

Mining Academic Expertise from Funded Research (Expert Search System)

Peeranat Fupongsiripan 2056647

March 11, 2014

Abstract

This project is concerned with development of mining academic expertise in Scottish universities from funded research. This system can be used to extract, analyse and find experts in particular areas in Scottish universities. <http://experts.sicsa.ac.uk/> is an existing academic search engine that assists in identifying the relevant experts within Scottish Universities, based on their recent publication output. The aim of this project is to develop mining tools for the data, and research ways to integrate it with existing deployed academic search engines to obtain the most effective search results. Most of the project's aims and requirements were satisfied. However, various difficulties were encountered in the late stages of the development life-cycle which ultimately led to a reduction of the system's scope.

Acknowledgements

Thanks to Dr. Craig Macdonald for his advice and supervision throughout this project.

Contents

1	Introduction	2
1.1	The Scottish Informatics and Computer Science Alliance(SICSA)	2
1.2	What is Expert Search?	3
1.3	Definition of Mining Academic Expertise from Funded Research and Aims	3
1.4	Context	3
1.5	Overview	4
2	Requirements Specification	4
2.0.1	Must Have	4
3	Background	5
3.1	Information Retrieval (IR)	5
3.1.1	Brief Overview of Information Retrieval System Architecture	5
3.1.2	Retrieving results in Information Retrieval	9
3.1.3	Evaluation	9

3.1.4	Precision and Recall	10
3.2	Search Engine	10
3.3	Learning to Rank	11
3.3.1	Query Dependent Feature	11
3.3.2	Query Independent Feature	11
3.3.3	Obtaining and Deploying a Learned Model	12
3.4	Tools	13
3.4.1	Terrier	13
3.4.2	RankLib	13
3.4.3	trec_eval	13
3.5	Expert Search	13
3.5.1	Determining Good Query Results	13
3.5.2	What makes an expert a good expert?	15
3.5.3	Voting Technique	15
3.6	Design and Implementation	16
3.6.1	Overview of Important Class Diagrams	16
3.6.2	Data Extraction	18
3.6.3	Indexing	20
3.6.4	Retrieving Documents (Experts) with respect to a Query	21
3.6.5	Producing a Learned Model	21
3.6.6	Applying a Learned Model	24
3.6.7	Differences Between Old and New Systems	24

1 Introduction

1.1 The Scottish Informatics and Computer Science Alliance(SICSA)

“The Scottish Informatics and Computer Science Alliance (SICSA) is a collaboration of Scottish Universities whose goal is to develop and extend Scotland’s position as a world leader in Informatics and Computer Science research and education” [7]. SICSA achieves this by working cooperatively rather than competitively, by providing support and sharing facilities, by working closely with industry and government and by appointing and retaining world-class staff and research students in Scottish Universities. A list of members of SICSA is given below.

- University of Aberdeen
- University of Abertay
- University of Dundee
- University of Edinburgh
- Edinburgh Napier University
- University of Glasgow
- Glasgow Caledonian University
- Heriot-Watt University

- Robert Gordon University
- University of St Andrews
- University of Stirling
- University of Strathclyde
- University of the West of Scotland

1.2 What is Expert Search?

With the enormous in the number of information and documents and the need to access information in large enterprise organisations, “collaborative users regularly have the need to find not only documents, but also people with whom they share common interests, or who have specific knowledge in a required area” [12, P. 388]. In an expert search task, the users’ need, expressed as queries, is to identify people who have relevant expertise to the need [12, P. 387]. An expert search system is an Information Retrieval [2] system that makes use of textual evidence of expertise to rank candidates and can aid users with their “expertise need”. Effectively, an expert search systems work by generating a “profile” of textual evidence for each candidate [12, P. 388]. The profiles represent the system’s knowledge of the expertise of each candidate, and they are ranked in response to a user query [12, P. 388]. In real world scenario, the user formulates a query to represent their topic of interest to the system; the system then uses the available textual evidence of expertise to rank candidate persons with respect to their predicted expertise about the query.

1.3 Definition of Mining Academic Expertise from Funded Research and Aims

<http://experts.sicsa.ac.uk/> [6] is a deployed academic search engine that assists in identifying the relevant experts within Scottish Universities, based on their recent publication output. However, integrating different kinds of academic expertise evidence with the existing one may improve the effectiveness of the retrieval system. The aim of this project is to develop mining tools for the funded projects, and research ways to integrate them with the existing academic search engines to obtain the most effective search results. The sources of the new evidence, funded projects, are from Grant on the Web [1] and Research Councils UK [5]. To integrate academic funded projects and publications together, Learning to Rank Algorithms for Information Retrieval (IR) are applied in this project.

1.4 Context

This project was initially developed by an undergraduate student a few years ago. It used academic’s publications as an expertise evidence to find experts. I have access to funded projects data in the UK. This data is integrated with existing data to improve the performance of <http://experts.sicsa.ac.uk/> [6]

1.5 Overview

In this dissertation, Section 2 discusses about Requirements Specification which is grouped into functional and non-functional requirements. Section 3 aims to explain the backgrounds of the project to readers. This section is necessary for readers to understand subsequent sections. Section 3.6 discusses about Designs and Implementations of the system which also references to deployed SICSA search system [6]. Section 5 provides results and analysis of the techniques used. This section can be viewed as the most important part of the whole project since it analyses adding expertise evidence improves the relevance of the search or not. The last section is the conclusion of the project.

2 Requirements Specification

Due to the scope of the project it would be impossible to commence without a clear vision of how the end product should function. These requirements were decided on during the early stages of the project. The reasoning behind them comes primarily from a number of sources:

- Research into a previous attempt at old SICSA system.
- Research into what makes a search result a good result.
- Research into learning to rank techniques to enhancing the performance of the retrieval system by integrating different kinds of expertise evidence.
- Discussion with Dr. Craig Macdonald.

However, developing an IR system is researched based. Utilizing learning to rank technique might not produce optimal ranking due to the lack of training dataset. Therefore, The requirements have been split into several sections depending on their necessity.

2.0.1 Must Have

- Extracting funded projects from different 2 sources: Grant on the Web [1] and Research Councils UK [5].
- Integrating publications and funded projects as expertise evidence.
- Mining expertise evidence.
- Utilizing learning to rank techniques with an attempt to enhance the performance of the retrieval system.
- Indexing and Retrieving documents (experts) effectively.
- Able to conclude that applying learning to rank techniques helps improve the performance of the retrieval system.

I will discuss with u about this section tomorrow

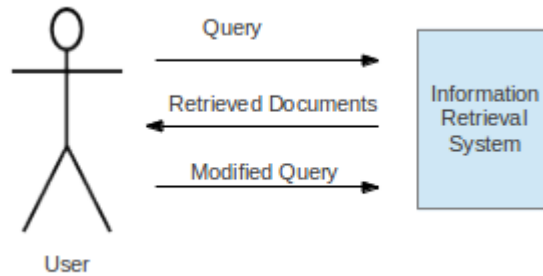


Figure 1: Search Process

3 Background

3.1 Information Retrieval (IR)

“Information Retrieval (IR) is the activity of obtaining information resources relevant to an information need from a collection of information resources” [2]. An information retrieval process begins when a user enters a query into the system. Queries are formal statements of information needs, for example search strings in web search engines. However, the submitted query may not give the satisfying results for the user. In this case, the process begins again. Figure 1 illustrates search process. In information retrieval, a query does not uniquely identify a single object in the collection. Instead, several objects may match the query, with different degrees of relevancy. In IR field, there are various types of retrieval models used to compute the degree of relevancy. This will be discussed in more details in later section.

3.1.1 Brief Overview of Information Retrieval System Architecture

In IR systems, two main objectives have to be met. First, the results must satisfy user - this means retrieving information to meet user’s information need. Second, retrieving process must be fast. This section is devoted to a brief overview of the architecture of IR system which makes readers understand how documents are retrieved and the data structure used in IR systems. To understand how retrieval process works, indexing process must be understood first. This process is done offline. There are 4 steps in indexing process:

- Tokenisation
- Stopwords Removal
- Stemming
- Index Structure Creation

Given a document containing Albert Einstein’s quote about life,

```

<DOC>
<DOCNO>1</DOCNO>
<CONTENT>
There are only two ways to live your life. One is as though nothing
is a miracle. The other is as though everything is a miracle.
</CONTENT>
</DOC>

```

Figure 2: Document

Term	frequency
there	1
are	1
only	1
two	1
ways	1
live	1
your	1
life	1
to	1
one	1
is	3
as	2
though	2
nothing	1
a	2
miracle	2
the	1
other	1
everything	1

Table 1: Terms and Frequency

There are only two ways to live your life. One is as though nothing is a miracle. The other is as though everything is a miracle.

it can be illustrated in a terms-frequency table.

Table 1 shows all the terms and frequency of each term in the document. It can be seen that there are some words in the document which occur too frequently. These words are not good discriminators. They are referred to as “stopwords”. They are articles, prepositions, and conjunctions etc.

Tokenisation is the process of breaking a stream of text into words called tokens(terms).

Stopwords Removal is the process of removing stopwords in order to reduce the size of the indexing structure. This also results in efficient lookup. Table 2 shows all the terms and frequency of each term after stopwords removal process.

Term	frequency
two	1
ways	1
live	1
life	1
one	1
nothing	1
miracle	2
everything	1

Table 2: Terms and Frequency After Stopwords Removal

Term	frequency
two	1
way	1
live	1
life	1
one	1
nothing	1
miracle	2
everything	1

Table 3: Terms and Frequency After Stemming

Now, the table representing the document has shorter size and contains only meaningful words.

Stemming is the process of reducing all words obtained from tokenisation process with the same root into a single root. A stem is the portion of a word which is left after the removal of its affixes(i.e. prefixes and suffixes). For example, connect is the stem for the variants connected, connecting, and connection. This process makes the size of the data shorter.

After stemming, all terms in the table are in its root forms. If a document is large in size, this process can reduce the size of the data. However, there is one drawback. That is, it prevents interpretation of word meanings. For instance, the root form of the term “gravitation” is “gravity”. But the meaning of “gravitation” is different from “gravity”.

Inverted Index Structure Creation is the process that creates an index data structure storing a mapping from terms(keys) to its locations in a database file, or in a document or a set of documents(values) [10]. The purpose of this data structure is to allow a full text searches. In IR, a value in a key-value pair is called posting. Figure 3 shows a simple inverted index. Given a query(a set of terms), it is now possible to efficiently search for documents containing those terms. However, each posting may contain additional information about a document such as the frequency of the term etc.

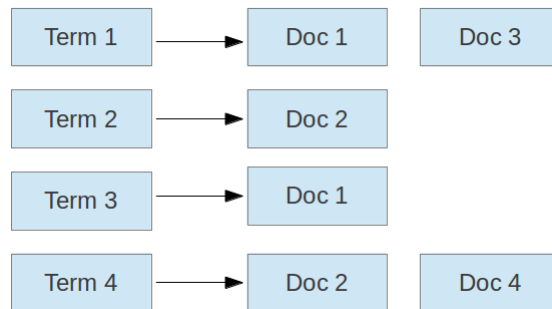


Figure 3: Simple Inverted Index

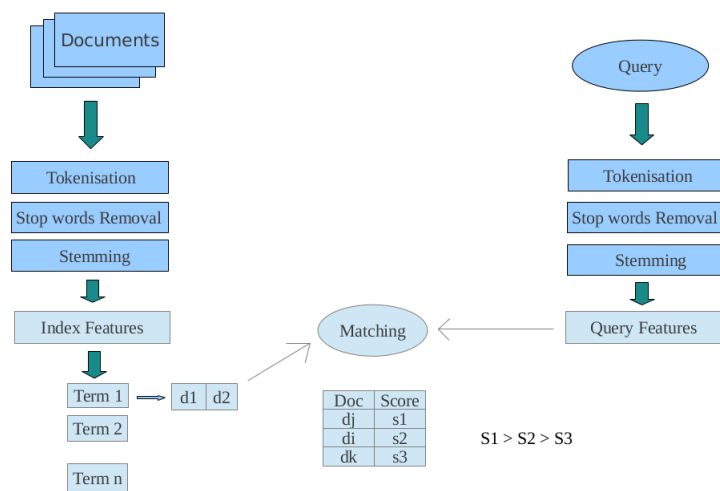


Figure 4: Retrieval Process

3.1.2 Retrieving results in Information Retrieval

Last section, basic IR systems were briefly explained. In this section, we will see how documents are retrieved. Figure 4 shows retrieval process of IR system. In IR, there are 2 phases in general: online and offline phases. The offline phase is the phase that all documents in the corpus are indexed and all features are extracted. It begins with tokenisation, stopwords removal and stemming as explained in the last section. Then for each document, features are extracted. The feature extracting process is not mandatory and application dependent. One most extracted feature is the frequency of the terms occurred in a document. This feature is very useful to weight how important the document is. For instance, if a query term appears in most of the documents in the corpus, it is not appropriate to give a document containing that term a high score. This feature is important if an IR system uses term frequency-inverse document frequency (tf-idf) as a weighting model. But the details of this is out of the scope of this project. Subsequently, inverted index is built.

Online phase begins after a user submits a query into the system. After that, tokenisation, stopwords removal and stemming processes are performed as same as the offline phase. Features can also be extracted from query terms as well. At this point, it comes to the process of matching and assigning scores to documents. Various models can be applied to obtain a score for each document. In this project, special model is used for expert search system. It will be discussed in section 3.5.3. Once scores have been assigned to relevant documents, the system rank all documents in order of decreasing scores and show to the user. Figure 4 gives a graphical representation of retrieval process.

3.1.3 Evaluation

This section is devoted to backgrounds of evaluation of IR systems. It is very important as it is a background for Evaluation Section. Since IR is research-based, understanding how evaluation is carried out will give a background for the readers to determine whether this project is achieved or not. In IR, there are 3 main reasons for evaluating IR systems [14, P. 3]:

1. Economic reasons: If people are going to buy the technology, they want to know how effective it is.
2. Scientific reasons: Researchers want to know if progress is being made. So they need a measure for progress. This can show that their IR system is better or worse than someone else's.
3. Verification: If an IR system is built, it is necessary to verify the performance.

To measure information retrieval effectiveness in the standard way, a test collection is required and it consists of 3 things [9]:

1. A document collection.
2. A test suite of information needs, expressible as queries.
3. A set of relevance judgments, standardly a binary assessment of either relevant or nonrelevant for each query-document pair.

The standard approach to information retrieval system evaluation revolves around the notion of relevant and nonrelevant documents. With respect to a user information need, a document in the test collection is given a binary classification as either relevant or nonrelevant [9]. However, this can be extended by using numbers as an indicator of the degree of relevancy. For example, documents labelled 2 is more relevant than documents labelled 1, or documents labelled 0 is not relevant. There are a number of test collection standards. In this project, Text Retrieval Conference (TREC) is used since it is widely used in the field of IR.

3.1.4 Precision and Recall

The function of an IR system is to [14, P. 10]:

- retrieve all *relevant documents* measured by **Recall**
- retrieve *no non-relevant documents* measured by **Precision**

Precision (P) is the fraction of retrieved documents that are relevant

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant}|\text{retrieved})$$

Recall (R) is the fraction of relevant documents that are retrieved

$$\text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved}|\text{relevant})$$

If a system has high precision but low recall, the system returns relevant documents but misses many useful ones. If a system has low precision but high recall, the system returns most relevant documents but includes lots of junks. Therefore, the ideal is to have both high precision and recall. To give a good example, consider Figure 5, since overall IR system A (blue) has higher precision than IR system B (red), system A is better than system B.

However, in certain cases, precisions of system A may be higher values than system B in some recall points or vice versa. Therefore, some other measure is used instead. Most standard among the TREC community is Mean Average Precision (MAP), which provides a single-figure measure of quality across recall levels. This approach works by averaging all precisions. In this project, MAP is used as a measure.

$$\text{MAP}(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} \text{Precision}(R_{jk})$$

3.2 Search Engine

A search engine is an information retrieval system designed to help find information stored on a computer system. The search results are usually presented in a list and are commonly called hits. Search engines help to minimize the time required to find information and the amount of information which must be consulted, akin to other techniques for managing information overload. The special kind of search engine is web search engine. It is a software system that is designed to search for information on the World Wide Web.

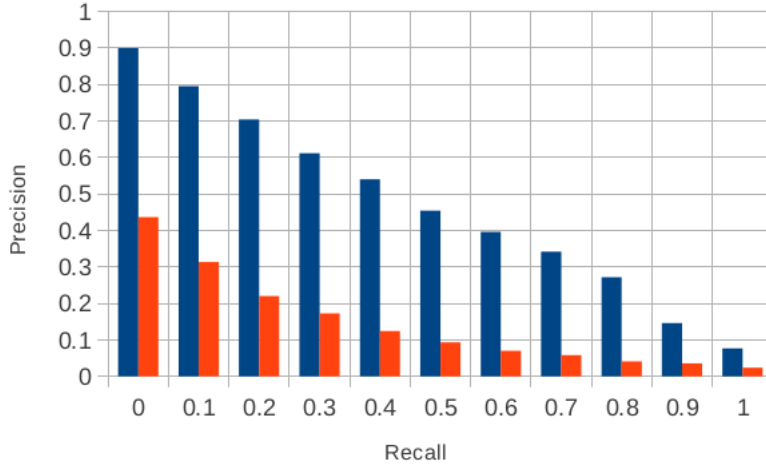


Figure 5: Precision-Recall Graph

3.3 Learning to Rank

Learning to rank or machine-learned ranking (MLR) is a type of supervised or semi-supervised machine learning problem in which the goal is to automatically construct a ranking model from training data [15]. Ranking is the central problem for information retrieval. However, employing learning to rank techniques to learn the ranking function is viewed as a promising approach to information retrieval [15]. In particular, many learning to rank approaches attempt to learn a combination of features (called the learned model) [13, P. 3]. The resulting learned model is applied to a vector of features for each document, to determine the final scores for producing the final ranking of documents for a query [13, P. 3]. In learning to rank, a feature is a binary or numerical indicator representing the quality of a document, or its relation to the query [13, P. 4]. This will be discussed in more details in later section.

3.3.1 Query Dependent Feature

Figure 1 shows a simple search process in IR. After a user submits a query into an IR system. The system ranks the results with respect to the query and returns a result set to the user. It can be clearly seen that the results obtained with respect to the query depends on the query the user submitted. In other words, document A can have 2 different degrees of relevancy if a user changes a query. In learning to rank, this is called query dependent feature.

3.3.2 Query Independent Feature

In contrast to Query Dependent Feature, a feature that does not depend on a user query is called query independent feature. This feature is fixed for each document. For now, it is better to not dig into great details about this because this will be focused in later section.

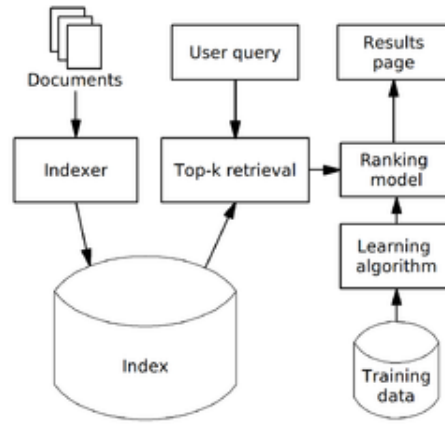


Figure 6: An architecture of a machine-learned IR system from http://en.wikipedia.org/wiki/Learning_to_rank

3.3.3 Obtaining and Deploying a Learned Model

The general steps for obtaining a learned model using a learning to rank technique are the following [13, P. 4]:

1. Top k Retrieval: For a set of training queries, generate a sample of k documents using an initial retrieval approach.
2. Feature Extraction: For each document in the sample, extract a vector of feature values.
3. Learning: Learn a model by applying a learning to rank technique. Each technique deploys a different loss function to estimate the goodness of various combination of features. Documents are labelled according to available relevance assessments.

Once a learned model has been obtained from the above learning steps, it can be deployed within a search engine as follows [13, P. 4]

4. Top k Retrieval: For an unseen test query, a sample of k documents is generated in the same manner as step (1).
5. Feature Extraction: As in step (2), a vector of feature values is extracted for each document in the sample. The set of features should be exactly the same as for (2).
6. Learned Model Application: The final ranking of documents for the query is obtained by applying the learned model on every document in the sample, and sorting by descending predicted score.

Figure 6 illustrates an architecture of a machine-learned IR system. The architecture will be discussed in more details in Design and Implementation Section.

3.4 Tools

3.4.1 Terrier

Every IR system requires programs that handle indexing, retrieving, ranking, etc. To build everything from scratch, it would be impossible within the time range. However, there are a number of search engine platforms that deal with IR functionalities effectively. Terrier [8] is a highly flexible, efficient, and effective open source search engine, readily deployable on large-scale collections of documents. Terrier implements state-of-the-art indexing and retrieval functionalities, and provides an ideal platform for the rapid development and evaluation of large-scale retrieval applications. It is open source, and is a comprehensive, flexible and transparent platform for research and experimentation in text retrieval. Research can easily be carried out on standard TREC and CLEF test collections.

3.4.2 RankLib

RankLib [4] is an open source library of learning to rank algorithms. It also implements many retrieval metrics as well as provides many ways to carry out evaluation. Currently eight popular algorithms have been implemented:

- MART (Multiple Additive Regression Trees, a.k.a. Gradient boosted regression tree)
- RankNet
- RankBoost
- AdaRank
- Coordinate Ascent
- LambdaMART
- ListNet
- Random Forests

However, AdaRank and Coordinate Ascent are only used in the project because other algorithms are too complex and not part of the scope of this project.

3.4.3 trec_eval

trec_eval is the standard tool used by the TREC community for evaluating a retrieval run, given the results file and a standard set of judged results.

3.5 Expert Search

3.5.1 Determining Good Query Results

Figure 1 illustrates the process of search. From this figure, the process will start again if a user is not satisfied with the results. In this section, the focus is on how to convince user that a query result is good. Suppose a user who is

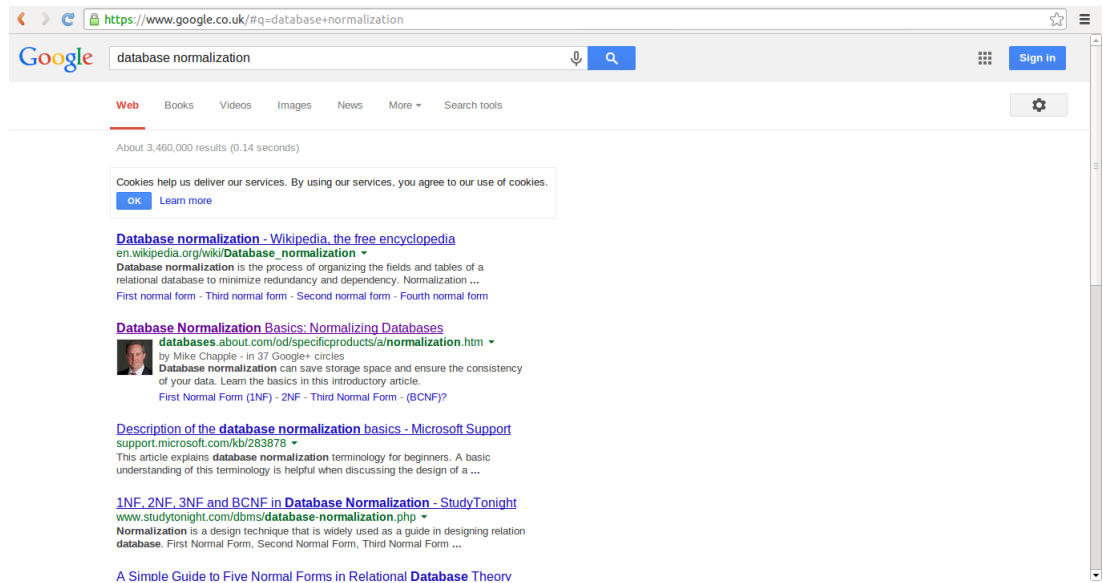


Figure 7: Sample Query

currently studying software engineering would like to know “how to normalize database tables”, first of all, he needs to interpret his information need into a query which is “database normalization”. He then types his query into his preferred search engine. After he submits the query, the search engine gives him a list of results ranked by the degree of relevancy. The question is how does he determine which result is what he is looking for?. Well, he could assume that the ranking provided by the search engine is correct. That is, the first result is what he is looking for. However, this is not always the case. He then explores each result and sees if it is the right one. But without exploring each result, could he be able to determine that which result is likely to satisfy his information need? Perhaps, there has to be some evidence to convince him by just looking at the result. The followings are evidence he could take into account [11]:

- URL
- Photo
- Author
- keywords of the article name

If a result in response to a query have all or some of these evidence, it has more credits than ones with no evidence at all. Figure 7 shows the results of the query “database normalization”. It is obvious that from the top 4 results, all of the article names include the keywords a user submitted, and the third result does not have query keywords included in the URL. Among all of which, the second result has more evidence than others. It has an author’s name, a photo of an author that other results do not.

3.5.2 What makes an expert a good expert?

Last section, it was obvious that some evidence could be used to convince users how reliable a document (link) is. In this section, the evidence within the context of this project will be discussed. As discussed in 1.3, this project makes use of publications and funded projects as expertise evidence. Based on these evidence, what makes an expert a good expert? Well, it is common sense to assume that an expert is professional if

- he has published a lot of publications.
- he has co-authored with a lot of other experts in publications.
- he has co-authored with a lot of other experts in funded projects.
- he has received a lot of funding.
- he has involved in a lot of projects.

It can be seen clearly that these assumptions or features in learning to rank are independent on the query a user submits. These features are query independent (section 3.3.2). However, query dependent features (section 3.3.1) should take into account as well. In this project, there are 2 query dependent features: funded project and publication features. To sum up, a good expert should have high scores in both query dependent and independent features.

3.5.3 Voting Technique

In section 3.1.2, we very briefly talked about weighting model. In other words, how documents are assigned scores using tf-idf. In this section, it aims to give an overview of voting technique used in this project. To understand this section, readers must understand what data fusion technique is. “Data fusion techniques also known as metasearch techniques, are used to combine separate rankings of documents into a single ranking, with the aim of improving over the performance of any constituent ranking” [12, P. 388]. Within the context of this project, expert search is seen as a voting problem. The profile of each candidate is a set of documents associated to him to represent their expertise. When each document associated to a candidate’s profile get retrieved by the IR system, implicit vote for that candidate occurs [12, P. 389]. Data fusion technique is then used to combine the ranking with respect to the query and the implicit vote. In expert search task, it shows that “improving the quality of the underlying document representation can significantly improve the retrieval performance of the data fusion techniques on an expert search task” [12, P. 387]. To give a simple example how data fusion technique works, take a look at this example

Let $R(Q)$ be the set of documents retrieved for query Q , and the set of documents belonging to the profile candidate C be denoted $profile(C)$. In expert search, we need to find a ranking of candidates, given $R(Q)$. Consider the simple example in Tables 4 and 5. The ranking of documents with respect to the query has retrieved documents $\{D1, D2, D3, D4\}$. Using the candidate profiles, candidate $C1$ has accumulated 3 votes, $C2$ 2 votes, $C3$ 2 votes and $C4$ no votes. If all votes are counted equally, and each document in a candidate’s profile is equally weighted, a possible ranking of candidates to this query could be $\{C1,$

Rank	Docs	Scores
1	D1	5.4
2	D2	4.2
3	D3	3.9
4	D4	2.0

Table 4: R(Q)

Profiles	Docs
C1	D3, D4, D2
C2	D1, D2
C3	D3, D2
C4	D5, D6

Table 5: Profiles

C2, C3}. However, in this project, the technique used is expCombMNZ and the formular is as follows

$$candScore(C, Q) = |R(Q) \cap profile(C)| \sum_{d \in R(Q) \cap profile(C)} score_d$$

where

$$|R(Q) \cap profile(C)|$$

is the number of documents from the profile of candidate C that are in the ranking R(Q).

3.6 Design and Implementation

3.6.1 Overview of Important Class Diagrams

As stated in section 1.3, this project extends from SICSA project which uses only publications as expertise evidence. The aim of this project is to integrate funded projects with publications to improve the performance of the retrieval. Figure 8 shows the relationship between Candidate, Project, Publication and University classes. Each component in the figure shows only important attributes and methods. The Candidate class represents an expert who lectures at a university and has a set of publications and projects.

As stated in section 3.5.2, the features extracted from funded projects and publications of each expert will be used in Learning to Rank to improve the retrieval performance of the system. But before we get to Learning to Rank, first of all, we need to understand AcademTechQuerying Class and RetrievedResults Class which are components used in retrieving results with respect to a query.

Figure 9 shows class diagrams of AcademTechQuerying and RetrievedResults. As stated in section 3.4.1, the search engine platform used in this project is Terrier. It handles indexing and retrieval processes discussed in section 3.1.1. The AcademTechQuerying Class makes use of 2 Terrier components as follows:

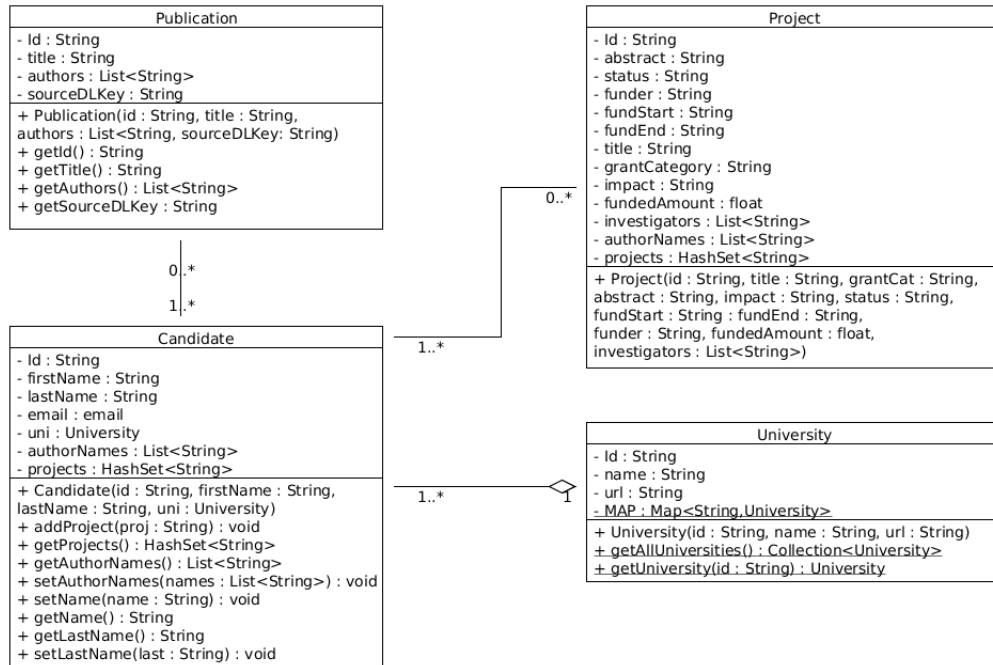


Figure 8: Class Diagram

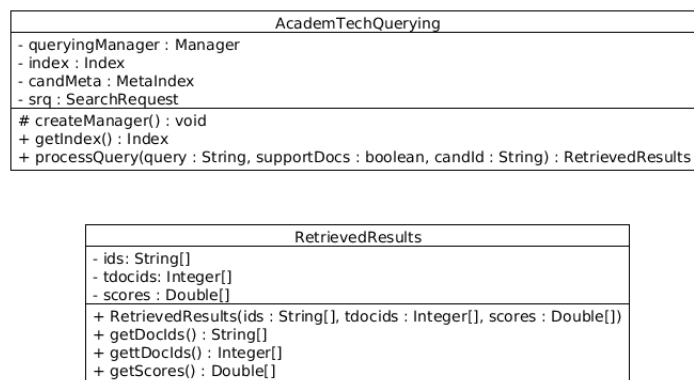


Figure 9: AcademTechQuerying Class and RetrievedResults Class

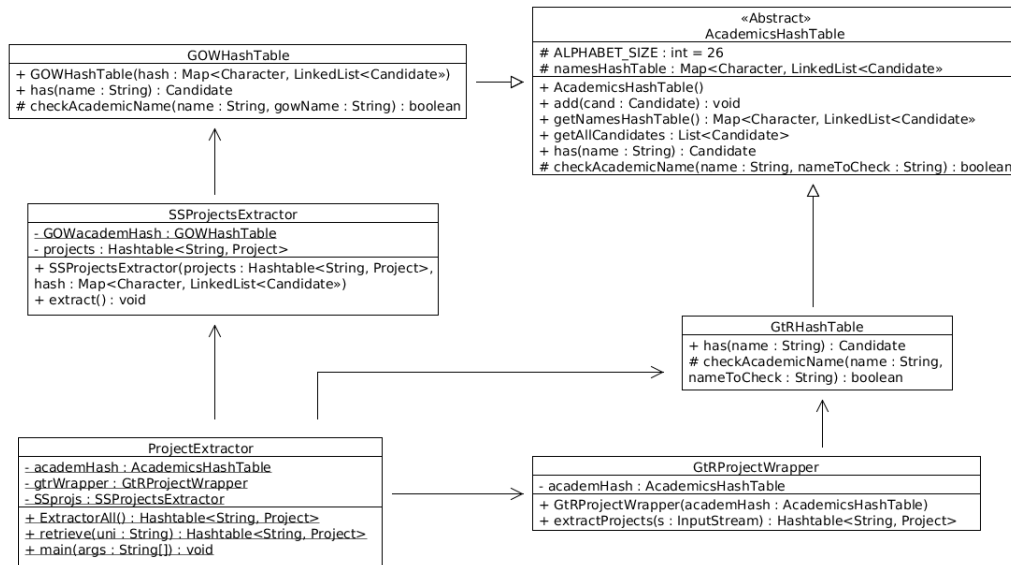


Figure 10: Projects Extraction Class Diagrams

- queryingManager of type Manager, responsible for handling/co-ordinating the main high-level operations of a query.
- srq of type SearchRequest, responsible for retrieving search result from Manager.

The important method of `AcademTechQuerying` is `processQuery()`. It returns `RetrievedResults`. This method takes a query string, `supportDocs` and `candId` as arguments. The first argument tells the system that a query is used to retrieve the result, the second and third arguments are used in case the system wants documents associated to a candidate with respect to a query.

The RetrievedResults Class has 3 attributes as follows:

- `ids`, an array of `String`, which is the ids of the documents(in this case, ids of experts) retrieved by the system.
- `tdocids`, an array of `Integer`, which is the ids used by Terrier
- `scores`, an array of `Double`, which is the scores of each document retrieved by the system.

3.6.2 Data Extraction

In section 1.3, it was stated that funded project information are obtained from Grant on the Web [1] and Research Councils UK [5]. This section shows classes used to extract funded projects from each source and gives statistics regarding the number of total funded projects from each source. Figure 10 shows class diagram related to projects extraction.

Source	Number of Funded Projects
Grant on the Web	32
Research Councils UK	337
	369 total

Table 6: The number of funded projects extracted from each source

AcademicsHashTable is an abstract class which makes use of HashMap data structure, Map<Character, LinkedList<Candidate>>, for efficient look up when matching candidate to our known candidates. It is keyed by the first character of candidate’s name. There are 2 abstract methods in this class: has() and checkAcademicName() methods. Both of them are used together to check if the candidate from a source is matched to our known candidates.

GtRHashTable extends from the abstract class AcademicsHashTable. The purpose of this class is the same as AcademicsHashTable Class but implements has() and checkAcademicName() methods which are suitable for data in Research Councils UK [5].

GOWHashTable extends from the abstract class AcademicsHashTable. Its purpose is similar to AcademicsHashTable Class but has different implementations of hash() and checkAcademicName() methods to GOWHashTable’s which is suitable for Grant on the Web [1] spreadsheet.

SSProjectsExtractor is a class that makes use of GOWHashTable Class to extract funded projects from Grant on the Web [1] spreadsheet.

GtRProjectWrapper is a class used GOWHashTable Class to extract funded projects from Research Councils UK [5].

ProjectExtractor is a class that makes use of SSProjectsExtractor and GtRProjectWrapper Classes to extract funded projects from both sources.

The most difficult part in extracting projects is to match the known expert’s names with the expert’s names in each source. This is because each source records expert’s name in different formats. For example, Prof. Joemon Jose may be recorded Jose JM in one source and Jose J in another. This is why Polymorphism [3] (different implementations of has() and checkAcademicName() methods between GtRHashTable and GOWHashTable) is required. However, pattern matching between expert’s names in Grant on the Web [1] spreadsheet is still impossible as this source records in the following format: [lastname], [role] [the first letter of name] such as Jose, Professor JM. The problem occurs with experts, Dr. Craig Macdonald lecturing at University of Glasgow and Prof. Catriona Macdonald lecturing at Glasgow Caledonian University as the source records Macdonald, Dr C for the former and Macdonald, Professor C for the latter. Although, role might be used to distinguish between them but for some expert the role is not known and there might be the situation that the role is also the same. Therefore, university is used as part of the matching as well.

Figure 6 shows the number of funded projects extracted from each source. There are 1569 known candidates in the system.

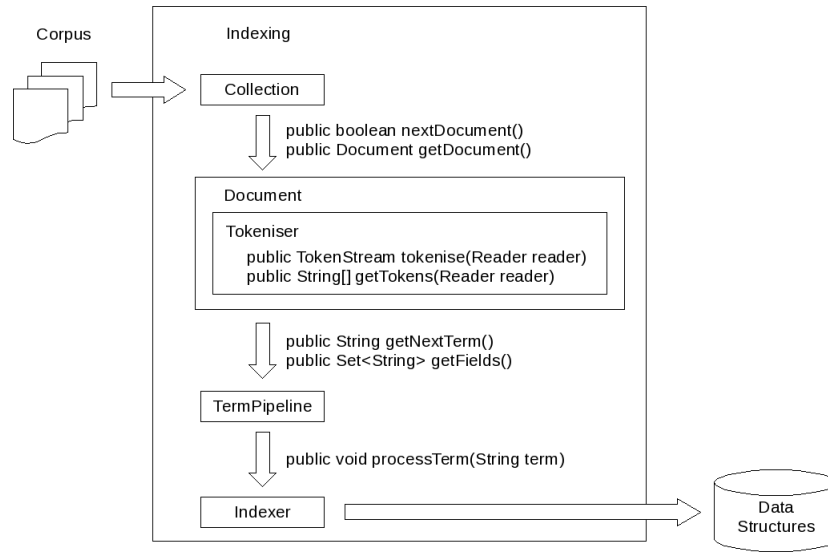


Figure 11: Indexing Architecture of Terrier from <http://terrier.org/docs/v3.5/basicComponents.html>

3.6.3 Indexing

In Section 3.1.1, indexing process is discussed. This process makes the retrieval process much more efficiently. In this section, steps towards indexing using Terrier [8] are discussed.

1. A corpus will be represented in the form of a **Collection** object. Raw text data will be represented in the form of a **Document** object. Document implementations usually are provided with an instance of a **Tokeniser** class that breaks pieces of text into single indexing tokens.
2. The indexer is responsible for managing the indexing process. It iterates through the documents of the collection and sends each found term through a **TermPipeline** component.
3. A TermPipeline can transform terms or remove terms that should not be indexed. An example for a TermPipeline chain is `termpipelines=Stopwords,PorterStemmer`, which removes terms from the document using the **Stopwords** object, and then applies Porter's Stemming algorithm for English to the terms.
4. Once terms have been processed through the TermPipeline, they are aggregated and the following data structures are created by their corresponding DocumentBuilders: `DirectIndex`, `DocumentIndex`, `Lexicon`, and `InvertedIndex`.
5. For single-pass indexing, the structures are written in a different order. Inverted file postings are built in memory, and committed to 'runs' when memory is exhausted. Once the collection had been

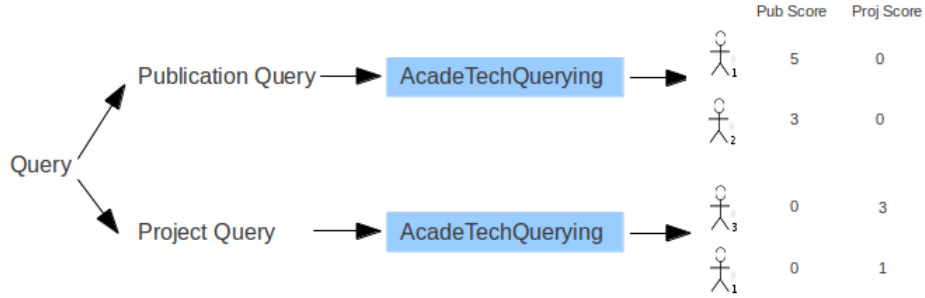


Figure 12: Querying

	Pub Score	Proj Score
Stick Figure 1	5	1
Stick Figure 2	3	0
Stick Figure 3	0	3

Figure 13: Results After Union

indexed, all runs are merged to form the inverted index and the lexicon.

I will discuss with u about this section today.

3.6.4 Retrieving Documents (Experts) with respect to a Query

Figure 12 shows the process of obtaining documents (experts) with respect to a query. First of all, the user query is transformed into publication query and project query. AcademTechQuerying Class then processes both queries to get results with respect to publication query and project query. The results of both queries are then unioned as shown in Figure 13.

3.6.5 Producing a Learned Model

Section 3.3.3 described general steps to produce a learned model. This section aims to discuss classes involved in producing a learned model. As section 3.3.3 suggested, the first step in producing a learned model is to generate a set of training queries. Then all features described in Section 3.5.2 for each document (expert) with respect to each training query are extracted and saved in a file. This file in learning to rank is called LETOR file. Figure 14 is a sample LETOR file. The lines preceded by hash key are ignored. They are just headers which describe features. There are 7 features as described section 3.5.2. The numbers after a hash key indicate features id. However, the lines not preceded by hash

```

1 # 1:publication_query_scores
2 # 2:project_query_scores
3 # 3:candidate_funding_scores
4 # 4:candidate_publication_coauthor_scores
5 # 5:candidate_project_coauthor_scores
6 # 6:candidate_total_project_scores
7 # 7:candidate_total_publication_scores
8 qld:quantum_computation 1:4.71288011140135E-5 2:0.0 3:0.0 4:53.0 5:0.0 6:0.0 7:41.0 # docno = UWS_00000117
9 qld:quantum_computation 1:1.1354922777300101E-5 2:3.246780992565791E-5 3:0.0984859921875 4:29.0 5:1.0 6:1.0 7:18.0 # docno = GLA_00000028
10 qld:quantum_computation 1:0.0 2:1.3775884186774861E-4 3:1.65500515625 4:32.0 5:4.0 6:3.0 7:33.0 # docno = GLA_00000027
11 qld:quantum_computation 1:3.4272271928222563E-4 2:0.0 3:0.0 4:41.0 5:0.0 6:0.0 7:44.0 # docno = GST_00000080
12 # 1:publication_query_scores
13 # 2:project_query_scores
14 # 3:candidate_funding_scores
15 # 4:candidate_publication_coauthor_scores
16 # 5:candidate_project_coauthor_scores
17 # 6:candidate_total_project_scores
18 # 7:candidate_total_publication_scores
19 qld:neural_network 1:0.0025657762336962274 2:2.1956492302765645E-4 3:0.901715 4:201.0 5:2.0 6:1.0 7:207.0 # docno = STI_00000017
20 qld:neural_network 1:0.0 2:0.0014732071662231284 3:0.00825125 4:5.0 5:1.0 6:1.0 7:7.0 # docno = STI_00000018
21 qld:neural_network 1:0.0 2:0.0033273449823246555 3:0.23846859375 4:0.0 5:2.0 6:1.0 7:0.0 # docno = STI_00000009
22 qld:neural_network 1:0.91772344534587E-5 2:0.005505326954511566 3:0.45389303515625 4:44.0 5:2.0 6:2.0 7:32.0 # docno = GST_00000052
23 # 1:publication_query_scores
24 # 2:project_query_scores
25 # 3:candidate_funding_scores
26 # 4:candidate_publication_coauthor_scores
27 # 5:candidate_project_coauthor_scores
28 # 6:candidate_total_project_scores
29 # 7:candidate_total_publication_scores
30 qld:parallel_logic_programming 1:8.772427802562802E-5 2:0.0 3:0.0 4:14.0 5:0.0 6:0.0 7:16.0 # docno = GLA_00000029
31 qld:parallel_logic_programming 1:0.0 2:4.5920546854416106E-4 3:0.0984859921875 4:29.0 5:1.0 6:1.0 7:18.0 # docno = GLA_00000028
32 qld:parallel_logic_programming 1:0.0 2:0.0019483794463139668 3:1.65500515625 4:32.0 5:4.0 6:3.0 7:33.0 # docno = GLA_00000027
33 qld:parallel_logic_programming 1:7.148613480987296E-4 2:0.004298622877033405 3:2.583915 4:34.0 5:7.0 6:4.0 7:34.0 # docno = EDI_00000080
34 qld:parallel_logic_programming 1:0.0 2:1.2250441979824126E-4 3:0.013901 4:45.0 5:1.0 6:1.0 7:44.0 # docno = EDI_00000082
35 qld:parallel_logic_programming 1:3.02577989898064E-4 2:0.0 3:0.0 4:63.0 5:0.0 6:0.0 7:47.0 # docno = NAP_00000053
36 qld:parallel_logic_programming 1:4.347270450240863E-4 2:0.0 3:0.0 4:49.0 5:0.0 6:0.0 7:28.0 # docno = HWU_00000051

```

Figure 14: Sample LETOR file

key are learned. They represent documents (experts) with scores of each feature attached to them. They are preceded by numbers. These numbers in learning to rank are labels which indicate the degree of relevancy. If the label is 0, it means the expert is irrelevant with respect to a query. If it is 1, the expert is relevant. However, the label needs not be binary. It could range from 0 to any positive number. The higher the number, the more relevant the expert is. In this project, labels range from 0 to 2 in order of the degree of relevancy. To the left of the label is query id and scores of each feature associated to that expert. Again, within these lines, anything after hash key is ignored.

34 training queries 7 are trained by the tool called RankLib 3.4.2 to get a learned model. This process is performed only one time.

Figure 15 shows a class diagram that is used to produce a learned model. There are 5 important methods as follows

- loadQueries() is a method that loads all training queries from a file.
- processQuery() is a method that retrieve documents (experts) with respect to publication query and project query.
- setScores() is a method that performs a union between results with respect to publication query and project query as discussed in Section 3.6.4
- setFeatures() is a method that extracts features for each document (expert) and write into a LETOR file.
- learnModel() is a method that produces a learned model using RankLib 3.4.2.

As stated in Section 3.4.2, learning to rank algorithms used in this project are AdaRank and Coordinate Ascent. The performance of each of them will be discussed in Evaluation Section. Figure 16 shows a sample model. Similar to LETOR file, lines preceded by hash key are ignored. They are just headers which describe parameters used in the learning to rank algorithm. Knowing

#	Query	#	Query
1	language modelling	18	game theory
2	manets	19	stable marriage
3	match	20	quantum computation
4	multimodal	21	constraint modelling
5	music as navigation cues	22	home networks
6	networking security	23	wireless sensor networks
7	neural network	24	distributed systems
8	older adults use of computers	25	operating system
9	parallel logic programming	26	terrier
10	query expansion	27	text searching
11	road traffic accident statistics	28	trec collection class
12	shoogle	29	usability
13	skill-based behavior	30	utf support terrier
14	sound in multimedia human-computer interfaces	31	visual impairment
15	statistical inference	32	wafer fab cost
16	suffix tree	33	database
17	mobile hci	34	programming languages

Table 7: Training Queries

LearningModel
<ul style="list-style-type: none"> - aq : AcademTechQuerying - existingQueryRS : Hashtable<String, HashMap<String, Integer> - allCands : String[] - pw : PrintWriter - finalProjScores : double[] - finalPubScores : double[] - candFundingScores : double[] - publicationCoAuthorScores : double[] - projectCoAuthorScores : double[] - totalProjectScores : double[] - totalPublicationScores : double[]
<ul style="list-style-type: none"> + LearningModel() + loadQueries() : void + processQueries() : void + setScores(query : String, pubRS : RetrievedResults, projRS : RetrievedResults) : void + setFeatures(query : String) : void + learnModel() : void + main(args : String[]) : void

Figure 15: Class Diagram of LearningModel

```

1 ## Coordinate Ascent
2 ## Restart = 2
3 ## MaxIteration = 25
4 ## StepBase = 0.05
5 ## StepScale = 2.0
6 ## Tolerance = 0.001
7 ## Regularized = false
8 ## Slack = 0.001
9 1:0.0890911958934695 2:0.068746011602406713 3:0.370221699921495E-4 4:-4.219105023867261E-7 5:-1.3896512810960643E-4 6:-0.0013678721378416296
7:-1.8511097618136217E-4

```

Figure 16: Sample Model

what each parameter does is out of the scope of this project. This model is obtained using Coordinate Ascent Algorithm. The numbers before colons are features id and after colons are scores of each feature. Section 3.6.6 will explain how learned model is applied to get optimal ranking.

3.6.6 Applying a Learned Model

Now that a learned model has been generated, this learned model can be applied to produce optimal ranking. Given a query, scores of query dependent features are computed for each expert as shown in figure 12. After that, scores of query independent features of experts obtained from querying are extracted. Then for each expert, the scores of each feature are multiplied by ones in the learned model and accumulated. Finally, the accumulated scores for each expert are sorted in descending order and experts with the high scores are ranked before those with low scores. This was briefly explained in section 3.3.3.

3.6.7 Differences Between Old and New Systems

Figures 17 18 19 are old interfaces of the system.

Figure 20 shows a new facet in response to a query “information retrieval”. Note that for each expert retrieved, the number of projects the expert has been involved is displayed. Figure 21 shows a new facet of profile page. There are 2 tabs for publications and funded projects and a drop down box for selecting result filter modes. The number of projects and total funding are shown added.

Figure 22 shows a facet when a project tab is selected. The right hand side of the page shows the most project collaborators and other filter modes: filter projects by amount of funding and by the status of projects.

Figure 23 shows a facet when user selects a project link. This shows descriptions of a selected project. The descriptions include:

- Funder
- Title
- Funded Value
- Funded Period
- Status
- Grant Category
- People involved in the project
- Abstract

- Impact

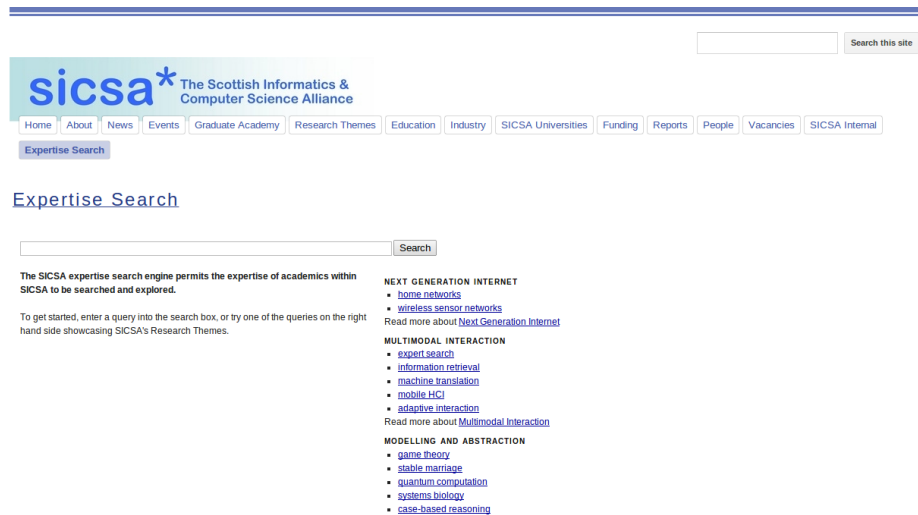


Figure 17: Home Page

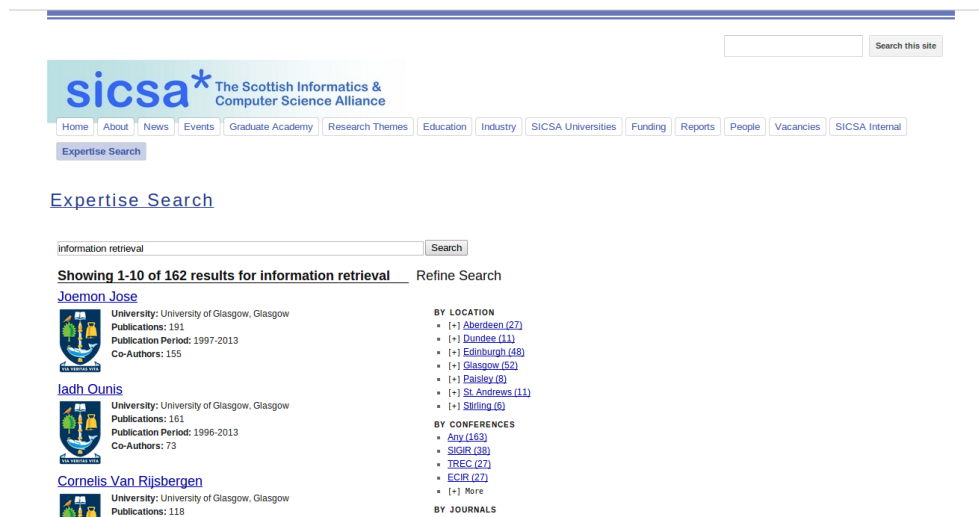


Figure 18: Experts In Response to “information retrieval” Query

Expertise Search

Joemon Jose

University: University of Glasgow, Glasgow
Role:
Homepage:
Publications: 191
Publishing Period: 1997-2013
Co-Authors: 154

Publications

Most Relevant **All**

- Adaptive image retrieval using a Graph model for semantic feature integration. Multimedia Information Retrieval, 2006 [\[Link\]](#)
[Jana Urban](#), [Joemon M. Jose](#)
- A Retrieval Mechanism for Semi-Structured Photographic Collections. DEXA, 1997 [\[Link\]](#)
[Joemon M. Jose](#), [David J. Harper](#)
- Evaluating the implicit feedback models for adaptive video retrieval. Multimedia Information Retrieval, 2007 [\[Link\]](#)
[Frank Hopfgartner](#), [Joemon M. Jose](#)
- A Simulated Study of Implicit Feedback Models. ECIR, 2004 [\[Link\]](#)
[Ryen W. White](#), [Joemon M. Jose](#), [C. J. van Rijsbergen](#), [Ian Ruthven](#)
- Exploiting external knowledge to improve video retrieval. Multimedia Information Retrieval, 2010 [\[Link\]](#)
[David Vallet](#), [Iva?n Cantador](#), [Joemon M. Jose](#)
- Facet-Based Browsing in Video Retrieval: A Simulation-Based Evaluation. MMM, 2009 [\[Link\]](#)
[Frank Hopfgartner](#), [Thierry Urruty](#), [Robert Villa](#), [Joemon M. Jose](#)
- ViGOR: a grouping oriented interface for search and retrieval in video libraries. ICDL, 2009 [\[Link\]](#)

Most Collaborated Co-Authors

- [Joemon M. Jose](#) (191)
- [Frank Hopfgartner](#) (24)
- [Martin Halvey](#) (21)
- [Robert Villa](#) (20)
- [Hideo Joho](#) (20)
- [Ian Ruthven](#) (19)

[\[All Authors\]](#)

Most Popular Terms


approach can evaluation
features feedback image
interfaces model need present
query relevance retrieval
semantic study system task
terms **user** video

[\[More\]](#)

Related Academics

- [Ian Ruthven](#)
- [Martin Halvey](#)


Figure 19: Old Expert's Profile Facet

 **The Scottish Informatics & Computer Science Alliance**


[Home](#) [About](#) [News](#) [Events](#) [Graduate Academy](#) [Research Themes](#) [Education](#) [Industry](#) [SICSA Universities](#) [Funding](#) [Reports](#) [People](#) [Vacancies](#) [SICSA Internal](#)

[Expertise Search](#)


Showing 1-10 of 274 results for information retrieval



Joemon Jose
University: University of Glasgow, Glasgow
Publications: 147
Publication Period: 1997-2013
Publication Co-Authors: 104
Projects: 4



Cornelis Van Rijsbergen
University: University of Glasgow, Glasgow
Publications: 118
Publication Period: 1971-2011
Publication Co-Authors: 70
Projects: 0



Mounia Lalmas
University: University of Glasgow, Glasgow
Publications: 158

Refine Search

BY LOCATION

- [\[+\] Aberdeen \(24\)](#)
- [\[+\] Dundee \(9\)](#)
- [\[+\] Edinburgh \(44\)](#)
- [\[+\] Glasgow \(46\)](#)
- [\[+\] Paisley \(6\)](#)
- [\[+\] St. Andrews \(11\)](#)
- [\[+\] Striding \(6\)](#)

BY CONFERENCES

- [\[+\] Auy \(142\)](#)
- [\[+\] SIGIR \(30\)](#)
- [\[+\] CIKM \(25\)](#)
- [\[+\] AAAI \(23\)](#)
- [\[+\] More](#)

BY JOURNALS

Figure 20: New Facets : Experts In Response to “information retrieval” Query

Joemon Jose

University: University of Glasgow, Glasgow
Role:
Homepage:
Projects: 4
Total Funding: £ 1112755.6
Publications: 147
Publishing Period: 1997-2013
Publication Co-Authors: 103

Publications

Projects

Most Relevant ▾

- Adaptive image retrieval using a Graph model for semantic feature integration. Multimedia Information Retrieval, 2006 [\[Link\]](#)
[Jana Urban](#), [Joemon M. Jose](#)
- A Retrieval Mechanism for Semi-Structured Photographic Collections. DEXA, 1997 [\[Link\]](#)
[Joemon M. Jose](#), [David J. Harper](#)
- Evaluating the implicit feedback models for adaptive video retrieval. Multimedia Information Retrieval, 2007 [\[Link\]](#)
[Frank Hopfgartner](#), [Joemon M. Jose](#)
- Exploiting external knowledge to improve video retrieval. Multimedia Information Retrieval, 2010 [\[Link\]](#)
[David Vallet](#), [Iva?n Cantador](#), [Joemon M. Jose](#)
- Facet-Based Browsing in Video Retrieval: A Simulation-Based Evaluation. MMM, 2009 [\[Link\]](#)
[Frank Hopfgartner](#), [Thierry Urruty](#), [Robert Villa](#), [Joemon M. Jose](#)
- Exploiting log files in video retrieval. JCDL, 2008 [\[Link\]](#)
[Frank Hopfgartner](#), [Thierry Urruty](#), [Robert Villa](#), [Nicholas Gildea](#), [Joemon M. Jose](#)
- Automatic topic detection strategy for information retrieval in spoken document. WIAMIS, 2009 [\[Link\]](#)

Most Collaborated Co-Authors

- [Joemon M. Jose](#) (147)
- [Ian Ruthven](#) (17)
- [Hideo Joho](#) (16)
- [Frank Hopfgartner](#) (16)
- [Robert Villa](#) (14)
- [Reede Ren](#) (13)

[\[All Authors\]](#)

Most Popular Terms

analysis approach can
evaluation features
feedback image models
need performance query
relevance retrieval
semantic study system task
users video web

[\[More\]](#)

Related Academics

- [Mounia Lalmas](#)
- [Martin Halvey](#)
- [Cornelis Van Rijsbergen](#)
- [Yi Yang](#)
- [Iadh Ounis](#)
- [Yuan Liu](#)

Figure 21: New Expert's Profile Facet

Joemon Jose

University: University of Glasgow, Glasgow
Role:
Homepage:
Projects: 4
Total Funding: £ 1112755.6
Publications: 147
Publishing Period: 1997-2013
Publication Co-Authors: 103

Publications

Projects

Most Relevant ▼

- Adaptive Search Models for Information Retrieval [\[Link\]](#). (Closed) Research Grant
[Joemon Jose](#)
- Towards Context-sensitive Information Retrieval Based on Quantum Theory: With Applications to Cross-media Search and Structured Document Access [\[Link\]](#). (Closed) Research Grant
[Joemon Jose](#)
- Community-generated media for the next billion [\[Link\]](#). (Closed) Research Grant
[Joemon Jose](#)
- Towards Context-sensitive Information Retrieval Based on Quantum Theory: With Applications to Cross-media Search and Structured Document Access [\[Link\]](#). (Closed) Research Grant
[Mounia Lalmas](#), [Joemon Jose](#)

Most Collaborators

- [Mounia Lalmas](#)(1)

Funded Amount

- [All](#)
- [Up to £100K](#)
- [£100K to £1M](#)
- [Above £1M](#)

Project Status

- [Active](#)
- [Closed](#)
- [Both](#)

Most Popular Terms

analysis approach can
evaluation features
feedback image models
need performance query
relevance retrieval
semantic study system task
users video web

Figure 22: Project Tab

Joemon Jose

University: University of Glasgow, Glasgow
Role:
Homepage:
Projects: 4
Total Funding: £ 1112755.6
Publications: 147
Publishing Period: 1997-2013
Publication Co-Authors: 103

Publications

Projects

Most Relevant ▼

Towards Context-sensitive Information Retrieval Based on Quantum Theory: With Applications to Cross-media Search and Structured Document Access

Overview

Abstract

People

Funder:	EPSRC
Funded Period:	Sep 1, 2008 - Apr 30, 2011
Funded Value:	£269015.0
Status:	Closed
Grant Category:	Research Grant

Most Collaborators

- Mounia Lalmas(1)

Funded Amount

- All
- Up to £100K
- £100K to £1M
- Above £1M

Project Status

- Active
- Closed
- Both

Most Popular Terms

analysis **approach** can
evaluation features
 feedback image **models**
need performance query
 relevance **retrieval**
 semantic study **system** task
users video web

[\[More\]](#)

Figure 23: Project Page

References

- [1] Engineering and physical sciences research council (epsrc). <http://gow.epsrc.ac.uk/>.
- [2] Information retrieval. http://en.wikipedia.org/wiki/Information_retrieval.
- [3] Polymorphism.
- [4] Ranklib. <http://sourceforge.net/p/lemur/wiki/RankLib/>.
- [5] Research councils uk - gateway to research. <http://gtr.rcuk.ac.uk/>.
- [6] The scottish informatics and computer science alliance expert search system. <http://experts.sicsa.ac.uk/>.
- [7] The scottish informatics and computer science alliance (sicsa). <http://www.sicsa.ac.uk/home/>.
- [8] Terrier. <http://www.terrier.org/>.
- [9] Information retrieval system evaluation, 2009.
- [10] Inverted index, 2012.
- [11] Conversation of dr. craig macdonald, 2014.
- [12] Iadh Ounis Craig Macdonald. Voting for candidates: Adapting data fusion techniques for an expert search task. pages 387, 388, 389, 2006.
- [13] Rdrygo L.T. Craig Macdonald and Iadh Ounis. About learning models with multiple query dependent features. 2012.
- [14] Joemon M Jose. Information retrieval - evaluation methodology. 2013.
- [15] Afshin Rostamizadeh Mehryar Mohri and Ameet Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2012.