

University of Sheffield

# William Briggs 6 Month Report



William Briggs

*Supervisor:* Dr Mark Stevenson

*Panel:* Dr Andreas Vlachos, Professor Richard Clayton

Department of Computer Science

May 18, 2018

## Declaration

All sentences or passages quoted in this report from other people's work have been specifically acknowledged by clear cross-referencing to author, work and page(s). Any illustrations that are not the work of the author of this report have been used with the explicit permission of the originator and are specifically acknowledged. I understand that failure to do this amounts to plagiarism and will be considered grounds for failure in this project and the degree examination as a whole.

Name:

---

Signature:

---

Date:

---

## **Abstract**

Medical data comes in large volumes, making it a challenge for systematic reviewers to process and find relevant information. Being able to apply automatic techniques to this field presents itself as a suitable application for natural language processing. We will look at two areas of this field. Biomedical information retrieval by efficiently indexing and querying a large medical database, can we return an optimum set of documents. Stopping criteria, given an existing set of rankings, what is the most efficient way to look at these documents and determine a cut-off (stopping) point.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Literature Survey</b>	<b>2</b>
2.1	Steps of a Systematic Review . . . . .	2
2.1.1	Question Definition . . . . .	2
2.1.2	Relevant literature search . . . . .	3
2.1.3	Data Filtering . . . . .	3
2.1.4	Data Extraction . . . . .	4
2.2	Indexing and Querying Medline and Automated Query Generation . . . . .	4
2.2.1	Automated Query Generation . . . . .	4
2.3	Stopping Criteria . . . . .	5
2.3.1	Evaluation Metrics for Finding Stopping Points . . . . .	5
2.3.2	Existing Stopping Methods . . . . .	6
2.4	Summary . . . . .	9
<b>3</b>	<b>Research Questions</b>	<b>10</b>
3.1	Automated Decision Making of Relevant Studies . . . . .	10
3.2	Information Extraction From Studies . . . . .	11
3.2.1	PICO Extraction . . . . .	11
3.3	Stopping Methods . . . . .	11
<b>4</b>	<b>Novel Work</b>	<b>12</b>
4.0.1	CLEF 2017 Runs . . . . .	12
4.1	Baseline Approaches to Stopping . . . . .	13
4.1.1	Percentage cut-off . . . . .	13
4.1.2	Similarity score cut-off . . . . .	14
4.2	Sample Method to Stopping . . . . .	16
4.2.1	Curve Fitting . . . . .	16
4.2.2	Gaussian Process Fitting . . . . .	18
4.2.3	Conclusion on Curve fitting and GP . . . . .	19
4.2.4	Gaussian Process Conclusion . . . . .	19
4.3	Indexing and Querying Medline with Limited Information . . . . .	19

<i>CONTENTS</i>	4
4.3.1 Acquiring Key Information from A Systematic Review Protocol . . . .	19
4.3.2 Indexing medline . . . . .	20
4.3.3 Results . . . . .	21
4.3.4 Medline automatic query Conclusion . . . . .	21
<b>5 Future Work</b>	<b>23</b>
5.1 Gnatt Chart . . . . .	24
<b>6 Doctorial Development Programme</b>	<b>25</b>

# List of Figures

2.1	Visualisation of target method last relevant document selection. $C$ is number of documents in collection. . . . .	7
2.2	Example of using knee method to find a stopping point. Image inspired from [20] . . . . .	9
4.1	Visualisation of using a confidence interval for predicting a stopping point. . .	17
4.2	Visualisation of using a confidence interval for predicting a stopping point using a gp. We can see the gp is significantly over fitting to the data. . . . .	18
5.1	Gnatt Chart . . . . .	24

# List of Tables

4.1	The 6 runs sampled from the CLEF 2017 task. . . . .	12
4.2	Lowest effort possible to find 70% of relevant documents. . . . .	13
4.3	Comparison of results between rankings when looking at a percentage of the ranked documents. . . . .	14
4.4	Similarity cut-off comparison for stopping for Test_Data_Sheffield-run-2. Using cosine similarity scores. . . . .	15
4.5	Similarity cut-off comparison for stopping for Test_Data_Sheffield-run-2. Using cosine similarity scores. Compares first document to succeeding documents .	15
4.6	Evaluation of curve fitting for different CLEF 2017 runs. lower = lower-bound confidence interval. Sample size = 3. Results are taken as averages over all topics for search method. with 0.5% cut-off . . . . .	17
4.7	Comparison of different of sample method using gp for different CLEF 2017 runs. lower = lower-bound confidence interval. Sample size = 3. Results are taken as averages over all topics for search method. with 0.5% cut-off . . . .	18
4.8	Comparing target method with curve fitting along with confidence interval. Using Sheffield-run-2 . . . . .	19
4.9	Results for IR medline system. Comparison for both 5000 and 25000 thresholds	21

# Chapter 1

## Introduction

Medical literature poses interesting challenges for Natural Language Processing (NLP) researchers. The sheer volume of medical data makes it difficult for humans to process efficiently.

Evidence-based medicine has become an important aspect of health care and policy making. One key task is the creation of systematic reviews. Systematic reviews are transparent reviews that aim to pull together and critically analyse and summarise relevant literature to a topical question [6]. The process of creating a systematic review is rigorous and time consuming with varying degrees of complexity in-between steps [17]. This report will look at the existing work done using NLP as part of the systematic review process as well as the novel work done by myself so far.

We first review the stages involved in creating a systematic review 2.1. We break the steps down by looking at the PICO strategy [16]. By breaking the steps down, it becomes easier to examine potential candidates for applying NLP techniques to the process. Areas for research are then identified.

We then move on to look at stopping methods for systematic reviews. 2.3 Stopping methods are about finding a suitable stopping point given a list of ranked documents. Two existing methods are examined; the target method and the knee method.

In the next section we identify some relevant research questions within the field 3. Three areas are focused on: Automated Decision Making of Relevant Studies 3.1, Information Extraction From Studies 3.2 and Stopping Methods

The work completed so far is then presented 4. We look at using curves to make predictions of finding a stopping point, including using a Gaussian process. We then present our work for the automatic query generation process by using systematic review protocols as a basis for inferring the query.

Finally we look at future work 5.



# Chapter 2

## Literature Survey

Systematic reviews have many different stages that propose themselves as a candidate for automation. This section is going to look at the techniques that have been applied for some of these stages in previous literature.

### 2.1 Steps of a Systematic Review

It is useful for us to break down the steps involved in creating a systematic review into subtasks. This way we can observe what techniques can be applied during the relevant subtasks to improve the efficiency of the process. The following definitions are derived task simplifications from the cochrane tutorial on systematic reviews: [15].

1. Question definition.
2. Relevant literature search.
3. Data Filtering.
4. Data Extraction.
5. Analysis and Data combination.

#### 2.1.1 Question Definition

One of the best known techniques for formulating a systematic review question is known as the PICO strategy [16]. This technique focuses on exposing 4 pieces of information in the systematic review question: patient population, intervention or exposure, comparison or control and outcome.

Example: (credit goes to [16])

”Is animal-assisted therapy more effective than music therapy in managing aggressive behaviour in elderly people with dementia?”

P	elderly patients with dementia
I	animal-assisted therapy
C	music therapy
O	aggressive behaviour

A potential point of interest would be attempting to generate these questions automatically given some literature context.

### 2.1.2 Relevant literature search

After formulating a question, systematic reviews need to search for the relevant literature that surrounds this question.

Large medical database-such as Pubmed [18] contain relevant studies that can be used to create the review. These databases are typically very large and require concise queries to efficiently retrieve data.

Naturally this can be modelled as an information retrieval problem. We have a large number of documents and we wish to retrieve the most relevant ones. One task for the 2017 CLEF conference was to produce a ranking of the most relevant documents for topics [11]. Many techniques have been proposed for ranking of relevant documents, with varying degrees of performance [1] [?] [13].

An important aspect of the relevant literature search step is the construction of the query. Query creators often apply filters (also known as hedges) to increase the effectiveness or/and the efficiency of the searching. Two key attributes for the query are the precision and the recall. By including synonymous phrases e.g: quality adjusted life or quality of well-being or disability adjusted life the sensitivity can be increased, but as expense of the precision. The creation of this query is a task that could potentially have some aspects of NLP applied to it.

### 2.1.3 Data Filtering

The data filter stage involves reducing the amount of documents returned by the initial query down to a smaller subset of relevant document. This is can also be referred to as the abstract screening phrase [11].

The length of this stage is highly dependant on how many documents were returned by the initial query, often in the excess of 5000 studies for a single query. In response to this, stopping criteria methods have been proposed that aim to optimize two key parameters; the effort and the recall. That is to say we want to get as many relevant documents as possible,

whilst looking at the fewest. Examples of approaches include the knee method [20] and the target method [4]. Other techniques could be applied and evaluated such as curve fitting, which is later examined in 2.3

#### **2.1.4 Data Extraction**

The data extraction phase involves pulling the relevant information from the filtered subset of studies. Examples of important information includes how many people took part in the study and what the results were.

Being able to extract the relevant information from studies presents itself as an information extraction problem. The task to automate the process of extracting relevant information would reduce time and complexities of manually reviewing studies [9].

## **2.2 Indexing and Querying Medline and Automated Query Generation**

Medline is a large collection of medical literature and data from around the world. Typically each entity will contain a title and an abstract containing some information on the study. Whilst Medline as a whole is very easy to access [14], the large size and complexity of the data makes it difficult to retrieve the relevant information.

### **2.2.1 Automated Query Generation**

Being able to automate query generation for literature searching would save systematic reviewers a significant amount of time. However, medical literature queries are typically complex and contain multiple levels of logical operators and synonymous term look ups. This makes the task of creating a query manually is a challenging task.

#### **Rapid Automatic Keyword Extraction Algorithm**

Rapid automatic keyword extraction (RAKE) is a keyword extraction algorithm [19]. This algorithm is used for taking the key pieces of information from text and is useful the domain of information extraction. This algorithm is of interest as it has potential usage within the field of query generation.

RAKE heavily relies on stop-words and punctuation separators as an indicator for the importance of a phrases and words. RAKE will iterate over sequences of words until a stop-word or separator is found, this phrase/word is then split and extracted. Frequency of

occurrence (tf) and word co-occurrence matrices can then be used to reduce the key-word set down further.

RAKE can be further optimized by specifying minimum term frequency rates to capture more prominent terms.

## 2.3 Stopping Criteria

Stopping criteria is about finding the optimum point to stop reviewing a set of documents. This could be useful in a decision making process. Consider having 100 relevant documents, where each document contains a binary value. If we looked at 1/3 of these documents and saw a trend of positive values, we could use this to infer the reliability of the remaining documents.

Another use of stopping criteria is when filtering through potentially relevant documents. Consider a query that returns 10000 documents, of which only a small sub-set of these are relevant. Reviewers would need to filter through each of these 10000 documents to pull out the relevant ones. Or it could be that the reviewers are happy to hit a 90% recall of relevant documents, and are happy to miss the remaining 10% in exchange for time-saved.

Two key methods have been proposed for finding stopping points so far, the target method [20] and the knee method. [4]. Both these methods are discussed below 2.3.2

### 2.3.1 Evaluation Metrics for Finding Stopping Points

In order to evaluate the suitability of our stopping method, we can use two evaluation metrics. The recall, which is simply the number of documents returned for a topic, and effort which is the number of documents we had to be examined

$$Recall = \frac{|R|}{|D|} \quad (2.1)$$

Where  $R$  is the set of relevant documents found and  $D$  is the set of all relevant documents.

Similarly, effort is computed as:

$$Effort = \frac{|L|}{|D|} \quad (2.2)$$

Where  $L$  is the set of documents that were examined.

Naturally we could exclusively optimized each of the parameters by either returning everything in the document collection ( $R = |D|$ ) or by just looking at a single document. ( $L = 1$ )

Therefore it becomes difficult for us to evaluate our stopping criteria as we need to consider both of these parameters adjacently.

In response to this we can make use of two more evaluation metrics that were proposed by Cormack and Grossman [4]:

$$reliability = P[acceptable(S) == 1] \quad (2.3)$$

reliability is computed over all searches and is read as the probability of the acceptability being 1. Where acceptability is calculated as:

$$acceptability(S) = \begin{cases} 1, & recall(S) \geq 0.7. \\ 0, & recall(S) < 0.7. \end{cases} \quad (2.4)$$

A stopping point is deemed to be acceptable if 70% of the relevant documents have been found. As such, the reliability is an average over a search method. In Cormack and Grossman [4] a target is set of achieving 95% reliability.

### 2.3.2 Existing Stopping Methods

#### Target

The target method is an approach that can guarantee a certain level of reliability 2.3.1. It was first proposed by Cormack and Grossman [4]. This method randomly samples returned documents until a target number  $T$  of relevant document are found.

The target  $T$  denotes how many documents we should randomly select from our initial query. A larger value of  $T$  will increase the effort required as we are more likely to select a document towards the end of the query set. Documents are looked at until the target point  $T$  has been reached.

We first compute a random target set of relevant documents. We then calculate the last document in the target set and mark that as our target point:

$$d_{last} = \underset{d \in T}{argmax} relrank(d) \quad (2.5)$$

It must hold that  $d$  is in the target set. *relrank* determines whether or not a document is relevant.

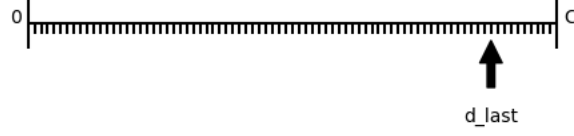


Figure 2.1: Visualisation of target method last relevant document selection.  $C$  is number of documents in collection.

Increasing our target set size is likely to increase the probability the last document being towards the end of the document collection.

We can calculate the recall of the point by looking at the relevance rank of the last document:

$$recall = \frac{relrank(d_{last})}{R} \quad (2.6)$$

Where  $R$  is the number of relevant documents.

For our method to be deemed reliable we must achieve 70% recall with a 95% average over all topics.

$$P\left(\frac{relrank(d_{last})}{R} \geq 0.7\right) \geq 0.95 \quad (2.7)$$

Assuming we have a large number of relevant documents  $R$  we need to determine cut-off  $c$

$$P\left(\frac{R - relrank(d_{last})}{R} > c\right) = 0.05 \quad (2.8)$$

This can be further simplified to:

$$P(R - relrank(d_{last}) > cR) = 0.05 \quad (2.9)$$

Which translates to the probability of the remaining relevant documents being higher than the cut-off point should be 0.05.

For this to hold,  $cR$  documents must be absent from  $T$ . This occurs with the probability:

$$\left(1 - \frac{10}{R}\right)^{cR} = 0.05 \quad (2.10)$$

Which can become:

$$c = \frac{\log(0.05)}{R \log(1 - \frac{10}{R})} \quad (2.11)$$

In cases where  $R$  has more than 10 relevant documents it follows:

$$c < \lim_{R \rightarrow \infty} \frac{\log(0.05)}{R \log(1 - \frac{10}{R})} = 0.299573 < 0.3 \quad (2.12)$$

Finally we have:

$$R \leq 10 \cup P\left(\frac{\text{relrank}_{last}(d)}{R} \geq 0.7\right) \geq 0.95 \quad (2.13)$$

Overall, while the target method is shown to acquire 95% reliability, the effort needed is often significantly highly, often requiring us to look at huge volume of documents.

## Knee Method

A different stopping kethod proposed by [20] is known as the knee method. This approach uses a curve to generate a 'knee', which is then used for predicting a stopping point. This approach is likely to be highly dependant on the quality of the initial rankings. This is because we need a curve that reaches a peak quickly, before flattening out.

We use a vertical line panning the length of the ranking set and use it to calculate the distance from the ranking at each point. The point with the maximum distance is chosen as a suitable stopping point.

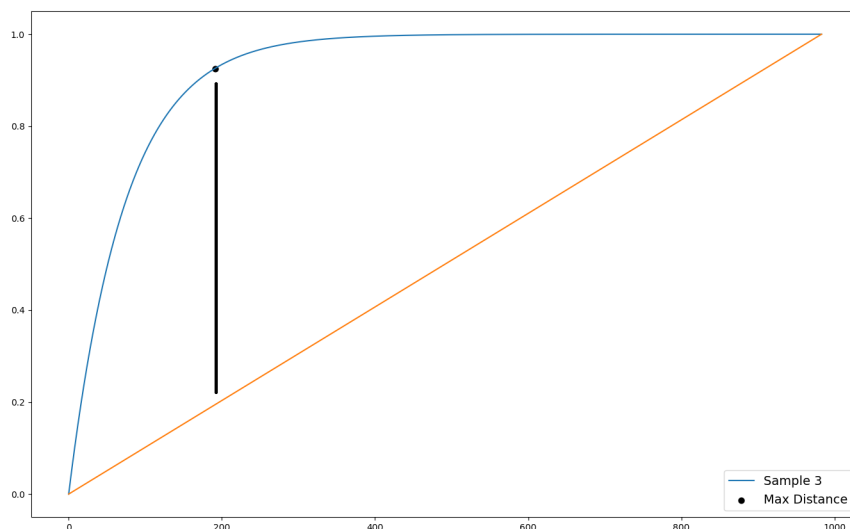


Figure 2.2: Example of using knee method to find a stopping point. Image inspired from [20]

We can see in the above example the method has predicted we look at around 200 of the 1000 documents to achieve a suitable stopping point.

This method also imposes an additional constraint for rankings of a large volume.

It was found that the knee method is a better approach for finding a stopping point than the Target method [4]. The recall was always found to be better and the reliability was found to be the same or higher for 6 out of 8 test collections.

## 2.4 Summary

We first examined the steps of a systematic review 2.1 and looked at potential areas to introduce NLP techniques.

Next we looked at how we can handle large scale databases like Medline 2.2. We also introduced the key-word extraction algorithm RAKE as a candidate for extracting relevant information from medical documents 2.2.1.

We then began to examine stopping criteria 2.3. We looked at the intuition behind looking for stopping points in ranked sets of documents. We looked into two existing approaches, the target method and the knee method



## Chapter 3

# Research Questions

In this chapter we will present some possible research areas for systematic reviews.

### 3.1 Automated Decision Making of Relevant Studies

Once relevant studies have been identified by their abstracts, reviewers are required to process the studies and extract useful information for the systematic review. Many of these studies will not be useful and can be discarded, but not without the cost of the reviewer having to look at the content.

A useful piece of research would be to determine if a study is relevant to the systematic review question.

There are two potential routes that could be taken for the investigating this task. The first approach would be to use existing information from the systematic review (ie the protocol) and determine if the study contains the information. A second approach would be to a semi-supervised learning method, having the reviewer look at a subset of the studies, so we can then build a classifier for the remaining studies. Both of these approaches can be applied to the abstract screening and the data filtering stage of the systematic review. For the abstract screening we would be trying to predict if a study is relevant, by looking at the abstract. For the data filtering stage we would be looking at the actual text of the study (typically a pdf file)

There are many challenges for this research topic. Not all studies are publicly available and are often protected by publishing licences. This makes it difficult to gather data. Another challenge is having to deal with such a broad range of data, as well as the pdf format studies.

## 3.2 Information Extraction From Studies

Much of the existing work in information extraction for systematic reviews is done using sentence-level information extraction [17] [10]. Often approached using classification methods to determine if a sentence does or does not contain relevant information [8].

### 3.2.1 PICO Extraction

A method for extracting PICO sentences was proposed [8]. This approach uses a Naive Bayes classifier against 3 of the 4 PICO elements. Training data was acquired using a bootstrapping technique, looking at structured abstracts that are contain key-word headings. The macro-averaged F-Measure was found to be 84%.

Further work [7] by the same author looked at whether the first sentences of in the structured abstracts were enough to train a classifier. This approach looks at 3 of the 4 PICO elements. The averaged macro average F-Measure was found to be 71%.

This existing work assumes extraction of PICO elements at abstract level. No work was found [10] in the extraction of PICO elements from the full texts. One interesting area of research, would be looking at learning a classifier from the abstracts and applying it to the full texts.

## 3.3 Stopping Methods

We examined existing working on finding on stopping methods in our literature survey 2.3. We can investigate further methods finding stopping points in systematic review rankings.

We can look at applying machine learning algorithms for plotting a regression based system of relevant documents. Techniques that could be interesting to experiment with include a generic curve fitting and a Gaussian Process (GP).

# Chapter 4

## Novel Work

In this section the work completed so far will be presented. Two main areas of the systematic review process has been focused on: stopping criteria and indexing/querying PubMed.

### 4.0.1 CLEF 2017 Runs

For the CLEF 2017 Technologically Assisted Reviews in Empirical Medicine [11], participants were expected to submit runs for ranking documents. Participants were given complex boolean queries that could be used for extracting relevant information to rank the documents. These runs were later released for public access <sup>1</sup>

Run	Description	Reference
Sheffield-run-2	Sheffield-2 used tfidf similarity along with standard pre-processing	[1]
Waterloo A-rank-cost	Waterloo used a baseline model implementation from the TREC Total Recall Track	[5]
Waterloo B-rank-cost	-	[5]
auth run-1	AUTH used a learning-to-rank approach and used both batch and active learning	[2]
auth run 2	-	[2]
ntu run-1	Used convolutional neural networks (CNN)	[12]
ucl full-text	Used a deep learning model architecture	[22]

Table 4.1: The 6 runs sampled from the CLEF 2017 task.

We have taken 6 ranking sets of which are of different quality. The Waterloo and AUTH ranks are the best rankings followed by Sheffield. The UCL and NTU submissions feature poorer quality rankings.

CLEF 2017 runs will follow a format similar to the example below. The most important information being the topic id and the document id.

CD010775 NF 19307324 1 0.27152011529138564 Test-Data-Sheffield-run-2

<sup>1</sup><https://github.com/CLEF-TAR/tar/tree/master/2017-TAR/participant-runs>

Results can be evaluate using qrel files, supplied as part of the CLEF 2017 data. These files contain relevant documents for each topic, and are formatted as follows:

```
CD008803 0 21467181 0
CD008803 0 20872357 1
CD008803 0 23837966 0
```

## 4.1 Baseline Approaches to Stopping

As discussed in 2.3.2 there are many baseline approaches we can take for finding a suitable stopping point. Most approaches are heavily dependant on the initial rankings of the document collection, and assume more relevant documents feature towards the start of the collection.

### 4.1.1 Oracle Scores

Looking at the oracle scores for each set of rankings will tell us the best we can possibly for do this task. We will assume we are satisfied with 70% recall. Results are taken as averages over all topics.

Submission	recall	reliability	effort
Test_Data_Sheffield-run-2	0.7	1.0	0.11
Waterloo A-rank-cost	0.7	1.0	0.07
Waterloo B-rank-cost	0.7	1.0	0.06
auth run-1	0.7	1.0	0.08
auth run-2	0.7	1.0	0.09
ntu run-1	0.7	1.0	0.4
ucl full-text	0.7	1.0	0.67

Table 4.2: Lowest effort possible to find 70% of relevant documents.

These scores highlight the importance of the ranking methods for finding a stopping point. The best performer, Waterloo B-rank-cost needs just 6% effort to hit 100% reliability. The worst performer, ucl full-text requires 67% effort for the same level of reliability.

### 4.1.2 Percentage cut-off method

One approach is we can simply take a cut of the document collection and evaluate how many relevant documents we have retrieved. All scores are averaged over the entire ranking set.

% of Documents	10%			25%			50%			75%			90%		
Run	Recall	Reliability	Effort	-	-	-	-	-	-	-	-	-	-	-	-
Sheffield-run-2	0.49	0.16	0.10	0.74	0.66	0.25	0.91	0.93	0.50	0.98	1.00	0.75	0.99	1.00	0.90
Waterloo A-rank-cost	0.80	0.6	0.10	0.91	0.93	0.25	0.98	1.00	0.50	0.99	1.00	0.75	0.99	1.00	0.90
Waterloo B-rank-cost	0.73	0.63	0.10	0.90	0.93	0.25	0.98	1.00	0.50	0.99	1.00	0.75	0.99	1.00	0.90
auth run-1	0.72	0.63	0.10	0.90	0.93	0.25	0.97	1.00	0.50	0.99	1.00	0.75	0.99	1.00	0.90
auth run 2	0.68	0.60	0.10	0.88	0.9	0.25	0.97	1.00	0.50	0.99	1.00	0.75	0.99	1.00	0.90
ntu run-1	0.19	0.00	0.10	0.39	0.06	0.25	0.62	0.43	0.50	0.83	0.83	0.75	0.91	0.93	0.90
ucl full-text	0.08	0.00	0.10	0.22	0.00	0.25	0.46	0.03	0.50	0.70	0.70	0.75	0.87	0.93	0.90

Table 4.3: Comparison of results between rankings when looking at a percentage of the ranked documents.

Sheffield has an average recall of 0.49 by the time 10% of the documents have been observed. Waterloo has achieved 80% recall by this point. After looking through 25% of the rankings Waterloo and AUTH have achieved around 90% recall of documents.

The limitation of this approach is the effort required looking through documents is still high. We want to lower this effort as much as possible, whilst only having to observe relevant documents.

#### 4.1.3 Similarity score cut-off method

The Sheffield document rankings will feature a similarity score between the document and the query. As we descend down the rankings, the score decreases in value as documents become less relevant.

We can make use of the similarity score to derive a stopping point. This method works by looking at documents  $D_i$  and  $D_{i+1}$  and determining if the difference between the similarity scores has become too large. This method will work on the basis that documents that are no longer relevant will have a sudden drop in score such that we can identify this as our stopping point.

$$Difference(D_i, D_{i+1}) > C \quad (4.1)$$

Where difference returns a score of how close document  $D_i$  and  $D_{i+1}$  are together and  $C$  is a cut-off constant. We can expand this to an example:

$$(1 - (0.73/0.75)) * 100 > 0.015 \quad (4.2)$$

Here we are saying if the two documents' scores are above 1.5% then we should stop looking down the rankings.

$dif(D_i, D_{i+1})$	recall	reliability	effort
0.01%	0.025	0.000	0.0023
0.05%	0.120	0.100	0.100
1%	0.359	0.333	0.333
2%	0.880	0.860	0.860
5%	1.000	1.000	1.0000

Table 4.4: Similarity cut-off comparison for stopping for Test\_Data\_Sheffield-run-2. Using cosine similarity scores.

These results are highly sensitive to the similarity score and show it's difficult to use this score as an affective measure for stopping. We found similarity scores rarely have sudden drops in values, making it difficult to use this method to identify a stopping point.

As an alternative approach, we can look at just the top document  $D_1$ , and compare it to the succeeding documents  $D_{1+i}$  in the rankings. We can formulate this as follows:

$$Difference(D_1, D_{1+i}) > C \quad (4.3)$$

$dif(D_1, D_{1+i})$	recall	reliability	effort
10%	0.048	0.000	0.004
20%	0.063	0.000	0.007
30%	0.113	0.000	0.152
40%	0.191	0.030	0.029
50%	0.319	0.060	0.063
60%	0.460	0.200	0.113
70%	0.638	0.433	0.210
80%	0.841	0.800	0.387
90%	0.979	1.000	0.679
100%	1.000	1.000	1.000
<b>85.5%</b>	<b>0.934</b>	<b>1.000</b>	<b>0.538</b>

Table 4.5: Similarity cut-off comparison for stopping for Test\_Data\_Sheffield-run-2. Using cosine similarity scores. Compares first document to succeeding documents

We found results to be much better when using the first document to look for a cut-off point. 85.5% difference in first and subsequent documents was found to be the point for hitting 1. reliability. This comes at the expense of making just over 50% effort.

## 4.2 Sample Method to Stopping

The limitations of the methods presented in 4.1 is that they they still require looking through a large volume of documents.

The approaches in this section will use sampling methods. These approaches assumes we have sensible ranking algorithm for returning documents for a query. They work by generating a sample set, which acts as our model for predicting a stopping point.

The first step is to generate a sample set. We used an interval method for generating our set, i.e select every  $N$ th document. The intuition being that the distribution of relevant document in the sample set, should be similar to that of the real set of relevant documents. This makes it suitable to use as a model.

### 4.2.1 Curve Fitting

Our first approach is to fit a non-linear curve against a sample set. We used a sample size of 3.

$$F(x) = n - a \exp^{-kx} \quad (4.4)$$

Where  $a$ ,  $k$  and  $n$  are learnt weights and  $x$  is an associated return rate for a document. We generate the curve using the non-linear least squares algorithm [21].

### Curve Predictions

The number of topics was reduced down from 30 to 23 using a cut-off parameter. This reduces some of topics that contain fewer relevant documents and do not generate suitable prediction curves. We used a cut-off parameter of 0.5%. This can be read as the number of relevant documents in a document collection must be atleast 0.5%

We also include a confidence interval evaluation for lower bounded range of a  $3\sigma$  confidence interval. The key advantages of using a curve as method of evaluating stopping criteria is being able to make use of this confidence interval in a real systematic review. In the context of a systematic reviewers at the data filtering stage, we could specify that the system is 95% certain that 70% of relevant documents have been found. At which point the reviewer can decide if its worth continuing to look at documents.

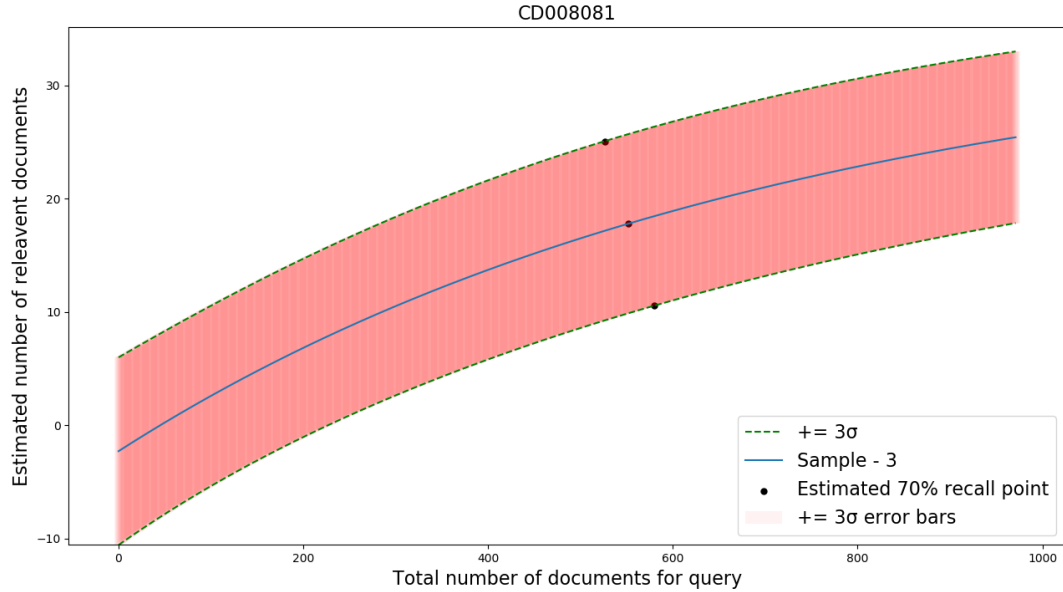


Figure 4.1: Visualisation of using a confidence interval for predicting a stopping point.

Submission	recall-lower	reliability-lower	effort-lower	topics sampled
Test_Data_Sheffield-run-2	0.69 0.74,	0.52 0.60	0.48 0.51	23
Waterloo A-rank-cost	0.71 0.73,	0.47 0.47	0.43 0.44	23
Waterloo B-rank-cost	0.71 0.75,	0.52 0.82	0.41 0.43	23
auth run-1	0.72 0.74,	0.52 0.60	0.41 0.42	23
auth run-2	0.70 0.72,	0.52 0.60	0.42 0.43	23
ntu run-1	0.76 0.74,	0.56 0.52	0.72 0.70	23
ucl full-text	0.86 0.94,	0.82 0.86	0.91 0.95	23

Table 4.6: Evaluation of curve fitting for different CLEF 2017 runs. lower = lower-bound confidence interval. Sample size = 3. Results are taken as averages over all topics for search method. with 0.5% cut-off

We have deliberately compared two of the better participant rankings (Waterloo and auth) and two of the lower performers (ntu and ucl). We can see the quality of the initial rankings significantly influences the performance of our stopping criteria. This suggests there is an important relationship between using a curve to predict a stopping point and how good the initial ranking of documents is.

Some of the ranking methods struggle to produce curves and when combined with a cut-off parameter produce become not worth considering in our evaluation. In this situation we simply returned everything for the given topic, resulting in 100% recall at the expense of



100% effort.

## 4.2.2 Gaussian Process Fitting

As an alternate approach to fitting a simple curve, we can apply a Gaussian Process (GP).

We will apply a constant kernel plus a squared-exponential kernel. One of the key challenges when applying a Gaussian Process was the amount of over-fitting.

Submission	recall-lower	reliability-lower	effort-lower	topics sampled
Test_Data_Sheffield-run-2	0.73 0.73,	0.73 0.73	0.50 0.50	23
Waterloo A-rank-cost	0.70 0.70,	0.40 0.40	0.42 0.42	23
Waterloo B-rank-cost	0.73 0.73,	0.71 0.71	0.41 0.41	23
auth run-1	0.74 0.75,	0.56 0.60	0.42 0.39	23
auth run-2	0.72 0.73,	0.52 0.52	0.42 0.42	23
ntu run-1	0.67 0.67,	0.40 0.40	0.64 0.64	23
ucl full-text	0.62 0.62,	0.52 0.56	0.82 0.82	23

Table 4.7: Comparison of different of sample method using gp for different CLEF 2017 runs. lower = lower-bound confidence interval. Sample size = 3. Results are taken as averages over all topics for search method. with 0.5% cut-off

We can see the gp does a better job at fitting to the topics than the curve method. However, for most of the lower ranking runs the effort is increases significantly, suggesting the gp is not fitting very well to the data.

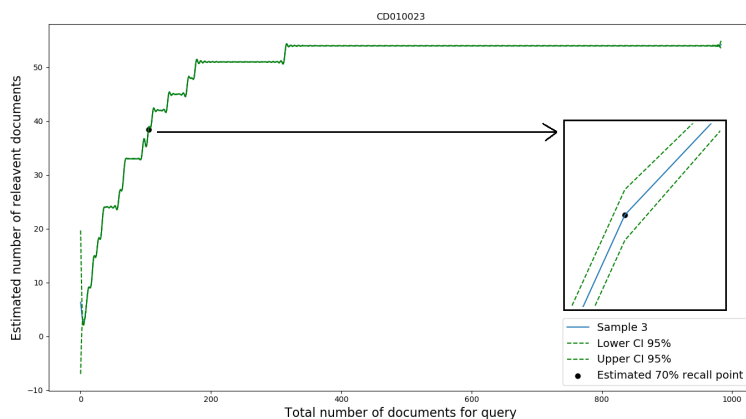


Figure 4.2: Visualisation of using a confidence interval for predicting a stopping point using a gp. We can see the gp is significantly over fitting to the data.

### 4.2.3 Conclusion on Curve fitting and GP

We implemented two methods for predicting stopping points in ranked medical studies. Our first approach used a general curve to estimate the point in which 70% recall is likely to have been hit. Our second method used a Gaussian Process in the same way. We used a sampling method to generate our curves to make predictions about the remaining studies.

#### Comparing with Target Method

Method	Target	recall	reliability	effort
Target Method	10	0.952	0.96	0.652
Sheffield-run2-curve	-	0.75	0.64	0.51

Table 4.8: Comparing target method with curve fitting along with confidence interval. Using Sheffield-run-2

We can evaluate how well this method does against an existing set of relevance rankings. We will use the Sheffield run data from the CLEF 2017 task [11]. As the target method allows us to specify our level of reliability, we needed a target  $T$  of 9 to hit 95% reliability. We can see that the reliability of our method does not come close to the Target method. Our curve, however has a much lower effort.

### 4.2.4 Gaussian Process Conclusion

So far, we found using a GP was not an appropriate choice due to the level of over-fitting. Further work is needed in determining an efficient kernel for our data.

## 4.3 Indexing and Querying Medline with Limited Information

CLEF 2018 [3] presents an appropriate sub-task for using a limited amount of information to retrieve relevant documents. Normally, reviewers are required to construct complex Boolean queries to retrieve data from Medline. The objective of CLEF 2018 Sub-Task 1: No Boolean Search [3] is to search effectively and efficiently bypassing the construction of the Boolean query.

### 4.3.1 Acquiring Key Information from A Systematic Review Protocol

A systematic review protocol is created before the systematic review process is started. A systematic review protocol describes the rationale, hypothesis, and planned methods of the

review. The Medline query is created manually with the help of the protocol. Here we are looking to generate a suitable query/relevant information from the protocol to then automatically query Medline.

We used RAKE [19] to extraction key-words from a protocol. The minimum word occurrence count is set to 1, as the protocols are typically small. We used a pubmed stop list as the phrase splitting parameter. Example shown below:

Topic: CD008122  
 Title: Rapid diagnostic tests for diagnosing uncomplicated *P. falciparum* malaria in endemic countries  
 Objective: To assess the diagnostic accuracy of RDTs for detecting clinical *P. falciparum* malaria (symptoms suggestive of malaria plus *P. falciparum* parasitaemia detectable by microscopy) in persons living in malaria endemic areas who present to ambulatory healthcare facilities with symptoms of malaria, and to identify which types and brands of commercial test best detect clinical *P. falciparum* malaria.

endemic countries objective|ambulatory healthcare facilities|rapid diagnostic tests|falciparum parasitaemia detectable|malaria endemic areas|diagnostic accuracy|falciparum malaria

The | symbol represents a separation between a phrase. The protocols are pre-processed as follows: Reference removal, lowercase, words less than  $N$  length removed, pubmed stoplist. We decided to not perform any stemming/additional manipulation at this stage, due to uncertainty of query format. The resulting content is stored in a separate file appended with '.kwq' (key-word query).

The key-word-query receives some final pre-processing prior to being loaded into our information retrieval (IR) system. We used a Lancaster stemmer to reduce words down to a base form. The result is as follows:

endem country object amb healthc facil rapid diagnost test falcipar parasitaem detect malar endem area diagnost acc falcipar malar

### 4.3.2 Indexing medline

Medline was downloaded from the online resource. We processed the xml files and retrieved the information for each study - title, id, abstract. To reduce the size, we store each record into a local database, containing only the relevant information for each study.

We used Apache Lucene to generate an index for the medline local database. The abstract

and title were concatenated together. Pre-processing was done using the same format as the query: pub-med stoplist, Lancaster stemmer and lower casing.

We will provide comparisons of various ranking methods as well as the evaluation scores for each.

### 4.3.3 Results

Results were generated using the eval script from the CLEF 2017/2018 task [11]. We calculated the top  $N$  results over the CLEF 2017 training set. We include a random baseline to provide a comparison between results.

Run	recall	ap	lastrel	wss100	wss95	normarea	$N$
Random-baseline	0.05	0.02	126.7	0.0	-0.0	0.024	-
Train-Data-Sheffield-bm25-Run1-objective-only	0.538	0.034	3039.051	0.101	0.108	0.431	5000
Train-Data-Sheffield-tfidf-Run1-objective-only	0.354	0.007	2633.718	0.021	0.023	0.247	5000
Train-Data-Sheffield-boolean-Run1-objective-only	0.313	0.034	3039.051	0.101	0.108	0.431	5000
Train-Data-Sheffield-bm25-Run1-objective-only	0.68	0.034	12310.231	0.169	0.172	0.592	25000
Train-Data-Sheffield-tfidf-Run1-objective-only	0.601	0.007	14883.744	0.13	0.136	0.455	25000
Train-Data-Sheffield-boolean-Run1-objective-only	0.471	0.007	12974.205	0.03	0.029	0.381	25000

Table 4.9: Results for IR medline system. Comparison for both 5000 and 25000 thresholds

As we increase the number of documents we return, the recall naturally increases. When we return 25000 documents for each topic, we are able to obtain a total recall rate of over 58%. However, the precision (ap, average precision) is very low, suggesting a significant amount of the documents are not useful. BM25 was found to be the best ranking method, followed by tfidf and boolean.

Improvements could certainly be made to this system:

- MeSH headings would be useful in expanding the range of the query to capture synonymous terms.
- Tokenization could be optimized to capture phrases of different sizes.
- Introducing a cost or stopping point to remove the amount of non-relevant documents. We can see for the 25000 documents set of results the last relevant document was around the 20000 point, meaning we could drop the last 5000 from our result set.

### 4.3.4 Medline automatic query Conclusion

We built an IR system using Apache Lucerne and compared three separate ranking methods. We found bm25 ranking gave the best results overall.

We found we were able to achieve fair results with a little optimization techniques to the index and query data.

We compared the performance of our system across different return thresholds, naturally finding as we increase the returned number of documents we get a higher recall. This comes at the expense of reduced precision.

We suggested further improvement to our system, such as including a phrase model for more robust features for both index and query.

## Chapter 5

# Future Work

This chapter will outline a plan for future work. This chapter will relate back to our research questions chapter. 3

We will look at feature extraction for systematic reviews. This will involve looking at content within a study and attempting to identify it automatically. The first step of this process is being able to obtain relevant studies for a systematic review. We will look at techniques for obtaining these in the form of pdf files. Work will also need undertaking in processing a diverse range of pdf documents. We will then need to look at bringing in annotators to mark when the relevant information occurs with a study, or look for common patterns that occur across all studies. This is likely to be a very challenging, but potentially very beneficial task.

We will investigate classification of studies and determine if we can make a binary decision on whether or not a study is relevant to a research question. To achieve this we will first need to process the content of the studies and decide on a sampling strategy. This might involve using a proportion of the relevant studies as a training examples, and then trying to classify the remainder of the studies. We will also look at using unsupervised methods, by trying to directly use the systematic review as an indicator of relevant studies. We will need to overcome the challenge of condensing large studies (e.g full pdf texts) into relevant chunks.

Finally we will look at alternate approaches to finding stopping points. We found that fitting a GP caused too much over fitting to our sampled rankings. We will look at alternate regression-based machine learning algorithms to finding a stopping point.

## 5.1 Gantt Chart

Task Names	Resource Names	Start	Duration	Finish	% Complete	Calendar days	Days Completed	Days Remaining
		01-Jan-2018	990	15-Oct-2021	11%	1384		
<b>Other Tasks</b>	<b>Other Tasks</b>	<b>01-Jun-2018</b>	<b>71</b>	<b>09-Sep-2018</b>	<b>0%</b>	<b>101</b>	<b>0</b>	<b>19</b>
Complete ethics course		01-Jun-2018	9	05-Jun-2018	0%	5	0	9
Attend 2018 CLEF		05-Sep-2018	10	09-Sep-2018	0%	5	0	10
<b>Stopping Methods</b>	<b>Stopping Methods</b>	<b>01-Jan-2018</b>	<b>43</b>	<b>28-Feb-2018</b>	<b>57%</b>	<b>59</b>	<b>98</b>	<b>74</b>
Investigate existing methods to stopping		01-Jan-2018	21	29-Jan-2018	100%	29	21	0
Implement Curve fitting for stopping		01-Jan-2018	43	28-Feb-2018	100%	59	43	0
Implement GP fitting for stopping		01-Jan-2018	43	28-Feb-2018	50%	59	21	22
Evaluation of stopping methods		01-Jan-2018	65	30-Mar-2018	20%	89	13	52
Investigate new methods for stopping		29-Apr-2018	45	01-Jul-2018	5%	64	2	43
<b>Pubmed Protocol Query</b>	<b>Pubmed Protocol</b>	<b>01-Apr-2018</b>	<b>9</b>	<b>12-Apr-2018</b>	<b>42%</b>	<b>12</b>	<b>85</b>	<b>101</b>
Download pubmed and pre process/store index		01-Apr-2018	3	04-Apr-2018	100%	4	3	0
Learn Apache Lucene		05-Apr-2018	4	10-Apr-2018	100%	6	4	0
Learn Rake		10-Apr-2018	3	12-Apr-2018	100%	3	3	0
Extract releavent information from protocol with Rake		11-Apr-2018	8	20-Apr-2018	100%	10	8	0
Create API to query with Rake key word content		11-Apr-2018	9	23-Apr-2018	100%	13	9	0
Run Training topics through IR system		23-Apr-2018	6	30-Apr-2018	100%	8	6	0
Run Test samples through IR system		23-Apr-2018	10	06-May-2018	0%	14	0	10
Evaluate system using trec-eval		23-Apr-2018	10	06-May-2018	80%	14	8	2
Submit for CLEF 2018 Task		06-May-2018	0	06-May-2018	0%	1	0	0
Write up for CLEF 2018 Task		06-May-2018	44	05-Jul-2018	0%	61	0	44
Investigate impact of phrase IR system		05-Jul-2018	45	05-Sep-2018	0%	63	0	45
<b>Text Classification</b>	<b>Text Classification</b>	<b>01-Apr-2018</b>	<b>88</b>	<b>01-Aug-2018</b>	<b>9%</b>	<b>123</b>	<b>22</b>	<b>220</b>
<b>Text Classification of Systematic Reviews</b>		<b>01-Apr-2018</b>	<b>45</b>	<b>01-Jun-2018</b>	<b>0%</b>	<b>62</b>	<b>0</b>	<b>45</b>
Investigate fesability of Task		01-Apr-2018	45	01-Jun-2018	40%	62	18	27
Retreive relevant documents for training		01-Apr-2018	45	01-Jun-2018	10%	62	4	41
Investigate different classificaton methods		10-Jun-2018	38	01-Aug-2018	0%	53	0	38
Classifical Approaches		10-Jun-2018	38	01-Aug-2018	0%	53	0	38
Deep Learning Approaches		20-Jun-2018	31	01-Aug-2018	0%	43	0	31
<b>Viva</b>	<b>Viva</b>	<b>05-Oct-2020</b>	<b>270</b>	<b>15-Oct-2021</b>	<b>0%</b>	<b>376</b>	<b>0</b>	<b>810</b>
<b>Viva Preperation</b>		<b>05-Oct-2020</b>	<b>270</b>	<b>15-Oct-2021</b>	<b>0%</b>	<b>376</b>	<b>0</b>	<b>270</b>
Writeup		05-Oct-2020	270	15-Oct-2021	0%	376	0	270
Review		05-Oct-2020	270	15-Oct-2021	0%	376	0	270

Figure 5.1: Gantt Chart

## Chapter 6

# Doctorial Development Programme

- Attended Healtex - UK HEALTHCARE TEXT ANALYTICS CONFERENCE
- Lab demonstrator for module COM4519 Cloud Computing
- Undertook module HAR6169 Study Design and Systematic Review Methods
- Marked assignments for COM3110 Text Processing
- Enrolled on FCE6100 Professional Behaviour and Ethical Conduct
- Completed TRAINING NEEDS ANALYSIS (TNA) form.
- Used Learning Management System (LMS) to attend 6 teacher training courses.
- Gave introduction talk to NLP group.
- Contributed to 2018 CLEF lab.
- Became member of Text Processing for Health Technology Assessment (TePHTA).



# Bibliography

- [1] ALHARBI, A., AND STEVENSON, M. Ranking abstracts to identify relevant evidence for systematic reviews: The university of sheffield’s approach to clef ehealth 2017 task 2. *CLEF 2017* (2017).
- [2] ANAGNOSTOU, A., LAGOPOULOS, A., TSOUMAKAS, G., AND VLAHAVAS, I. A cost-effective hybrid ltr approach for document ranking. *Working Notes of CLEF* (2017).
- [3] CLEF. Clef. <https://sites.google.com/view/clef-ehealth-2018/task-2-technologically-assisted-reviews-in-en>
- [4] CORMACK, G. V., AND GROSSMAN, M. R. Engineering quality and reliability in technology-assisted review.
- [5] CORMACK, G. V., AND GROSSMAN, M. R. Technology-assisted review in empirical medicine: Waterloo participation in clef ehealth 2017. *Working Notes of CLEF* (2017), 11–14.
- [6] GOUGH, D., OLIVER, S., AND THOMAS, J. *An Introduction to Systematic Reviews*. Sage, London, 2012.
- [7] HUANG, K.-C., CHIANG, I.-J., XIAO, F., LIAO, C.-C., LIU, C. C.-H., AND WONG, J.-M. Pico element detection in medical text without metadata: Are first sentences enough? *Journal of Biomedical Informatics* 46, 5 (2013), 940 – 946.
- [8] HUANG, K.-C., LIU, C., YANG, S.-S., XIAO, F., WONG, J.-M., LIAO, C.-C., AND CHIANG, I.-J. Classification of pico elements by text features systematically extracted from pubmed abstracts. 279–283.
- [9] JONNALAGADDA, S. R., GOYAL, P., AND HUFFMAN, M. D. Automating data extraction in systematic reviews: a systematic review.
- [10] JONNALAGADDA, S. R., GOYAL, P., AND HUFFMAN, M. D. Automating data extraction in systematic reviews: a systematic review. *Systematic Reviews* 4, 1 (Jun 2015), 78.
- [11] KANOULAS, E., LI1, D., AZZOPARDI, L., AND SPIJKER, R. Clef 2017 technologically assisted reviews in empirical medicine overview.

- [12] LEE, G. Medical document classification for systematic reviews using convolutional neural networks. *Working Notes of CLEF* (2017), 11–14.
- [13] LEE, G. E. A study of convolutional neural networks for clinical document classification in systematic reviews: Sysreview at clef ehealth 2017. *CLEF 2017* (2017).
- [14] MEDLINE. medline. <https://www.nlm.nih.gov/bsd/pmresources.html>.
- [15] NUNN, J. cochrane. <http://cccr.org.cochrane.org/animated-storyboard-what-are-systematic-reviews>.
- [16] OF TASMANIA, U. pico. <https://utas.libguides.com/SystematicReviews/FormulateQuestion>.
- [17] O’MARA-EVES, A., THOMAS, J., MCNAUGHT, J., MIWA, M., AND ANANIADOU, S. Using text mining for study identification in systematic reviews: a systematic review of current approaches. *Systematic Reviews* 4, 1 (Jan 2015), 5.
- [18] PUBMED. pubmed. <https://www.ncbi.nlm.nih.gov/pubmed/>.
- [19] ROSE, S., ENGEL, D., AND CRAMER, N. Automatic keyword extraction from individual documents.
- [20] SATOPAA, V., ALBRECHT, J., IRWIN, D., AND RAGHAVAN, B. Finding a “kneedle” in a haystack: Detecting knee points in system behavior.
- [21] SCIPY. leastsquares. [https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.least\\_squares.html#scipy.optimize.least\\_squares](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.least_squares.html#scipy.optimize.least_squares).
- [22] SINGH, G., MARSHALL, I., THOMAS, J., AND WALLACE, B. Identifying diagnostic test accuracy publications using a deep model. *Working Notes of CLEF* (2017), 11–14.