" : "2017-05-18T10:42:38Z",
11    "isLabelled" : true
12 }
```

The purpose of the article table is for storage of raw article data that has been brought in by the article fetcher. Most fields are self explanatory. The *isLabelled* field is defaulted to false when an article is first brought in, and is then changed to true when the Topic Labeller has labelled it. This is present so that articles that either weren't labelled, or caused an error when being labelled, can be found with a simple search in the database.

5.1.3 Topics

A sample document for a topic can be seen below:

Listing 2: A sample document in the Topics table

```

1 {
2     "_id" : ObjectId("591df830acea820aebf55465"),
3     "Label" : "Opposition (politics)",
4     "Articles" : [
5         ObjectId("591e172dacea820aebf55a74"),
6         ObjectId("5925f959acea821a75b831bc"),
7         ObjectId("59281cddacea821ab611d3ed"),
8         ObjectId("59291a47acea82594af9ef53"),
9         ObjectId("592aa5daacea8264125210d3")
10    ],
11    "Clusters" : [
12        ObjectId("591f8bc7acea821e031f4dc6")
```

```

13   ],
14   "NeedsClustering" : true,
15   "imageUrl" : "https://upload.wikimedia.org/wikipedia/
      commons/d/d1/Stand_in_opposition_city_hall_boston.
      jpg"
16 }
```

A topic object is a core concept in the back end. It keeps references to all articles that have been labelled with that topic, and also holds references to all the clusters (summaries) that are produced from these articles. The *imageUrl* field holds either the lead image for the corresponding Wikipedia article, or a placeholder if the Wikipedia article doesn't have an image.

The *NeedsClustering* flag performs a similar function to the *isLabelled* flag in the Articles table. It is defaulted to true, and is changed to false once all the clustering (but not necessarily summarisation) is completed. It is used again to find topics that may have 'slipped through the cracks', either due to an error, or for some other reason.

5.1.4 Clusters

A sample document for a cluster can be seen below:

Listing 3: A sample document in the Summaries table

```

1 {
2   "_id" : ObjectId("591e47f7acea8221dcb97c42"),
3   "Articles" : [
4     ObjectId("591df7cfacea820abe2970a9"),
5     ObjectId("591df7cfacea820abe2970b2")
6   ],
7   "Summaries" : "{\"[daily-mail, the-guardian-uk]\": [{\"
      sentence\": \"President Donald Trump says he's 'very
      close' to naming an FBI director.\", \"
      sentencePosition\": 0.0, \"absoluteSentencePosition\":
      0, \"identifier\": 0, \"relatedNodes\": [], \"source\": \"
      daily-mail\"}, {\"sentence\": \"Trump confirmed
      Thursday afternoon that former Democratic Sen. Joe
      Lieberman is a favorite for the position.\", \"
      sentencePosition\": 0.019230769230769232, \"
      absoluteSentencePosition\": 1, \"identifier\": 1, \"
      identifier2\": 1}]} }
```

```

relatedNodes\":[],"source\":\"daily-mail\"},{"
  sentence\":\"We're very close to an FBI director,'
  Trump told reporters shouting questions at him as he
  met with the Colombian president in the Oval Office
  .\", "sentencePosition":0.038461538461538464,"
  absoluteSentencePosition":2,"identifier":2,"
  relatedNodes\":[],"source\":\"daily-mail\"},{"
  sentence\":\"Soon,' he added.\", "sentencePosition"
  :0.057692307692307696,"absoluteSentencePosition"
  :3,"identifier":3,"relatedNodes\":[],"source\":
  \"daily-mail\"}

8     ...
9   ]}"
10 }

```

An object in this table represents a Summary, and corresponds to the Clusters array in a Topic. It contains an array of references to the original articles used for the summarisation.

The *Summaries* field is a string representation of a map object, with the values representing to a summary, and the corresponding key representing a combination of articles used to make that summary (represented using their sources). For example, in the full version of the sample document, there would be three keys in the string representation:

- [daily-mail,the-guardian-uk], for a summary using sources from *The Daily Mail* and *The Guardian*.
- [daily-mail], for a summary using sources from *The Daily Mail* and *The Guardian*.
- [the-guardian-uk], for a summary using sources from *The Daily Mail* and *The Guardian*.

These represent all the possible permutations (of size greater than zero) of the articles in the cluster. Its primary use is in the customisation offered to users, that can be seen in more detail in section 5.8.

5.1.5 Users

A sample document for a User can be seen below:

Listing 4: A sample document in the Users table

```

1 {
2   "_id" : ObjectId("592c5288acea82147ca7f0d1"),
3   "emailAddress" : "fake@gmail.com",
4   "topicIds" : "[{\\"topicId\\":\\"591e3c74acea821e2b7e6c02
5   \\",\\\"sources\\":[\\\"wikipedia\\\",\\\"business-insider-uk
  \\",\\\"daily-mail\\\",\\\"espn-cric-info\\\",\\\"metro\\\",\\"
  \"mirror\\\",\\\"newsweek\\\",\\\"sky-sports-news\\\",\\\"the-
  telegraph\\\",\\\"the-times-of-india\\\",\\\"bbc-news\\\",\\"
  \"bbc-sport\\\",\\\"bloomberg\\\",\\\"cnn\\\",\\\"cnbc\\\",\\\"espn\\\",
  \\\"four-four-two\\\",\\\"the-washington-post\\\",\\\"the-wall-
  street-journal\\\",\\\"associated-press\\\",\\\"the-
  guardian-uk\\\"]},\\\"digests\\\":false}]"
}

```

A user object is very simple. The only field that is not necessarily self-explanatory is the *topicIds* field.

This is an array of strings, with each string forming a representation of a user's topic subscription (the topic being represented by the field marked *topicId*), and the settings for that topic. The two core settings for each subscription is the *digests* setting (a boolean flag indicating whether a user wants to receive email digests that include this topic), and a list of sources that the user would like their news to default to (further information on this feature can be seen in section 5.8). By default, when a user subscribes to a topic, *digests* is set to false, and *sources* is represented by a list comprising all possible sources.

5.1.6 Digests

A sample document for a digest can be seen below:

Listing 5: A sample document in the Digests table

```

1 {
2   "_id" : ObjectId("591ec5c433625db581d77cc6"),
3   "id" : "591ec5c433625db581d77cc6",
4   "articleHolders" : "[{\\"topicId\\":\\"593020de5826a72aa69
  7a201\\\",\\\"articleId\\\":\\"59322bc5acea820dd0c7ef37\\\",
  \\\"title\\\":\\\"Elon Musk 'intrigued' by India's
  objective of all-electric cars by 2030 - Times of
  India\\\",\\\"imageUrl\\\":\\\"http://timesofindia.

```

```

1     indiatimes.com/photo/msid-58964385/58964385.jpg?5372
2\", \"lastPublished\": \"2017-06-03T08:40:00Z\", { \
3     topicId\": \"5930210d5826a72aa697a227\", \"articleId\"
4     : \"59324084acea820dd0c7f65b\", \"title\": \"Kathy
5     Griffin loses ALL of her tour gigs in wake of Trump
6     scandal\", \"imageUrl\": \"http://i.dailymail.co.uk/i/
7     pix/2017/06/03/05/410C672100000578-0-image-a-41_1496
8     464334669.jpg\", \"lastPublished\": \"2017-06-03T04:39
9     :14Z\", { \"topicId\": \"593020de5826a72aa697a201\", \
10    \"articleId\": \"59323dddacea820dd0c7f566\", \"title\": \
11    \"Putin says hacking of Democratic Party may have
12    been CIA false flag op\", \"imageUrl\": \"http://i.
13    dailymail.co.uk/i/pix/2017/06/03/05/410C6AD500000578
14    -0-image-a-92_1496463783859.jpg\", \"lastPublished\":
15    \"2017-06-03T04:28:04Z\" }
16
17     ...
18
19     },
20
21     "emailAddress" : "specialk109@gmail.com"
22 }
```

A digest document represents a single daily digest for a user. Thus, a user who gets daily digests would have a single object for each day. Having said a user would have a single object for each day, it's worth pointing out that the digests are not linked to the User object. Whilst there is no concrete justification for this to be the case either way, a reason for the digests table and users table to not be linked in some way is that they simply don't need to be. A user is never called when a digest is retrieved, and similarly a digest is never called when a user is retrieved.

The *articleHolders* field in the document above represents the articles in a digest. It is an array of strings, with each string a JSON representation of an object containing a topic id, article id, and basic information about the cluster (title, image url, and a publishing date). The reasoning behind having the topic id in addition to the cluster is due to the pattern matching employed in the URL for the article viewer on the website (this is explained in further detail in section 5.8).

5.2 Potentially Useful APIs

5.2.1 News Outlets

Readership statistics for selected outlets in the United Kingdom

Statistics below are taken from the *Digital News Report 2016*[23] by the Reuters Institute for the Study of Journalism[24]. Viewership is given as a percentage of people who say they use the outlet at least weekly.

Outlet	Viewership	Political Stance
BBC	51%	Neutral
Mail Online	17%	Right
Huffington Post	14%	Left
The Guardian	14%	Left
Sky News Online	11%	Neutral
The Telegraph Online	9%	Right
The Independent	6%	Centrist

Table 5.1: Readership statistics for selected outlets in the United Kingdom

The Guardian[28]

The Guardian is a UK newspaper that commands about 14% of the viewership as outlined in Table 5.1. It is also known for having an API that allows a user full and free access to all articles, including the article text itself. The API can access any article dated on or after 1999, and as long as not used for a commercial purpose, is completely free to use.

To obtain an article from *The Guardian* there are two calls that need to be made. First, we make a call to the search endpoint, which returns a JSON object similar to this:

Listing 6: A sample response to an API call to The Guardian

```

1 {
2   "response": {
3     "status": "ok",
4     "userTier": "free",
5     "total": 1,
6     "startIndex": 1,
7     "pageSize": 10,
8     "currentPage": 1,
9     "pages": 1,
10    "orderBy": "newest",
11    "results": [
12      {
13        "id": "1234567890",
14        "url": "https://www.theguardian.com/technology/2018/jan/12/the-best-new-smartphones-for-2018-reviewed",
15        "title": "The best new smartphones for 2018, reviewed",
16        "author": "By Matt Bell Last updated on 12 Jan 2018 at 12:00 UTC",
17        "published": "2018-01-12T12:00:00Z",
18        "content": "The best new smartphones for 2018, reviewed",
19        "image": "https://i.guim.co.uk/img/media/2018/01/12/1510753333535/100_1000x500_q85_box_0.jpg?w=1000&h=500&t=1510753333535&s=1510753333535&r=0&f=false&p=0&q=95&a=0&v=1510753333535&u=https://www.theguardian.com/technology/2018/jan/12/the-best-new-smartphones-for-2018-reviewed",
20        "type": "article"
21      }
22    ]
23  }
24}
```

```

13         "id": "politics/blog/2014/feb/17/alex-
14             salmond-speech-first-minister-scottish-
15                 independence-eu-currency-live",
16             "sectionId": "politics",
17             "sectionName": "Politics",
18             "webPublicationDate": "2014-02-17T12:05:47Z
19                 ",
20             "webTitle": "Alex Salmond speech - first
21                 minister hits back over Scottish
22                     independence - live",
23             "webUrl": "https://www.theguardian.com/
24                 politics/blog/2014/feb/17/alex-salmond-
25                     speech-first-minister-scottish-
26                         independence-eu-currency-live",
27             "apiUrl": "https://content.guardianapis.com
28                 /politics/blog/2014/feb/17/alex-salmond-
29                     speech-first-minister-scottish-
30                         independence-eu-currency-live"
31         }
32     ]
33 }
34 }
```

From this, we extract the apiUrl, and call this, with the API key in the request headers. This will return the raw text for the article.

News API

A major issue is that most newspapers (for example, *The New York Times*[29]) only give the first paragraph to an article in their API, and don't give out their full articles as standard. There are solutions, such as Webhose[32], that offer feeds into this outlets, but at either a cost, or with a very limited number of requests per month.

News API[21] is a company that provides an API to get headlines for over seventy different news outlets. A call returns data in the JSON format as follows:

Listing 7: A sample response to an API call to News API

```

1 {
2 "status": "ok",
3 "source": "the-next-web",
```

```
4 "sortBy": "latest",
5 "articles": [
6 {
7     "author": "Abhimanyu Ghoshal",
8     "title": "Are these ridiculous headphones the way
9         forward for music tech?",
10    "description": "All the amazing tech we now have at
11        our disposal for enjoying music is closing us
12        off from other people instead of bringing us
13        together. Is there hope yet?",
14    "url": "https://thenextweb.com/gear/2017/01/26/can-
15        music-tech-make-us-sociable-again/",
16    "urlToImage": "https://cdn3.tnwcdn.com/wp-content/
17        blogs.dir/1/files/2017/01/Vinci-hed-1.jpg",
18    "publishedAt": "2017-01-26T13:12:10Z"
19 }
20 ]
21 }
```

We can use the url provided for each of these results to get to the webpage of the corresponding article. However at this stage I'll need to create a scraper to extract the raw content of the article. There could be copyright implications for this step, even though the final production summary won't look the same as the raw content. These implications will be fully examined in the final report.

News API counts amongst its 70 sources the following:

- BBC News
- Mail Online
- *The Telegraph*
- *The New York Times*
- Associated Press
- *The Independent*
- *The Daily Mirror*

News API also has *The Financial Times* and *The Guardian*, although these are left off the list above as they have their own fully accessible APIs.

Wikipedia

Wikimedia (the parent company for Wikipedia[36]) provides documentation for the Wikipedia API on its website MediaWiki website. Here, I can search using the search term that a user has used (perhaps in searching for a topic), and get a result that is as follows:

Listing 8: A sample response to an API call to Wikipedia's API

```

1 {
2     "query": {
3         "searchinfo": {
4             "totalhits": 4152
5         },
6         "search": [
7             {
8                 "ns": 0,
9                 "title": "Albert Einstein",
10                "snippet": ""<span class=\"searchmatch>Einstein</span>" redirects here.  
For other uses, see <span class=\"searchmatch>Albert</span> <span class  
=>Einstein</span> (disambiguation) and <span class=\"  
searchmatch>Einstein</span> (disambiguation). <span class=\"  
searchmatch>Albert</span> <span class  
=>Einstein</span> (/?alb?  
rt ?a?n?ta?n/; German:",  
"size": 124479,  
"wordcount": 13398,  
"timestamp": "2015-05-10T12:37:14Z"
14 },
15 ...
16 }
```

The search query can be altered to provide different information, such as the URL for the lead image, which can be used to show previews of various items to users. The Wikipedia API can also be used in the topic labelling aspect of the Machine Learning techniques (see Section 3.1.2).

5.3 Potentially useful libraries

5.3.1 Mallet

Mallet (Machine Learning for Language Toolkit) is a java library that provides an interface for various natural language-based machine learning tasks. These tasks include:

- Document Classification
- Sequence Tagging
- Clustering
- Topic Modeling
- Information Extraction

Mallet in Topic Modeling

Mallet is particularly useful for Topic Modeling. Mallet uses a corpus in the form of a list of strings in order to train a set of topics. The library trains the topics using Gibbs Sampling and Latent Dirichlet Allocation.

Mallet also provides an inferencer that allows a user to submit a new document and receive a list of probabilities. Each of these probabilities shows the likelihood that the corresponding topic is an accurate classification of the document we are testing. Thus the user can then infer the actual topics that the document is about by simply taking the topics with the highest probabilities.

5.3.2 Natural Language Processing

There are several Natural Language Processing libraries available. The following are designed for use in Java:

Apache OpenNLP

OpenNLP is an open source library developed by Apache. The aim of the project is to support the basic aspects of natural language processing.

OpenNLP supports all of the basic tasks that were specified in section 3.3.1 with varying degrees of success, plus ‘document categorisation’.

Document Categorisation, as the name suggests, takes a document and categorises it. This could potentially be used in the topic modelling phase, but has the disadvantage that it doesn't come with a model for training, and the user of the library has to create their own. This is a contrast to the Mallet library, which can create a model when the user passes in the raw data.

When it comes to the other key tasks of Natural Language Processing, OpenNLP provides flexibility on models. They provide a default model in a variety of languages, including English. However, there is at least one different model file for each task, so space constraints need to be considered.

Name Finder

For using the name finder, there are different model files, depending on what types of proper nouns a user is looking for. These are:

- Date
- Location
- Money
- Organisation
- Percentage
- Person
- Time

Space Constraint Table

Table 5.2 shows each model and their respective sizes. As Apache's recommended method to bring the files in to the program is through the Java class `FileInputStream` their sizes can be important, as they make use of the Java heap.

Model Name	Task	Size (MB)
en-chunker	Chunking	2.6
en-ner-date	Name Finder: Dates	5.0
en-ner-location	Name Finder: Locations	5.1
en-ner-money	Name Finder: Money	4.8
en-ner-organization	Name Finder: Organisations	5.3
en-ner-percentage	Name Finder: Percentages	4.7
en-ner-person	Name Finder: People	5.2
en-ner-time	Name Finder: Time	4.7
en-pos-maxent	POS Tagging	5.7
en-sent	Sentence Detection	0.01
en-token	Tokenisation	0.44
Total if all used		43.55

Table 5.2: A table showing the space required for the various models provided for use with OpenNLP

Stanford CoreNLP

The Stanford CoreNLP is a toolkit provided for Java by the Stanford Natural Language Processing Group. Like Apache OpenNLP, CoreNLP also offers the same functionalities that were discussed in section 3.3.1. However, Stanford go above and beyond this list in what they can provide. Their coreference and name recognition solutions are ‘competition winning’, whilst they also provide functionality for semantic analysis of text.

As a result, CoreNLP could have a major benefit for abstractive summarisation, through its coreference resolution.

Extracting elements from a document

The following is a code snippet, taken from the Stanford CoreNLP documentation that demonstrates how to perform certain Natural Language tasks from a given annotated document.

Listing 9: Analysing an annotated document using Stanford’s CoreNLP

```

1
2 // these are all the sentences in this document
3 // a CoreMap is essentially a Map that uses class objects
   ↵ as keys and has values with custom types

```

```

4 List<CoreMap> sentences =
    ↪ document.get(SentencesAnnotation.class);
5
6 for(CoreMap sentence: sentences) {
    // traversing the words in the current sentence
7     // a CoreLabel is a CoreMap with additional
    ↪ token-specific methods
8     for (CoreLabel token:
9         ↪ sentence.get(TokensAnnotation.class)) {
10        // this is the text of the token
11        String word = token.get(TextAnnotation.class);
12        // this is the POS tag of the token
13        String pos = token.get(PartOfSpeechAnnotation.class);
14        // this is the NER label of the token
15        String ne = token.get(NamedEntityTagAnnotation.class);
16    }
17
18    // this is the parse tree of the current sentence
19    Tree tree = sentence.get(TreeAnnotation.class);
20 }
21
22 // This is the coreference link graph
23 // Each chain stores a set of mentions that link to each
    ↪ other,
24 // along with a method for getting the most representative
    ↪ mention
25 // Both sentence and token offsets start at 1!
26 Map<Integer, CorefChain> graph =
    document.get(CorefChainAnnotation.class);

```

WordNet

WordNet is a large scale database developed by Princeton that offers a lexical database of the English language. It can primarily be used as a thesaurus, and could therefore be useful in the natural language generation of abstractive summarisation, as well as in finding similarities between sentences.

5.4 Machine Learning

5.4.1 Topic Modelling

For the Topic Modelling task I used the Mallet library, which is described in full detail in section 5.3.1.

There are two sections to the topic modelling phase: Training, and Estimating.

Training

For training the model, I used the ParallelTopicModel class provided by the Mallet library. The class provides a parallel implementation of Latent Dirichlet Allocation. There are several key parameters that need to be passed in to the model before it can be estimated:

- **numTopics**, the number of topics to model the corpus into.
- **alphaSum** is distributed evenly across all topics. For example, an alphaSum of 1 over 100 topics would result in each topic having an alpha score of 0.01. An alpha score is viewed as a smoothing term. It ensures that the probability of a given topic being in a document is never zero.
- **beta** is similar to alpha, and is the corresponding smoothing factor ensuring that the probability of a word being in a given topic is never zero.
- **numIterations** is the number of iterations the trainer should run through before presenting a final model.

Initially I attempted to model topics by feeding in articles and seeing what topics came up. However, it was clear that this wasn't ideal, as frequently occurring words, such as 'a', 'the' and 'and' were too prominent in the list of words in each topic, and thus were skewing any results from the estimation stage.

I then tried to use the provided 'stoplists' file. With this file, Mallet aims to filter out common words with no real significance out of the given texts, before training topics. It became noticeable quickly though that the stoplists file wasn't extensive enough, as it was missing contracted words, such as 'it's'. Given the frequency of these words, I dismissed this idea and set about a new solution.

As a final solution, I decided to remove all insignificant words by removing all non-nouns from the article bodies. I did this using the Apache OpenNLP library, and found that this made (to the naked eye) the topic models look a lot better.

My final parameters can be seen in table 5.3.

Parameter	Final Value
alphaSum	1.0
beta	0.01
numTopics	500
numIterations	1500

Table 5.3: A table showing the final parameters chosen for the topic modelling training

Initially, the articles would be used to retrain the model each time a set of articles were about to be labelled. However, in the optimisation phase of the project, I changed this so that the training would happen once a week, and a model saved on file. Full details and justification can be found in section 6.1.1.

Estimation

The body of the article to be estimated is passed in to the model. This was initially the whole article, but after I decided to keep only nouns for training, I did the same for estimation.

There are three key parameters for estimation. These are:

- **numIterations**, which is the number of times to iterate through the model before producing a final estimate.
- **thinning**, which is the number of iterations to run before saving an interim model.
- **burn-in** is the number of iterations to perform before beginning the estimation. The reasoning behind this is that since the first iteration begins with a random distribution, the first set of iterations are unlikely to be truly representative of the actual model. The burn-in period could potentially be compared to the warm-up of an athlete.

After the topic model for the new article has been estimated, I prepare the result for topic labelling. To do this, I parse the data into a new custom class I created, that stores the word and the ‘distribution’. The distribution is calculated by the Mallet library, and is an indication of the likelihood that an article belongs to a certain topic.

The final parameters chosen can be seen in table 5.4.

Parameter	Final Value
numIterations	2500
thinning	40
burn-in	50

Table 5.4: A table showing the final parameters chosen for the topic modelling estimation

5.4.2 Topic Labelling

When I first implemented the Topic Labelling phase of the backend, I aimed to emulate the algorithm of Lau, Grieser, Newman and Baldwin as closely as possible. This algorithm is detailed in section 3.1.2. However, once this algorithm had been implemented, it was clear that there were some potential flaws when it applied to my project.

A key issue in the topic labelling phase was that of speed. The algorithm required dozens (and potentially hundreds) of calls to the Wikipedia API. There are ten calls in step two (searching Wikipedia for the original top ten topic terms), followed by calls on every single noun chunk found in step three of the algorithm. Calculating the Related Article Conceptual Overlap scores (RACO, see section 3.1.2 would also require every single ‘outlink’ to be found, and every category of it. Given this needs to be done for every single noun chunk found earlier, it was easy to see why this could become prohibitively costly in terms of speed.

A simple test run using five articles from *The Guardian* of varying lengths resulted in each taking at least fifteen minutes to complete the topic labelling process. Given the potential volume of articles expected per day, I reluctantly dismissed the algorithm as too expensive for the project.

In order to streamline the labelling process I removed major aspects of the algorithm. This included removing the creation of noun chunks, and the RACO calculation. I also looked to some simple heuristics, such as isolating names using the OpenNLP

library and automatically adding them as labels.

Brief overview of final labelling algorithm

1. **Find all proper names** in the article body using the Apache OpenNLP library. I used the models for Locations, Organisations and People in this phase.
2. **Search Wikipedia** for the top fifteen results for each of the ten topic words provided by the modelling algorithm.
3. **Retrieve text from Wikipedia articles** for all results in step two that aren't in the list of names found in step one.
4. **Perform Candidate Ranking.** This step is described in more detail below.
5. **Add top 18 results to candidates from step one.** The number 18 was found via trial and error, with the aim of striking a balance between too many potential labels and too few. Too many could result in a very slow clustering process, as there are too many topics that need clustering, causing a backlog. Too few could reduce accuracy.

This significantly increased the speed of the algorithm, and articles that were taking fifteen minutes or more on the previous algorithm were now taking less than a minute, which I decided would be fast enough for the expected volume of articles.

A full analysis of the speed of the new topic labelling algorithm, and its role in the bigger picture of the summarisation process is outlined in the evaluation section.

Brief overview of candidate ranking algorithm

1. **Isolate the nouns from the Wikipedia articles of every label**
2. **Use nouns from step one as a corpus in a TF-IDF model**
3. **Isolate the nouns from the original article**

4. For each label:

- Calculate the sum of performing TF-IDF on each individual noun from the label's article. The TF-IDF algorithm is explained in full detail in section 3.1.3.
- Find the total number of nouns that are present in the original article that are also present in the candidate article. This is known as the crossover.
- Normalise the crossover by dividing it by the total number of nouns in the candidate article and the original article.
- Multiply the crossover by the sum of TF-IDF operations found earlier to give a final total.

5. Order the candidates by score, highest to lowest.**5.4.3 Clustering****5.5 Summarisation****5.5.1 Extractive Summarisation****5.5.2 An attempt at Abstractive Summarisation****5.6 Restlet****5.7 Server Tasks****5.8 Front End****5.9 Key Classes**

6 Optimisation

6.1 Speed Optimisations

6.1.1 Topic Modelling

6.1.2 Topic Labelling

6.1.3 Clustering

6.2 Memory Optimisations

7 Evaluation

7.1 Machine Learning and Summarisation

7.2 Summary Analysis

7.3 User Interface Evaluation

8 Conclusion

9 Future Work

9.1 Foreign Languages

9.2 Further Optimisations

9.3 Other Apps

References

- [1] *Apple*. www.apple.com.
- [2] *Apple News*. www.apple.com/uk/news/.
- [3] *BBC*. www.bbc.co.uk.
- [4] *Cluster Analysis*. URL: https://en.wikipedia.org/wiki/Cluster_analysis (Retrieved Dec. 29, 2016).
- [5] Gunes Erkan and Dragomir R Radev. ‘LexRank: Graph-based Lexical Centrality as Salience in Text Summarization’. In: *Journal of Artificial Intelligence Research* (December 2004).
- [6] *Flipboard*. www.flipboard.com.
- [7] Pierre-Etienne Genest and Guy Lapalme. ‘Framework for Abstractive Summarization using Text-to-Text Generation’. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics* (June 2011).
- [8] Pierre-Etienne Genest and Guy Lapalme. ‘Fully Abstractive Approach to Guided Summarization’. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics* (July 2012).
- [9] *Google*. www.google.com.
- [10] *Google Forms*. www.forms.google.com.
- [11] *Google News*. www.news.google.com.
- [12] Karl Grieser et al. ‘Using Ontological and Document Similarity to Estimate Museum Exhibit Relatedness’. In: *ACM Journal ACM Journal of Computing and Cultural Heritage* (2011).
- [13] Vishal Gupta and Gurpreet Lehla. ‘A Survey of Text Summarization Extractive Techniques’. In: *Journal of Emerging Technologies in Web Intelligence* (2010).
- [14] *Hierarchical Clustering*. URL: https://en.wikipedia.org/wiki/Hierarchical_clustering (Retrieved Dec. 29, 2016).
- [15] *How does Google News cluster stories?* URL: <https://www.quora.com/How-does-Google-News-cluster-stories/answer/Bharath-Kumar-M?srId=Qord> (Retrieved Jan. 8, 2017).
- [16] Statistic Brain Research Institute. *Attention Span Statistics*. URL: <http://www.statisticbrain.com/attention-span-statistics/>.

- [17] N. R. Kasture et al. ‘A Survey on Methods of Abstractive Text Summarization’. In: *International Journal for Research in Emerging Science and Technology* (November 2014).
- [18] Atif Khan and Naomie Salim. ‘A review on abstractive summarization methods’. In: *Journal of Theoretical and Applied Information Technology* (January 2014).
- [19] *Latent Semantic Analysis*. URL: https://en.wikipedia.org/wiki/Latent_semantic_analysis#Latent_semantic_indexing (Retrieved Dec. 22, 2016).
- [20] Jey Han Lau et al. ‘Automatic Labelling of Topic Models’. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics* (June 2011).
- [21] *News API*. URL: www.newsapi.org (Retrieved Jan. 27, 2017).
- [22] *Ofcam*. www.ofcom.org.uk.
- [23] *Reuters Institute Digital News Report 2016*. Reuters Institute for the Study of Journalism, 2016.
- [24] *Reuters Institute for the Study of Journalism*. www.reutersinstitute.politics.ox.ac.uk.
- [25] *Statistic Brain Research Institute*. www.statisticbrain.com.
- [26] *Summly*. www.summly.com.
- [27] *Term Frequency - Inverse Document Frequency*. URL: https://en.wikipedia.org/wiki/Tf%E2%80%93idf#Example_of_tf.E2.80.93idf (Retrieved Dec. 29, 2016).
- [28] *The Guardian*. www.theguardian.com.
- [29] *The New York Times*. www.nytimes.com.
- [30] *University of Canberra*. www.canberra.edu.au.
- [31] *University of Oxford*. www.oxford.ac.uk.
- [32] *Webhose*. URL: <https://webhose.io> (Retrieved Jan. 27, 2017).
- [33] Harald Weinreich et al. ‘Not Quite the Average: An Empirical Study of Web Use’. In: *ACM Transactions on the Web* (February 2008).
- [34] *What is a good explanation of Latent Dirichlet Allocation?* URL: <https://www.quora.com/What-is-a-good-explanation-of-Latent-Dirichlet-Allocation/answer/Edwin-Chen-1?srid=Qord> (Retrieved Jan. 8, 2017).
- [35] *What's the difference between Latent Semantic Indexing and Latent Dirichlet Allocation*. URL: <https://www.quora.com/Whats-the-difference-between-Latent-Semantic-Indexing-LSI-and-Latent-Dirichlet-Allocation-LDA/answer/Joseph-Turian?srid=Qord> (Retrieved Dec. 22, 2016).

- [36] *Wikipedia*. www.wikipedia.org.
- [37] *WWDC Apple Design Award Winners for 2014*. URL: <http://www.macworld.com/article/2358481/wwdc-apple-design-awards-winners-for-2014.html> (Retrieved Dec. 29, 2016).
- [38] *Yahoo News Digest*. www.uk.mobile.yahoo.com/newsdigest/.
- [39] *YouGov*. www.yougov.com.

Appendices

A Source Code

B API

C User Guide

Index

- Apple, 25
 - Design award, 25
 - News, 23
- Architecture, 37
- Associated Press, 50
- Automatic Labelling of Topic Models, 27
- Back End, 44
- BBC, 13, 47, 50
- Centroid Clustering, 29
- Clustering, 22, 29, 31, 44
 - Centroid, 29
 - Density, 29
 - Hierarchical, 29
 - k-means, 29
- Concept
 - Multimodal Summarisation, 34
- Connectivity, 29
- Daily Mirror, The, 50
- Density Clustering, 29
- Dependency Tree, 32
- Dice's coefficient, 28
- Digital News Report, 11, 13, 14, 46
- Domain Ontology, 34
- Evaluation, 17
- Financial Times, The, 48, 50
- Flipboard, 22
- Front End, 37
- Gibbs Sampling, 26
- Google, 13, 27
 - Forms, 14
 - News, 21, 22
 - Search Engine, 13
- Guardian, The, 47, 49, 50
- Hierarchical Clustering, 29
 - Agglomerative, 29
 - Divisive, 29
- Huffington Post, 47
- Independent, The, 47, 50
- Latent Dirichlet Allocation, 26
- Latent Semantic Analysis, 26
- Latent Semantic Indexing, 26
- LDA, 26
- LexRank, 31
- LSI, 26
- Machine Learning, 26
- Mail Online, 47, 50
- New York Times, 13, 49, 50
- News
 - Outlets, 46
- News
 - Personalised, 12
- News Aggregator, 17, 18
 - Reasons for use, 11
- News API, 49
- News Digests, 19
- Not Quite the Average: An Empirical Study of Web Use, 13
- Ofcom, 13
- PageRank, 31
- React, 46
- Redux, 46
- Related Article Conceptual Overlap, 27, 28
- Reuters Institute for the Study of Journalism, 11–13, 46
- Rich Semantic Graph, 35
- Sky News, 47
- Social Media, 11
- Statistic Brain Research Institute, 13
- Summarisation, 13, 31, 45

- Abstractive, 31, 32, 34
- Extractive, 31, 35
- Information Item, 35
- Multimodal Semantic, 34
- Ontology based, 34
- Rule based, 33
- Semantic Graph, 35
- Tree based, 32
- Summarised, 17
- Survey, 45
- Telegraph, The, 47, 50
- Term Frequency-Inverse Document Frequency, 30, 31
 - Inverse Document Frequency, 30
 - Term Frequency, 30
- TextRank, 31
- TF-IDF, 30, 31
- Topic Labelling, 27, 31, 44
- Topic Modelling, 26, 31
- Latent Dirichlet Allocation, 26
- Latent Semantic Indexing, 26
- United Kingdom, 11–13, 46
- University of Canberra, 13
- University of Oxford, 11
- User Survey, 13
- Using Ontological and Document Similarity to Estimate Museum Exhibit Relatedness, 28
- Webhose, 49
- Wikimedia, 50
- Wikipedia, 27, 28, 50
 - wikipedia, 39
- Wireframe, 37
- Yahoo News Digest, 23
- YouGov, 13