

Graphs and Information Retrieval

Proefschrift

ter verkrijging van de graad van doctor
aan de Radboud Universiteit Nijmegen
op gezag van de rector magnificus prof. dr. J.H.J.M. van Krieken,
volgens besluit van het college van decanen
in het openbaar te verdedigen

op woensdag 22 maart 2023
om 12:00 uur precies

door

Chris Frans Henri Kamphuis
geboren op 22 maart 1993 te Oldenzaal, Nederland

Promotor:

prof. dr. ir. A.P. (Arjen) de Vries

Manuscriptcommissie:

Person A (Affiliation)

Person B (Affiliation)

Person C (Affiliation)

This work is part of the research program Commit2Data with project number 628.011.001 (SQIREL-GRAPHS), which is (partly) financed by the Netherlands Organisation for Scientific Research (NWO).

Printed by Drukkerij

Typeset using L^AT_EX

ISBN: 111-11-11111-11-1

Copyright © Chris Kamphuis, 2023

Contents

1	Introduction	1
1.1	Background and Context	1
1.2	Problem Description and Research Questions	2
1.3	Thesis Contributions and Structure	2
1.4	Publications	2
2	Related Work	3
2.1	Information Retrieval	3
2.1.1	Inverted Indexes	3
2.1.2	Ranking methods	3
2.1.3	Similarity Search	3
2.2	Relational Databases	3
2.3	Graphs	4
2.4	Reproducible Science	4
3	IR using Relational Databases	5
3.1	Introduction	6
3.1.1	Boolean retrieval	6
3.1.2	Probabilistic Relational Algebra	6
3.1.3	BOW retrieval models in SQL	6
3.2	Reproducibility	7
3.3	Variants of BM25	8
3.4	Experiments	14
3.5	Results	15
3.6	Do BM25 variants matter?	17
4	From Tables to Graphs	19
4.1	Introduction	19

4.1.1	Current solutions	20
4.1.2	Research Questions	20
4.2	GeeseDB	20
4.3	Design	21
4.3.1	Backend	22
4.4	Graph Query Language	24
4.5	Usage	26
4.5.1	Indexing and Search	28
4.6	Conclusion	33
5	A Graph of Entities	35
5.1	Introduction	35
5.2	Entity Linking	35
5.3	REL	36
5.3.1	Mention Detection	36
5.3.2	Candidate Selection	36
5.3.3	Entity Disambiguation	36
5.4	From REL to REBL	37
5.4.1	Mention Detection	38
5.4.2	Candidate Selection and Entity Disambiguation	39
5.5	Effects on Execution	40
5.6	MMEAD	41
5.7	Results	41
5.8	Conclusion	41
6	Conclusion	45
	Bibliography	47
	Summary	57
	Samenvatting	59
	Acknowledgements	61
	Research Data Management	63
	Curriculum Vitae	65

Chapter 1

Introduction

I *propose* to consider the question, “Can machines think?”

Alan Turing - 1950

1.1 Background and Context

I also propose to consider the question, “Can machines think?” Instead of approaching this through a thought experiment like Turing did, nowadays one can approach this question by asking it to a search engine. When issuing this query to popular web search engines we get varying results: the first result on Google is a passage generated from the article written by Turing, while the first result on Bing is a passage generated from a website that concludes machines can not think¹.

We use systems that process queries every day in our lives when we are looking for *information*. While Google and Bing are all purpose web engines that mainly focus on finding and retrieving information from the internet, people also used specialized search systems in their day-to-day lives: Amazon and EBay for product search, Scholar and ResearchGate for scientific resources, Youtube and TikTok for Videos, or Facebook and LinkedIn when we are searching for people. It might

¹However, if a machine can not think, can we trust the result presented by this algorithm?

even be possible that you are reading this text, after you found this document through search.

When searching for the query “Can machines think?”, the approach of searching through text documents only, might be sufficient for the person who searches. However in many cases when searching today, only considering text is not sufficient. For example, when one wants to buy a product on Amazon, aspects other than text also need to be considered. Lets say for example you want to buy an iPhone; What is the price, which edition is the most recent, or which color does it have.

When someone searches for people on LinkedIn, they are generally more interested in persons that have connections in common compared to complete strangers. If you are looking for someone to do a job, it is ideal that a shared connection can vouch for them. In this case, how people relate to each other in their network might be an indication of *relevance*. Not only the structure on how people relate to each other determines relevance; other examples are their experience, where they work, or reviews on their previous work might matter.

Although it might be possible to encode all this information in written form, often it is more convenient to save this information in a more structured approach. In this work

1.2 Problem Description and Research Questions

1.3 Thesis Contributions and Structure

Chapter 3 will present ..., we will discuss content which has been previously described in the following published works: [27, 25]

Chapter 4 will present ..., we will discuss content which has been previously described in the following published works: [24, 26]

Chapter 5 will present ..., we will discuss content which has been previously described in the following published works: [29]

1.4 Publications

Chapter 2

Related Work

2.1 Information Retrieval

Everything that is needed to process a query like, “Can machines think?”, is subject to research by the field of information retrieval.

2.1.1 Inverted Indexes

2.1.2 Ranking methods

Boolean Retrieval

Vector Space Models

Probabilistic ranking Models

Language Models

Learning to Rank

Vector Space Models revisited

2.1.3 Similarity Search

2.2 Relational Databases

Relational databases are usually used to store structure data.

2.3 Graphs

Instead of using columnar data, it might be more attractive to model your data using graphs.

2.4 Reproducible Science

Chapter 3

IR using Relational Databases

“Is this new question a worthy one to investigate?” This latter question we investigate without further ado, thereby cutting short an infinite regress.

Alan Turing - 1950

Abstract

In this chapter we first look at the history of using relations databases for information retrieval. We see that there have been many attempts to express information retrieval problems using relational databases. We will highlight one of the latter attempts that revived the idea of expressing bag-of-words ranking functions using SQL. A prototype system that makes use of these expressions is presented. We show that this system can be used for rapid IR prototyping, and is especially helpful in the context of reproducible information retrieval research. We implement many different BM25 variants that have been proposed in the literature, and we do not find the expected effectiveness gains we would expect based on the results presented in the work that proposed them.

3.1 Introduction

Where commonly information retrieval researchers use inverted indexes as data structures, there is also history of researchers using relational databases for representing the data in information retrieval systems.

3.1.1 Boolean retrieval

Perhaps the earliest work on using relational databases for information retrieval is the work by Schek and Pistor [51]. In this work the authors recognize that the relational data model is widely accepted as an interface to query structured data, however in cases of unstructured data, like text, it is inconvenient to use. They proposed an extension for the relational model by allowing Non First Normal Form (NF^2) relations. This allows for text queries to be more easily expressed, however the systems that can be build in this language are basically boolean retrieval systems. Which at the time worked well, but scoring was not a feature implemented.

In as similar fashion, Macleod [41], compared the relational model with the inverted index model. He showed how queries of the IBM STAIRS system could be expressed using the relational model. These were however still boolean queries, so scoring using uncertainty was not considered.

3.1.2 Probabilistic Relational Algebra

Fuhr [16] recognized that where databases contain formatted data, IR systems deal with unformatted data and that for this kind of data uncertain inference is required. He proposes to express this uncertainty using a probabilistic relational algebra [17] (PRA). PRA can be considered an extension of the standard relational algebra.

3.1.3 BOW retrieval models in SQL

In a more recent work by Mühleisen et al. [43] showed that the common used BM25 ranking function can also be easily expressed using relational tables. Their work specifically focused on the retrieval efficiency of several systems. They argue that instead of using a custom

Table 3.1: Results presented by Mühleisen et al. [43]; MAP and P@5 on the ClueWeb12 collection are reported for five different systems that run BM25. As shown in the table, only the two database systems achieve the same effectiveness score.

System	MAP	P@5
Indri	0.246	0.304
MonetDB & VectorWise	0.225	0.276
Lucene	0.216	0.265
Terrier	0.215	0.272

build information retrieval system using an inverted index, researchers could just simply store their data representations in a column-oriented relational database, and formulate the ranking functions using SQL. They show that their implementation of BM25 in SQL is on par in efficiency and effectiveness compared to systems that use an inverted index.

There was however an interesting observation in the paper that I would like to highlight: All the systems evaluated in this paper implement BM25, there was however a substantial difference between the effectiveness scores produced by these systems, as shown in table 3.1. The only two systems that achieved the exact same effectiveness score were the two database systems. I should however note that both these systems were however developed by the same research group.

These results came out as quite a surprise as the authors took specific care to keep document pre-processing identical for all systems, but the observed difference in MAP of 3% absolute was the largest deviation in score reported.

3.2 Reproducibility

Not only did we observe the differences in effectiveness scores for BM25 in the paper by Mühleisen et al. [43]. In the SIGIR 2015 Workshop on Reproducibility, Inexplicability, and Generalizability of Results (RIGOR) [3] and the Open-Source IR Replicability Challenge (OSIRRC) workshop [10] similar results are observed. See tables 3.2

Table 3.2: Results from the RIGOR workshop[3], MAP@1000 on the .GOV2 collection is reported for four different systems that run BM25. As shown in the table, all four implementations report a different effectiveness score.

System	MAP@1000
ATIRE	0.2902
Lucene	0.3029
MG4J	0.2994
Terrier	0.2697

Table 3.3: Results from the OSIRRC workshop[10], AP, P@30, and NDCG@20 on the robust04 collection are reported for seven different systems that run BM25. As shown in the table, all implementations report (again) a different effectiveness score.

System	AP	P@30	NDCG@20
Anserini (Lucene)	0.2531	0.3102	0.4240
ATIRE	0.2184	0.3199	0.4211
ielab	0.1826	0.2605	0.3477
Indri	0.2388	0.2995	0.4041
OldDog [25]	0.2434	0.2985	0.4002
Pisa	0.2534	0.3120	0.4221
Terrier	0.2363	0.2977	0.4049

and 3.3, respectively.

It is not clear

3.3 Variants of BM25

We [27] examined several BM25 variants which will be introduced, how the variant varies from the original formulation as proposed by Robertson et al. is marked in red.

Robertson et al. [49]

The original formulation of BM25 consists of three parts, the first part is derived from the binary independence relevance model [47], which results in an approximation of the classical *idf*:

$$w_i^{\text{IDF}} = \log \frac{N - df_t + 0.5}{df_t + 0.5} \quad (3.1)$$

where N is the size of the collection, and df_t are the number of documents in the collection that contain query term t .

Equation 3.2 shows the original formulation of BM25: N is the number of documents in the collection, df_t is the number of documents containing term t , tf_{td} is the term frequency of term t in document d . Document lengths L_d and L_{avg} are the number of tokens in document d and the average number of tokens in a document in the collection, respectively. Finally, k_1 and b are free parameters that can be optimized per collection.

$$\sum_{t \in q} \log \left(\frac{N - df_t + 0.5}{df_t + 0.5} \right) \cdot \frac{tf_{td}}{1 - b + b \cdot \left(\frac{L_d}{L_{avg}} \right) + tf_{td}} \quad (3.2)$$

Lucene (default)

Equation 3.5 shows the variant implemented in Lucene (as of version 8), which introduces two main differences highlighted in red. First, since the IDF component of Robertson et al. is negative when:

$$df_t > \frac{N}{2} \quad (3.3)$$

This happens for terms that appear in more than half of all documents, e.g.: “the” or “a”. Many systems do not consider these terms when searching by keeping a list of common words that can be ignored (stop words). However, if these words are considered, having a negative IDF would decrease the relevance score of documents that have the query term in the document. To avoid negative values in all possible cases, Lucene adds a constant one before calculating the *log* value.

Second, the document length used in the scoring function is compressed (in a lossy manner) to a one byte value, denoted L_{lossy} . With

only 256 distinct document lengths, Lucene can pre-compute the value of

$$k_1 \cdot \left(1 - b + b \cdot \left(\frac{L_{d\text{lossy}}}{L_{avg}} \right) \right) \quad (3.4)$$

for each possible length, resulting in fewer computations at query time.

$$\sum_{t \in q} \log \left(\textcolor{red}{1} + \frac{N - df_t + 0.5}{df_t + 0.5} \right) \cdot \frac{tf_{td}}{k_1 \cdot \left(1 - b + b \cdot \left(\frac{L_{d\text{lossy}}}{L_{avg}} \right) \right) + tf_{td}} \quad (3.5)$$

Lucene (accurate)

Equation 3.6 represents our attempt to measure the impact of Lucene’s lossy document length encoding. We implemented a variant that uses exact document lengths, but is otherwise identical to the Lucene default.

$$\sum_{t \in q} \log \left(\textcolor{red}{1} + \frac{N - df_t + 0.5}{df_t + 0.5} \right) \cdot \frac{tf_{td}}{k_1 \cdot \left(1 - b + b \cdot \left(\frac{L_d}{L_{avg}} \right) \right) + tf_{td}} \quad (3.6)$$

ATIRE [55]

Equation 3.7 shows BM25 as implemented by ATIRE; it implements the IDF component of BM25 as $\log(N/df_t)$, which also avoids negative values. The TF component is multiplied by $k_1 + 1$ to make it look more like the classic Robertson/Sparck Jones (RSJ) weight [48]; this has no effect on the resulting ranked list, as all scores are scaled linearly with this factor.

$$\sum_{t \in q} \log \left(\frac{\textcolor{red}{N}}{\textcolor{red}{df}_t} \right) \cdot \frac{\textcolor{red}{(k_1 + 1)} \cdot tf_{td}}{k_1 \cdot \left(1 - b + b \cdot \left(\frac{L_d}{L_{avg}} \right) \right) + tf_{td}} \quad (3.7)$$

BM25L [38]

BM25L, as shown in equation 3.10, builds on the observation that BM25 penalizes longer documents too much compared to shorter ones. The

IDF component differs, to avoid negative values. The TF component is reformulated as:

$$\frac{(k_1 + 1) \cdot c_{td}}{k_1 + c_{td}} \quad (3.8)$$

with

$$c_{td} = \frac{tf_{td}}{1 - b + b \cdot \left(\frac{L_d}{L_{avg}}\right)} \quad (3.9)$$

The c_{td} component is further modified by adding a constant δ to it, boosting the score for longer documents. The authors report using $\delta = 0.5$ for highest effectiveness.

$$\sum_{t \in q} \log \left(\frac{\mathbf{N} + 1}{\mathbf{df}_t + 0.5} \right) \cdot \frac{(\mathbf{k}_1 + 1) \cdot (c_{td} + \delta)}{k_1 + (c_{td} + \delta)} \quad (3.10)$$

BM25+ [37]

BM25+, as shown in equation 3.11, encodes a general approach for dealing with the issue that ranking functions unfairly prefer shorter documents over longer ones. The proposal is to add a lower-bound bonus when a term appears at least one time in a document. The difference with BM25L is a constant δ to the TF component. The IDF component is again changed to a variant that disallows negative values.

$$\sum_{t \in q} \log \left(\frac{\mathbf{N} + 1}{\mathbf{df}_t} \right) \cdot \left(\frac{(\mathbf{k}_1 + 1) \cdot tf_{td}}{k_1 \cdot \left((1 - b) + b \cdot \left(\frac{L_d}{L_{avg}} \right) \right) + tf_{td}} + \delta \right) \quad (3.11)$$

BM25-adpt [36]

BM25-adpt is an approach that varies k_1 per term (i.e., uses term specific k_1 values). In the original formulation of BM25, k_1 can be considered a hyperparameter that regulates the increase of score for additional occurrences of a term; k_1 ensures that every additional occurrence gets discounted as it provides less information. However, Lv and Zhai argued that this does not necessary have to be the case. If there are much fewer documents that have $t + 1$ occurrences versus t , it should provide more information compared to when the number of

documents are almost equal. In order to find the optimal term-specific k_1 value, the authors want to maximize the information gain for that particular query term. This is done by first identifying the probability of selecting a document randomly from the collection that contains term q at least once in a document as:

$$p(1|0, q) = \frac{df_t + 0.5}{N + 1} \quad (3.12)$$

The probability of a term occurring one more time is defined as:

$$p(t + 1|t, q) = \frac{df_{t+1} + 0.5}{df_t + 1} \quad (3.13)$$

In both these formulas, 1 and 0.5 are added for smoothing to avoid zero probabilities. Then the information gain from t to $t + 1$ occurrences is computed as, subtracting the initial probability:

$$G_q^t = \log_2 \left(\frac{df_{t+1} + 0.5}{df_t + 1} \right) - \log_2 \left(\frac{df_t + 0.5}{N + 1} \right) \quad (3.14)$$

Here df_t is not defined as a normal document frequency, but based on the length normalized term frequency:

$$df_t = \begin{cases} |D_{t|c_{td} \geq t-0.5}| & t > 1 \\ df(q) & t = 1 \\ N & t = 0 \end{cases} \quad (3.15)$$

In this case $df(q)$ is the “normal” document frequency, and c_{td} is the same as in BM25L (pivoted method for length normalization [52]):

$$c_{td} = \frac{tf_{td}}{1 - b + b \cdot \left(\frac{L_d}{L_{avg}} \right)} \quad (3.16)$$

This means the following: df_t is equal to number of documents in the collection when $t = 0$, it is equal to the “normal” document frequency when $t = 1$, and otherwise it will be the number of documents that have at least t occurrences of the term (rounded up) using the pivoted method c_{td} .

Then, the information gain is calculated for $t \in \{0, \dots, T\}$, until $G_q^t > G_q^{t+1}$. This threshold is chosen as a heuristic: When t becomes large, the estimated information gain can be very noisy. So T is chosen

as the smallest value that breaks the worst burstiness rule [9] (the information gain starts decreasing). The optimal value for k_1 is then determined by finding the value for k_1 that minimizes the following equation:

$$k'_1 = \arg \min_{k_1} \sum_{t=0}^T \left(\frac{G_q^t}{G_q^1} - \frac{(k_1 + 1) \cdot t}{k_1 + t} \right)^2 \quad (3.17)$$

Essentially, this gives a value for k_1 that maximizes information gain for that specific term; k_1 and G_q^1 are then plugged into the BM25-adpt formula:

$$\sum_{t \in q} \mathbf{G}_q^1 \cdot \frac{(\mathbf{k}'_1 + 1) \cdot t f_{td}}{\mathbf{k}'_1 \cdot \left((1 - b) + \left(\frac{L_d}{L_{avg}} \right) \right) + t f_{td}} \quad (3.18)$$

We found that the optimal value of k_1 is actually not defined for about 90% of the terms. A unique optimal value for k_1 only exists when $t > 1$ while calculating G_q^t . For many terms, especially those with a low df , $G_q^t > G_q^{t+1}$ occurs before $t > 1$. In these cases, picking different values for k_1 has virtually no effect on retrieval effectiveness. For undefined values, we set k_1 to 0.001, the same as Trotman et al. [56]

TF $l \circ \delta \circ p \times \text{IDF}$ [50]

$\text{TF}l \circ \delta \circ p \times \text{IDF}$, as shown in equation 3.21, models the non-linear gain of a term occurring multiple times in a document as:

$$1 + \log(1 + \log(t f_{td})) \quad (3.19)$$

To ensure that terms occurring at least once in a document get boosted, the approach adds a fixed component δ , following BM25+. These parts are combined into the TF component using the pivoted method for length normalization [52]:

$$c_{td} = \frac{t f_{td}}{1 - b + b \cdot \left(\frac{L_d}{L_{avg}} \right)} \quad (3.20)$$

The same IDF component as in BM25+ is used, which gives us $\text{TF}l \circ \delta \circ p \times \text{IDF}$:

$$\sum_{t \in q} \log \left(\frac{\mathbf{N} + 1}{\mathbf{df}_t} \right) \cdot (\mathbf{1} + \log(\mathbf{1} + \log(\mathbf{c}_{td} + \delta))) \quad (3.21)$$

3.4 Experiments

Our experiments were conducted using Anserini (v0.6.0) on Java 11 to create an initial index, and subsequently using relational databases for rapid prototyping, which we dub “OldDog” [25] after Mühleisen et al. [43]; following that work use MonetDB as well. Evaluations with Lucene (default) and Lucene (accurate) were performed directly in Anserini; the latter was based on previously-released code that we updated and incorporated into Anserini.¹ The inverted index was exported from Lucene to OldDog, ensuring that all experiments share exactly the same document processing pipeline (tokenization, stemming, stopword removal, etc.). While exporting the inverted index, we precalculate all k_1 values for BM25- adpt as suggested by Lv and Zhai [36]. As an additional verification step, we implemented both Lucene (default) and Lucene (accurate) in OldDog and compared results to the output from Anserini. We are able to confirm that the results are the same, setting aside unavoidable differences related to floating point precision. All BM25 variants are then implemented in OldDog as minor variations upon the original SQL query provided in Mühleisen et al. [43]. The term-specific parameter optimization for the adpt variant was already calculated during the index extraction stage, allowing us to upload the optimal (t, k) pairs and directly use the term-specific k values in the SQL query. The advantage of our experimental methodology is that we did not need to implement a single new ranking function from scratch.

The experiments use three TREC newswire test collections: TREC Disks 4 and 5, excluding Congressional Record, with topics and relevance judgments from the TREC 2004 Robust Track (Robust04); the New York Times Annotated Corpus, with topics and relevance judgments from the TREC 2017 Common Core Track (Core17); the TREC Washington Post Corpus, with topics and relevance judgments from the TREC 2018 Common Core Track (Core18). Following standard experimental practice, we assess ranked list output in terms of average precision (AP) and precision at rank 30 (P@30). The parameters shared by all models are set to $k_1 = 0.9$ and $b = 0.4$, Anserini’s defaults. The parameter δ is set to the value reported as best in the corresponding source publication.

¹<http://searchivarius.org/blog/accurate-bm25-similarity-lucene>

Table 3.4: Effectiveness scores different BM25 variants, all were implement as SQL queries, so the underlying data representations are exactly the same.

	Robust04		Core17		Core18	
	AP	P@30	AP	P@30	AP	P@30
Robertson et al.	.2526	.3086	.2094	.4327	.2465	.3647
Lucene (default)	.2531	.3102	.2087	.4293	.2495	.3567
Lucene (accurate)	.2533	.3104	.2094	.4327	.2495	.3593
ATIRE	.2533	.3104	.2094	.4327	.2495	.3593
BM25L	.2542	.3092	.1975	.4253	.2501	.3607
BM25+	.2526	.3071	.1931	.4260	.2447	.3513
BM25-adpt	.2571	.3135	.2112	.4133	.2480	.3533
TF $l \circ \delta \circ p \times$ IDF	.2516	.3084	.1932	.4340	.2465	.3647

All experiments were run on a Linux desktop (Fedora 30, Kernel 5.2.18, SELinux enabled) with 4 cores (Intel Xeon CPU E3-1226 v3 @ 3.30 GHz) and 16 GB of main memory; the MonetDB 11.33.11 server was compiled from source using the `--enable-optimize` flag.

3.5 Results

Table 3.4 shows the effectiveness scores of the different BM25 variants.

The observed differences in effectiveness are very small and can be fully attributed to variations in the scoring function; our methodology fixes all other parts of the indexing pipeline (tag cleanup, tokenization, stopwords, etc.). Both an ANOVA and Tukey’s HSD show no significant differences between any variant, on all test collections. This confirms the findings of Trotman et al. [56]: effectiveness differences are unlikely an effect of the choice of the BM25 variant. Across the IR literature, we find that differences due to more mundane settings (such as the choice of stopwords) are often larger than the differences we observe here. Although we find no significant improvements over the original Robertson et al. [49] formulation, it might still be worthwhile to use a variant of BM25 that avoids negative ranking scores.

You might have caught that the effectiveness scores of both ATIRE

Table 3.5: Average retrieval time per query in ms: Anserini (top) and OldDog (bottom)

	Robust04	Core17	Core18
Lucene (default)	52	111	120
Lucene (accurate)	55	115	123
Robertsen et al.	158 ± 25	703 ± 162	331 ± 96
Lucene (default)	157 ± 24	699 ± 154	326 ± 90
Lucene (accurate)	157 ± 24	701 ± 156	324 ± 88
ATIRE	157 ± 24	698 ± 159	331 ± 94
BM25L	158 ± 25	697 ± 160	333 ± 96
BM25+	158 ± 25	700 ± 160	334 ± 96
BM25-adpt	158 ± 24	700 ± 157	330 ± 92
TF _{l_od_op} × IDF	158 ± 24	698 ± 158	331 ± 96

and Lucene (accurate) are exactly the same. This is not a mistake. As explained, the $k_1 + 1$ in ATIRE scales the scores linearly and does not effect the ranking. So the only difference that can change the effectiveness scores are the different IDF functions, but these are practically the same, especially when a collection has a large number of documents (N):

$$\log\left(\frac{N}{df_t}\right) = \log\left(\frac{N - df_t + df_t}{df_t}\right) \quad (3.22)$$

$$= \log\left(\frac{N - df_t}{df_t} + \frac{df_t}{df_t}\right) \quad (3.23)$$

$$= \log\left(\frac{N - df_t}{df_t} + 1\right) \quad (3.24)$$

$$\approx \log\left(\frac{N - df_t + 0.5}{df_t + 0.5} + 1\right) \quad (3.25)$$

Table 3.5 presents the average retrieval time per query in milliseconds (without standard deviation for Anserini, which does not report time per query). MonetDB uses all cores for both inter- and intra-query parallelism, while Anserini is single-threaded.

Comparing Lucene (default) and Lucene (accurate), we find negligible differences in effectiveness. However, the differences in retrieval time are also negligible, which calls into question the motivation behind the original length approximation. Currently, the similarity function and thus the document length encoding are defined at index time. Storing exact document lengths would allow for different ranking functions to be swapped at query time more easily, as no information would be discarded at index time. Accurate document lengths might additionally benefit downstream modules that depend on Lucene. We therefore suggest that Lucene might benefit from storing exact document lengths.

3.6 Do BM25 variants matter?

In summary, the previous sections describe a double reproducibility study: it methodologically validates the usefulness of databases for IR prototyping claimed by Mühleisen et al. [43] and performed a large-scale study of BM25 to confirm the findings of Trotman et al. [56]. Returning to our original motivating question regarding the multitude of BM25 variants: “Does it matter?”, we can conclude that the answer appears to be “no, it does not”.

Chapter 4

From Tables to Graphs

Abstract

This is the abstract

4.1 Introduction

In recent years there has been a lot of exciting new information retrieval research that makes use of non-text data to improve the effectiveness of search systems. Consider for example dense representations for retrieval [18, 35, 33], knowledge graphs to leverage entity information [21, 5, 11], and non-textual learning-to-rank features [12, 39]. All these research directions have improved the effectiveness of search systems by making use of more diverse data. Despite the fact that search systems consider more diverse sources of data, the usage of this data is often implemented through the use of a coupled architecture. In particular, first-stage retrieval is often carried out with different software compared to later retrieval stages where these novel reranking techniques tend to be used. In our view, researchers could benefit from a system where retrieval stages are more tightly coupled, that facilitates the exploration on how to use non-content data for ranking, and serves the data in a format suitable for reranking with e.g. transformers or tree based methods.

4.1.1 Current solutions

PyTerrier was developed by Macdonald et al. [40], which is a Python extension of Terrier [44]. Pyserini was developed by Lin et al. [31], which is a Python extension of Anserini [59].

4.1.2 Research Questions

In order to fulfill these needs we propose GeeseDB¹, a prototype Python toolkit for information retrieval that leverages graphs as data structures, allowing metadata and graph oriented data to be easily included in the ranking pipeline. The toolkit is designed to quickly set up first stage retrieval, and make it easy for researchers to explore new ranking models quickly.

4.2 GeeseDB

In short, GeeseDB aims to provide the following functionalities:

- GeeseDB is an easy to install, self-contained Python package available through `pip install` with as few as possible dependencies. It contains topics and relevance judgements for several standard IR collections out-of-the-box, allowing researchers to quickly start developing new ranking models.
- First stage (sparse) retrieval is directly supported. In only a few lines of code it is possible to load documents and create BM25 rankings.
- Data is served in a usable format for later retrieval stages. GeeseDB allows to directly run queries on Pandas data frames for efficient data transfer to sequential reranking algorithms.
- Easy data exploration is supported through querying data with SQL, but more interestingly, also using a graph query language (based on the Cypher query language), making the exploration of new research avenues easier. This prototype supports a subset of the graph query language Cypher, similar to the property graph database model query language as described by Angles [2].

¹<https://github.com/informagi/geesedb>

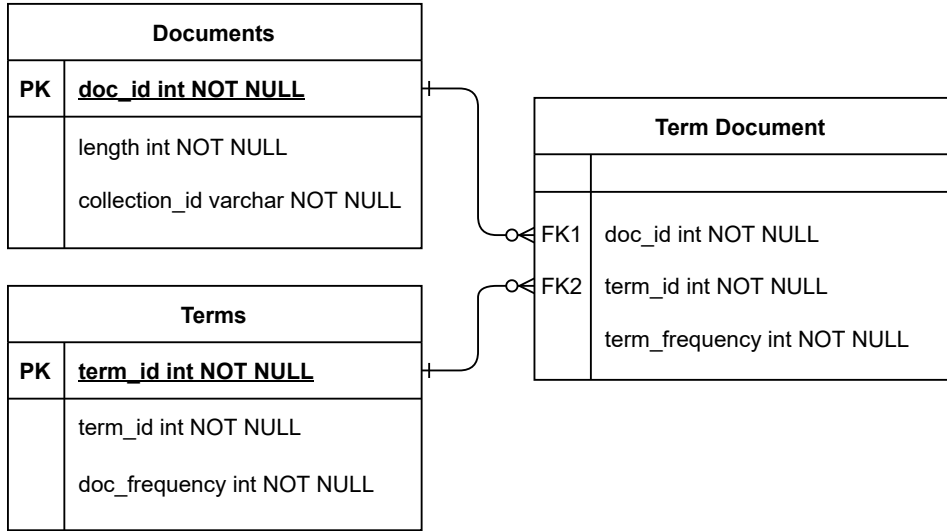


Figure 4.1: Database schema by Mühleisen et al. for full text search in relational databases

4.3 Design

At the core of GeeseDB lies the full text search design presented by Mühleisen et al. [43]. In this work, a column-store database for IR prototyping is proposed, which uses the database schema described in Figure 4.1, consisting of three database tables. (One for all term information, one for all document information, and one that contains the information on how terms relate to documents; the information that is found in a posting list of an inverted index). Using these three tables they show that BM25 can be easily expressed as a SQL query, with latencies that are on par with custom-build IR engines. In GeeseDB we use the exact same relational schema for full text search. Instead of seeing the document data and term data as tables that relate to each other through a many-to-many join table, it is also possible to consider this schema as a bipartite graph. In this graph both documents and terms are considered as nodes, connected to each other through edges. Basically, if a term occurs in a document there exists an edge between that term and document. GeeseDB uses the data model of property graphs; labeled multigraphs where both edges and nodes can have property-value pairs. The database schema as described in

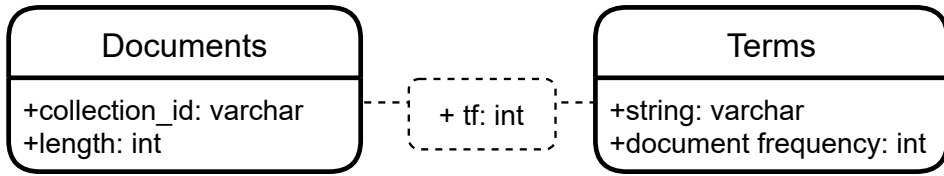


Figure 4.2: Graph schema representing bipartite document-term graph

Figure 4.1 would then translate to the property graph schema shown in Figure 4.2. A small example of a graph represented by this schema is shown in Figure 4.3, document nodes contain document specific information (i.e. document length and the collection identifier), term nodes contain information relevant to the term (i.e. the term string and the term’s document frequency), and the edges between document and terms nodes contain term frequency information (i.e. how often is the term mentioned in the document represented the respective nodes it connects). If one wants to for example also store position data, this graph can easily be changed to a graph where the edges store the position of a term. If a term would appear multiple times in a document, the property graph model would allow for multiple edges to exist between two nodes. The graph schema that we described by Figure 4.2 maps one-to-one to the relational database schema described by Figure 4.1, so nodes are represented by normal relational tables that represent specific data units (terms, documents), while edges are represented by many-to-many join tables. So even though we think of the data as graphs, in the backend they are represented as relational tables. When using GeeseDB for search we at least expect the document-term graph to be present, of course new node types can be introduced in order to explore new search strategies.

4.3.1 Backend

GeeseDB is built on top of DuckDB [45], an in-process SQL OLAP (analytics optimized) database management system. DuckDB is designed to support analytical query workloads, meaning that it specifically aims to process complex long-running queries where a significant portion of the data is accessed, conditions matching the case of IR research. DuckDB has a client Python API which can be installed using `pip`, afterwards it can be used directly. DuckDB has a separated API

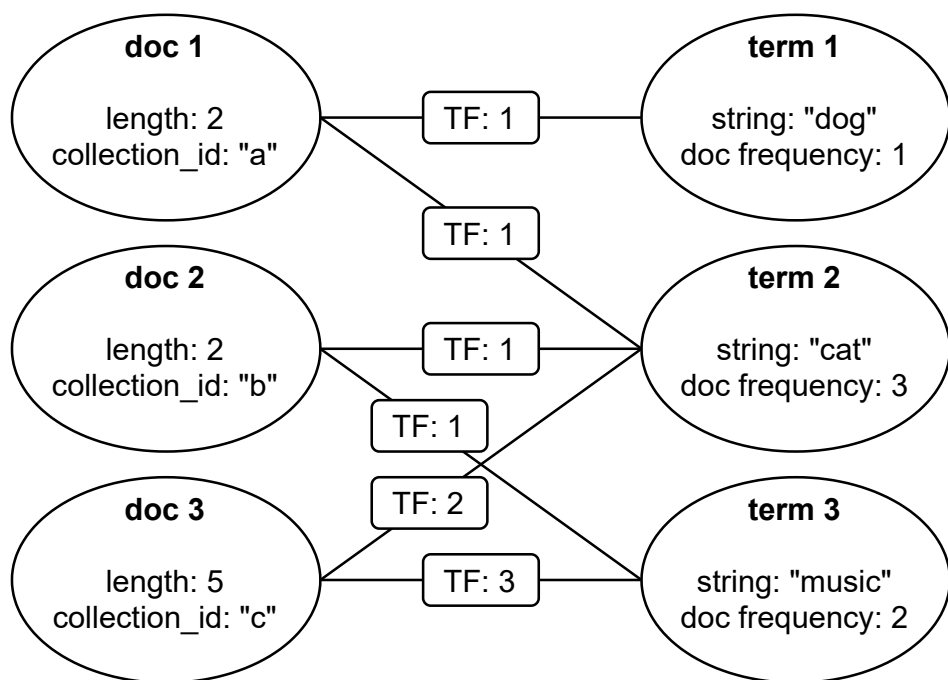


Figure 4.3: Example term-document graph that maps to relational database schema

```
1 MATCH (d:docs)-[]-(a:authors)-[]-(d2:docs)
2 WHERE d.collection_id = "96ab542e"
3 RETURN DISTINCT d2.collection_id
```

Figure 4.4: An example cypher query that finds all documents that were written by the same author that wrote the document with the `collection_id` “96ab542e”

built around both NumPy and Pandas, providing NumPy/Pandas views over the same underlying data representation, without incurring data transfer (usually referred to as “zero-copy” reading). Pandas DataFrames can be registered as virtual tables, allowing to directly query the data present in Pandas DataFrames. GeeseDB inherits all these functionalities from DuckDB.

As DuckDB is a SQL database management system, we can execute analytical SQL queries on the tables that contain our data, including the BM25 rankings described by Mühleisen et al. [43]. By default, the BM25 implementation provided with GeeseDB implements the disjunctive variant of BM25, instead of the conjunctive variant they used. Although the conjunctive variant of BM25 can be calculated more quickly, we chose to use the disjunctive variant as it is more commonly used by IR researchers and the differences between effectiveness scores are noticeable on smaller collections. For now we only support the original formulation of BM25 by Robertson et al. [49], however support of or adding other versions of BM25 [27] is trivial.

4.4 Graph Query Language

What distinguishes GeeseDB from alternatives, database-backed (Old-Dog) [25] or native systems (Anserini [59], Terrier [44]) is the graph query language, based on Cypher [15]. For now, GeeseDB implements Cypher’s basic graph pattern matching queries for retrieving data. An example of a graph query supported by GeeseDB is presented in Figure 4.4. This query finds all documents written by the same authors as those who wrote document “96ab542e”. For comparison, Figure 4.5 illustrates the same query represented in SQL; much more complex than the Cypher version, due to the join conditions that have to be

```

1 SELECT DISTINCT d2.collection_id
2 FROM docs AS d2
3 JOIN doc_author AS da2 ON (d2.collection_id = da2.doc)
4 JOIN authors AS a2 ON (da2.author = a2.author)
5 JOIN doc_author AS da3 ON (a2.author = da3.author)
6 JOIN docs AS d ON (d.collection_id = da3.doc)
7 WHERE d.collection_id = '96ab542e'

```

Figure 4.5: SQL query that corresponds to the graph query described in Figure 4.4.

```

1 MATCH (d:docs {d.collection_id: "96ab542e"})
2 RETURN d.len

```

Figure 4.6: Graph query where the length of document with collection_id is returned.

made explicit. In order to connect the “docs” table with the “authors” table 2 joins are needed, reconnecting the “docs” table again introduces two more joins.

At the moment of writing, GeeseDB supports the following Cypher keywords: **MATCH**, **RETURN**, **WHERE**, **AND**, **DISTINCT**, **ORDER BY**, **SKIP**, and **LIMIT**. Instead of using **WHERE** to filter data, it is also possible to use graph matching, as shown in Figure 4.6; the query returns the length of document “96ab542e”. We plan to support the other keywords of Cypher in the future, as well as directed edges. Everything that is not yet directly supported yet by our implementation can of course still be expressed in SQL, which is fully supported.² In order to know how to join nodes to each other if no edge information has been provided, GeeseDB stores information on the schema. This way GeeseDB knows how nodes relate to each other through which edges. GeeseDB has a module for updating the graph schema, allowing researchers to easily set up the graph they want represented in the database.

²GeeseDB supports the graph queries by translating them to their corresponding SQL queries, both nodes and edges are after all just tables in the backend.

4.5 Usage

GeeseDB comes as an easy-to-install Python package that can be installed using pip, the standard package installer for Python:

```
$ pip install geeseadb==0.0.1
```

After installing GeeseDB we can immediately start using it. All examples we show in this paper were run on version v0.0.1 of GeeseDB. However, as GeeseDB is actively being developed we advise readers to use the latest version of GeeseDB, which can be installed when not specifying a package version. It is also possible to install the latest commit by installing the latest version directly from GitHub. As an example, we will show how to use GeeseDB for the background linking task of the TREC News Track [53]. The goal of this task is: *Given a news story, find other news articles that can provide important context or background information.* These articles can then be recommended to the reader to help them understand the context in which these news articles take place. The collection used for this task is the Washington Post V3 collection³ released for the 2020 edition of TREC. It contains 671,945 news articles published by the Washington Post published between 2012 and 2020, and 50 topics with relevance assessments (topics correspond to collection identifiers of documents for which relevant data has to be found). The articles in this collection contain useful metadata; in particular, we will use authorship information. We extracted 25,703 unique article authors, where it is possible that multiple authors co-wrote a news article. We also annotate documents with entity information which was obtained by using the Radboud Entity Linker [57]. In total 31,622,419 references to 541,729 unique entities were found, the links also contain mention and location information, as well as the `ner_tag` found by the linker’s entity recognition module (The `ner_tag` is part of a link, as the entity linker can assign different tags to the same entity).⁴ Figure 4.7 illustrates the data schema that we use for the background linking task.

³<https://trec.nist.gov/data/wapost/>

⁴The annotated data will be made publicly available.

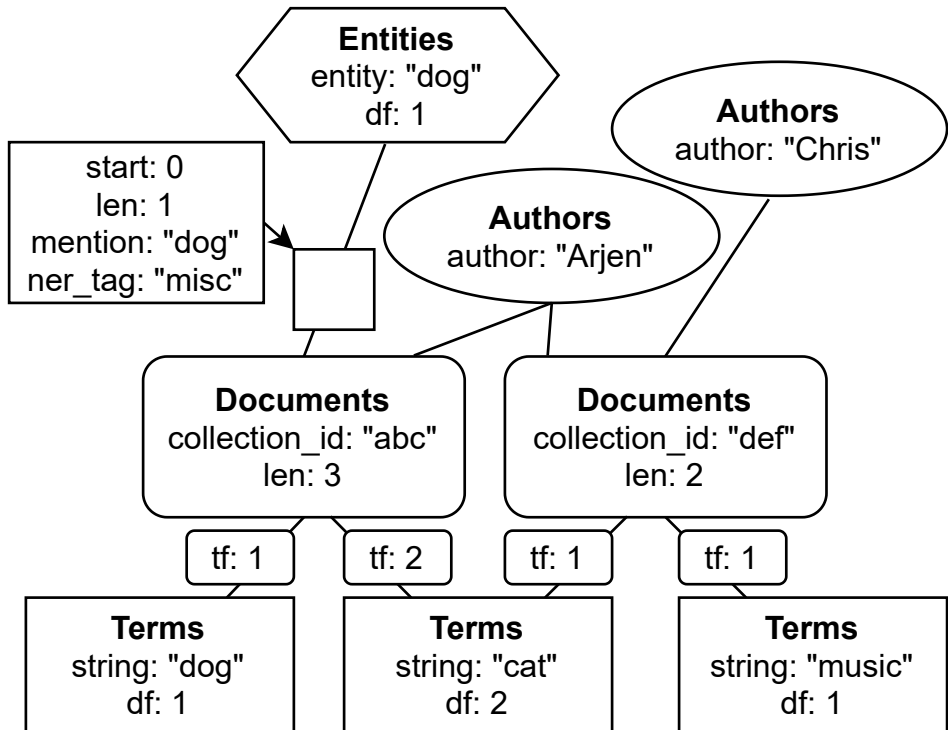


Figure 4.7: Example property graph for the TREC News Track's background linking task. The node types are authors, entities, terms and documents. Edges connect document nodes to other types of nodes. Both edges and nodes can have properties (following the property graph model). Multiple edges may exist between one entity node and one document node, as one entity can be linked multiple times to one document.

```
1 from geesedb.index import FullTextFromCSV
2
3 index = FullTextFromCSV(
4     database='/path/to/database',
5     docs_file='/path/to/docs.csv',
6     term_dict_file='/path/to/term_dict.csv',
7     term_doc_file='/path/to/term_doc.csv'
8 )
9 index.load_data()
```

Figure 4.8: Load text data from the WashingtonPost collection formatted as csv files in the format as described by Mühleisen et al. [43]

```
1 from geesedb.search import Searcher
2
3 searcher = Searcher(
4     database='/path/to/database',
5     n=10
6 )
7 hits = searcher.search_topic('obama and trump')
```

Figure 4.9: Example on how to create a BM25 ranking for the query “obama and trump” that returns the top 10 documents.

4.5.1 Indexing and Search

In order to start, a database containing at least the document and term information needs to be created. Figure 4.8 shows how the data can be easily loaded using CSV files.

Instead of loading the data from CSV files it is also possible to load the text data directly using the CIFF format for data exchange [32]. GeeseDB also has functionalities to create the CSV files used here from the CIFF format. Authorship information and entity links can be loaded similarly. After loading the data we can quickly create a BM25 ranking for ad hoc search in the Washington Post collection as shown in Figure 4.9.

For the background linking task however, we do not have regular topics; we only have the collection identifiers of the documents we need


```
1 MATCH (d:docs {collection_id: ?})-[]-(t:term_dict)
2 RETURN string
3 ORDER BY tf*log(671945/df)
4 DESC
5 LIMIT 5
```

Figure 4.10: Prepared Cypher statement that finds the top-5 TF-IDF terms in a document.

```
1 SELECT term_dict.string
2 FROM term_dict
3 JOIN term_doc ON
4     (term_dict.term_id = term_doc.term_id)
5 JOIN docs ON
6     (docs.doc_id = term_doc.doc_id)
7 WHERE docs.collection_id = ?
8 ORDER BY term_doc.tf * log(671945/term_dict.df)
9 DESC
10 LIMIT 5;
```

Figure 4.11: Prepared SQL statement that finds the top-5 TF-IDF terms in a document.

to find relevant background info for. In order to search for relevant background reading, queries that represent our information need to be constructed. A common approach is to use the top- k TF-IDF terms of the source article. These can easily be found using the Cypher statement shown in Figure 4.10. Instead of using Cypher it would also be possible to use SQL, as shown in Figure 4.11; however this example shows again the Cypher query is more elegant.

Processing Cypher queries depends on the schema information that needs to be loaded as well. We have a supporting class for this, and the schema data used in this paper will be available via GitHub. Using the terms found with Cypher, we can construct queries that we can pass to the searcher, and create a BM25 ranking. The code that generates the rankings for all topics is presented in Figure 4.12. As you can see, with only a limited number of lines of Python code it is quite easy to

create rankings. From this point it is quite trivial to write the content of `hits` to a runfile, and evaluate using `trec_eval`.

Instead of “just” ranking with BM25, using e.g. the metadata in order to adapt the ranking is straightforward. In the case of background linking, it makes sense to consider authorship information when recommending articles that might be suitable as background reading. As journalists are often specialized in certain news topics (e.g. politics, foreign affairs, tech), the stories they write often share context. Also, when journalists collaborate on stories they write together on topics they specialize in as well. As authorship information is available to us, we can decide to use the information whether an article is written by the authors of the topic article, or by someone they have collaborated with in the past. Finding the articles that are written by this group of people can easily be done using a graph query, the query that finds these articles is shown in Figure 4.13.

Depending on the number of documents found by this query, different rescoring strategies can be decided upon. If the set of documents written by the authors or their co-authors is large, perhaps it is possible to only consider these documents, but if the set is small, a score boost might be more appropriate. Figure 4.14 shows an example on how to only consider documents found with the query in Figure 4.13, in this particular case we ensure that at least 2000 documents are found before filtering.

To give another example; the graph query language is also useful when considering entities. When journalists write news articles, the articles relate to events concerning e.g. people, organisations, or countries. In other words, the basis of news articles lay the entities as they are often the subject of news. So, instead of using the most informative terms in a news article, it could be useful to consider the entities identified in the article instead. Important entities tend to be mentioned in the beginning of a news article [28]; Figure 4.15 shows the Cypher query to retrieve the text mentions of the first five mentioned entities.

Before it is possible to search using the text describing the first five entity mentions, the text needs to be processed. The term data loaded in GeeseDB was already processed, as it was data loaded from CSV files built from a CIFF file created from an Anserini [59] (Lucene) index. Anserini has an easy to use Python extension, Pyserini [31],

```

1  from geesedb.search import Searcher
2  from geesedb.connection import get_connection
3  from geesedb.resources import
   ↪  get_topics_backgroundlinking
4  from geesedb.interpreter import Translator
5
6  db_path = '/path/to/database'
7          searcher = Searcher(
8              database=db_path,
9              n=1000
10         )
11
12  translator = Translator(db_path)
13  c_query = """cypher TFIDF query"""
14
15  query = translator.translate(c_query)
16  cursor = get_connection(db_path).cursor
17  topics = get_topics_backgroundlinking(
18      '/path/to/topics'
19  )
20  for topic_no, collection_id in topics:
21      cursor.execute(query, [collection_id])
22      topic = ' '.join(cursor.fetchall()[0])
23      hits = searcher.search_topic(topic)
24

```

Figure 4.12: Create a BM25 ranking for all background linking topics using the top-5 TFIDF terms. Note that in this case a processed topic file was used that only contains the topic identifier and the topic article id. The topic file in this format is provided on our GitHub.

```

1  MATCH (d:docs)-[]-(a1:authors)-[]-(d2:docs)-[]-(a2:authors)-
    ↪ []-(d2:docs {collection_id:
    ↪ ?})
2  RETURN DISTINCT d.collection_id

```

Figure 4.13: Cypher query to find documents written by co-authors of the authors of the topic article.

```

1  # import and first lines the same as previous example
2
3  author_c_query = """cypher authorship query"""
4  author_query = t.translate(author_c_query)
5
6  cursor = get_connection(db_path).cursor
7  topics = get_topics_backgroundlinking(
8      '/path/to/topics'
9  )
10 for topic_no, collection_id in topics:
11     cursor.execute(query, [collection_id])
12     topic = ' '.join(cursor.fetchall()[0])
13     hits = searcher.search_topic(topic)
14
15     cursor.execute(author_query, [collection_id])
16     docs_authors = {
17         e[0] for e in cursor.fetchall()
18     }
19     if len(docs_authors) > 2000:
20         hits =
            ↪ hits[hits.collection_id.isin(docs_authors)]

```

Figure 4.14: Find documents written by all authors that collaborated with the authors of the topic article, if there are more than 2000 documents found only consider these documents as background reading candidates.

```
1 MATCH (d:docs {collection_id: ?})-[]-(e:entities)
2 RETURN mention
3 ORDER BY start
4 LIMIT 5
```

Figure 4.15: Retrieve the first five entities mentioned in the topic article; and return the terms used to mention the entity.

that can be used to tokenize the text in the same way as the documents were tokenized. Figure 4.16 shows the Python code where we extract the mentions, process them such that they become a usable query for GeeseDB, and then BM25 ranking is created with this query.

4.6 Conclusion

```

1  from geesedb.search import Searcher
2  from geesedb.connection import get_connection
3  from geesedb.resources import
   ↪ get_topics_backgroundlinking
4  from geesedb.interpreter import Translator
5  from pyserini.analysis import Analyzer,
   ↪ get_lucene_analyzer
6
7  db_path = '/path/to/database'
8      searcher = Searcher(
9          database=db_path,
10         n=1000
11     )
12
13  analyzer = Analyzer(get_lucene_analyzer())
14
15  translator = Translator(db_path)
16  c_query = """cypher entity query"""
17  query = translator.translate(c_query)
18
19  cursor = get_connection(db_path).cursor
20  topics = get_topics_backgroundlinking(
21      '/path/to/topics'
22  )
23
24  for topic_no, collection_id in topics:
25      cursor.execute(query, [collection_id])
26      topic = ' '.join([e[0] for e in cursor.fetchall()])
27      topic = ' '.join(analyzer.analyze(topic))
28  hits = searcher.search_topic(topic)

```

Figure 4.16: Create a BM25 ranking for all background linking topics using the mention text of the first five linked entities in the source article.

Chapter 5

A Graph of Entities

Abstract

This is the abstract

5.1 Introduction

5.2 Entity Linking

Entity linking concerns the task of automatically identifying entity mentions in the text and linking them to the corresponding entities in a knowledge-base (KB). It fulfils a key role in knowledge-grounded understanding of text and has been proven effective for diverse tasks in information retrieval [20, 19, 58, 22, 6, 46, 8], natural language processing [34, 14], and recommendation [60]. Utilizing entity annotations in these downstream tasks depends upon the annotation of text corpora with a method for entity linking. Due to the complexity of entity linking systems, this process is often performed by a third-party entity linking toolkit, examples including DBpedia Spotlight [42], TAGME [13], Nordlys [23], GENRE [7], and REL [57].

A caveat in existing entity linking toolkits is that they have not been designed for batch processing large numbers of documents. Existing entity linking toolkits are primarily optimized to annotate individual documents, one at a time. This severely restricts utilization of state-of-the-art entity linking tools such as REL and GENRE, that employ neural approaches and require GPUs for fast operation. Annotating

millions of documents incurs significant computational overhead, to the extent that annotation of a large text corpus becomes practically infeasible using modest computational power resources. Batch entity linking is however necessary to build today’s data-hungry machine learning models, considering large text corpora like the new MS MARCO v2 (12M Web documents) [4].

5.3 REL

This paper describes our experience with optimizing the Radboud Entity Linking (REL) toolkit for batch processing large corpora. REL annotates individual documents efficiently, requiring only modest computational resources, while performing competitively when compared to the state-of-the-art methods on effectiveness. It considers entity linking as a modular problem consisting of three stages:

5.3.1 Mention Detection

The goal of this step is to identify all possible text spans in a document that might refer to an entity. If a text span that refers to an entity is not identified properly in this stage, the system will not be able to correctly link the entity in later stages.

5.3.2 Candidate Selection

For every detected mention, REL considers up to $k_1 + k_2 (= 7)$ candidate entities. $k_1 (= 4)$ candidate entities are selected based on their prior occurrence probability $p(e|m)$ (for entity e given mention m). These priors are pre-calculated from Wikipedia hyperlinks and the CrossWiki [54] corpus. The other $k_2 (= 3)$ entities are chosen based on the similarity of their embeddings to the contextual embedding of the mention (considering a context of maximum 200 word tokens).

5.3.3 Entity Disambiguation

The goal of this final step is to map the mention to the correct entity in a knowledge base. The candidate entities for each mention are

obtained from the previous stage and REL implements the Ment-norm method proposed by Le and Titov [30].

This paper explains the challenges of batch processing in REL and presents the approaches we found to overcome these challenges. We show that our updated REL toolkit, REBL, improves REL efficiency 9.5 times, decreasing the processing time per document (excluding mention detection) on a sample of 5000 MS MARCO documents from 1.23 seconds to 0.13 seconds. We demonstrate that REBL enables the annotation of a large corpus like MS MARCO v2, given modest computational resources. We discuss potential improvements that can be made in order to further improve efficiency of batch entity linking. The REBL code and toolkit are available publicly at <https://github.com/informagi/REBL>.

5.4 From REL to REBL

The objective that led to this paper was to link the MS MARCO v2 collection [4]. This collection contains 11,959,635 documents split into 60 compressed files, totaling roughly 33GB in size. Decompressed, these files are in JSON line format (where every line represents a JSON document). Documents have five fields: *url*, *title*, *headings*, *body*, and *docid*. For our experiments we wanted to link the title, headings, and body of the documents. We use the 2019-07 Wikipedia dump to link to, which is one of the two dumps REL was initially developed on. It is, however, straightforward to take another dump of Wikipedia and develop another REL instance.

In order to ease linking this size of data, we separated the GPU heavy mention detection stage from the CPU heavy candidate selection and entity disambiguation stages; the modified code can be found on GitHub.¹ The inputs for mention detection are the compressed MS MARCO v2 document files, and its output consists of the mentions found and their location in the document, in Apache Parquet format.² These files together with the source text are the input for the subsequent phases (candidate selection and entity disambiguation). The final output consists of Parquet files containing spans of text and their

¹<https://github.com/informagi/REBL>

²<https://github.com/apache/parquet-format>

linked entities. In the following, we discuss what is changed for mention detection, candidate selection, and entity disambiguation steps to make REL more suited to link the MS MARCO v2 collection.

5.4.1 Mention Detection

REL [57] uses Flair [1] for mention detection, a state-of-the-art named entity recognition system. Flair uses the `segtok`³ package to segment an (Indo-European) document in sentences, internally represented as **Sentence** objects. These sentences are split into words / symbols represented as **Token** objects. When creating these representations however, it is not possible to recreate the source text properly, as Flair removes multiple whitespace characters when occurring after each other. REL corrects for this to preserve the correct span data with regard to its location in the source text, which is an inefficient process. We set out to construct the underlying data structures ourselves for REBL. To do this, we used the `syntok`⁴ package, a follow-up version of `segtok`. The author of both packages claims that the `syntok` package segments sentences better than `segtok`.

When constructing the sentences from the token objects, we ran into another issue originated from data handling procedure in Flair: Flair removes various zero width Unicode characters from the source text: zero width space (U+200B), zero width non-joiner (U+200C), variation selector-16 (U+FE0F), and zero-width no-break space (U+FEFF). These characters occur rarely, but in a collection as big and diverse as MS MARCO v2, these characters are found in some documents. When encountering these characters, the token objects were constructed such that the span and offset of the token still referred to that of the source text.

For the case of the zero width space, we updated the `syntok` package; although zero width space is not considered a whitespace character according to the Unicode standard, it should be considered a character that separates two words. For the other Unicode characters removed by Flair, we manually update the span in the **Token** objects created by Flair such that they refer correctly to the positions in the source text. Now, when Flair identifies a series of tokens as a possible

³<https://github.com/fnl/segtok>

⁴<https://github.com/fnl/syntok>

mention, we can directly identify the location in the source text from the `Token` objects.

Flair supports named entity recognition in batches; this way multiple batches of text can be sent to the GPU for faster inference time. Because REL had been designed to tag one document at a time, it did not use this functionality. REBL exploits this feature, allowing the user to specify the number of documents to be tagged simultaneously.

5.4.2 Candidate Selection and Entity Disambiguation

REL makes use of a $p(e|m)$ prior, where e is an entity, and m is a mention. These priors are saved in a (SQLite) database, and up to 100 priors per mention are considered. Data conversion between client and the representation stored in the database incurred however a large serialization cost. We updated this to a format that is faster to load, with the additional benefit of a considerably decreased database size.⁵ We experimented with data storage in the DuckDB column oriented database as an alternative, but found that SQLite was (still) more efficient as key-value store, at least in DuckDB’s current state of development.

We found that the entity disambiguation stage took much longer than reported in the original REL paper. This difference is explained by the length of the documents to be linked. The documents evaluated by Van Hulst et al. [57] were on average 323 tokens long with an average of 42 mentions to consider. The number of tokens in an MS MARCO v2 document is on average 1800, with 84 possible mentions per document.⁶ Per mention, 100 tokens to the left, and 100 tokens to the right are considered as the context for the disambiguation model. The larger documents result in a larger memory consumption per context and per document, with higher processing costs as a consequence.

We improve the efficiency of the entity disambiguation step such that it could be run in a manageable time. REL recreates database cursors for every transaction. We rewrote the REL database code such that one database cursor is created for the entity disambiguation

⁵The table that represents the priors shrank from 9.6GB to 2.2GB.

⁶These figures are calculated over the body field; we also tagged the shorter title and headers fields.

module. Within a document, the same queries were issued to the database multiple times. This happens for example when a mention occurs multiple times within a document. By caching the output of these queries, we were able to significantly lower the number of database calls needed. We cached all database calls per every segment in the collection, as we ran the process for every segment separately.

The default setting of REL is to keep embeddings on the GPU after they are loaded. This, however, slowed down disambiguation when many documents are being processed consecutively, because operations like normalization were carried out over all embeddings on the GPU. By clearing these embeddings as soon as a document is processed, a significant speed up has been achieved.

Finally, after retrieving the embeddings from the database, REL puts them in a python list. We rewrote the REL code such that the binary data is directly loaded from NumPy, a data format that Pytorch operates on.

5.5 Effects on Execution

In the mention detection stage, we improved tokenization and applied batching. In the MS MARCO v2 collection, 411,906 documents have tokens that were automatically removed by Flair, which are 3.4% of all documents in the collection. The MS MARCO v1 collection does not have documents that contained these characters; the documents in that version of the collection are (probably) sanitized before publishing. Batching documents in the mention detection stage decreased the average time for finding all named entities. We used batches of size 10, as the documents are relatively large. The optimal batch size will depend on the available GPU memory.

A few documents in the MS MARCO v2 collection could not be linked. This happened only in extraordinary cases, where linking with entities did not make sense in the first place; an example being a document consisting of numbers only.⁷ Here, the `syntok` package created one long `Sentence` object from this file that could not fit in GPU memory.

⁷The source document was a price list in PDF format.

Table 5.1 shows the improvements we made to the candidate selection and entity disambiguation step, and describes how much time is saved in REBL. The code improvements to create the database cursor only once and to load the data directly from NumPy had no noticeable effect on the overall run time of entity disambiguation and are not reported in this table. Note that the large standard deviations are primarily due to the differences in processing costs between long and short documents.

5.6 MMEAD

5.7 Results

5.8 Conclusion

We introduced REBL, an extension for the Radboud Entity Linker. We utilize REL’s modular design to separate the GPU heavy mention detection stage from the CPU heavy candidate selection and entity disambiguation stages, as many researchers have dedicated GPU and CPU machines. The mention detection module has been made more robust and reliable, using a better segmenter and preserving location metadata correctly. The candidate selection step and entity disambiguation step were updated to improve their runtime, especially for longer documents.

Although it is now possible to run REL [57] on MS MARCO v2 [4] in a (for us) somewhat reasonable time, we identified further improvements to implement, that we work on actively.

Found mentions are compared to all other mentions during the candidate selection step, the complexity of this step is $O(n^2)$, with n being the number of mentions found in a document, which is especially problematic for longer documents. As we are only interested in mentions that are similar, we expect that it might be worthwhile to implement a locality sensitive hashing algorithm to decrease the number of comparisons needed in this stage. However, we would need to run additional experiments to ensure the effectiveness of the model does not suffer.

REBL now implements a two step approach that writes intermediate

Table 5.1: Efficiency improvements for Candidate Selection and Entity Disambiguation. Improvements are calculated over a sample of 5000 documents using a machine with an Intel Xeon Silver 4214 CPU @ 2.20GHz using 2 cores, that has 187GB RAM memory, and a GeForce RTX 2080 Ti (11GB) GPU. Improvements are cumulative; the times shown include the previous improvement as well.

Improvement	Seconds	Explanation
Old Candidate Selection + Entity Disambiguation	1.23 ± 2.09	Average time it takes to select candidates and disambiguate per document
No embedding reset	0.26 ± 1.60	The default setting of REL was to keep embeddings in GPU memory after they were loaded, by clearing them from GPU memory after every document a speed up was achieved.
Cache database calls	0.15 ± 1.31	When an entity occurs within a document, there is a high probability of it occurring multiple times. By caching the calls, we increase the memory usage but are able to lower the time needed for candidate selection + entity disambiguation.
Representation change candidates	0.13 ± 1.19	By representing the candidates better in the database, we were able to save on conversion time lowering the time needed for candidate selection.

results to the file system in Parquet format. A streaming variant would be preferable. We have also kept SQLite as database backend, but will consider specialized key-value stores to speed up candidate selection and entity disambiguation. We will revisit DuckDB upon progress in the implementation of zero-cost positional joins.

The candidate selection stage considers the context of a mention. This context has to be constructed from the source document. As a result, we load the source data a second time during candidate selection. Alternatively, we may output mention context in the mention detection stage, which could then speed up the remaining. However, this would significantly increase the size of the mention detection output. More experiments are needed to strike the right balance here.

Overall, it has become clear that a data processing oriented perspective on entity linking is necessary for efficient solutions. Having made explicit quite a few implicit design choices, re-evaluating these might lead to more effective entity linking as well.

Chapter 6

Conclusion

We can only see a short distance ahead, but we can see plenty there that needs to be done.

Alan Turing - 1950

Bibliography

- [1] AKBİK, A., BERGMANN, T., BLYTHE, D., RASUL, K., SCHWETER, S., AND VOLLGRAF, R. FLAIR: An easy-to-use framework for state-of-the-art NLP. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)* (2019), pp. 54–59.
- [2] ANGLES, R. The Property Graph Database Model. In *Proceedings of the 12th Alberto Mendelzon International Workshop on Foundations of Data Management* (Aachen, 2018), AMW '18, CEUR-WS.org.
- [3] ARGUELLO, J., CRANE, M., DIAZ, F., LIN, J., AND TROTMAN, A. Report on the SIGIR 2015 Workshop on Reproducibility, Inexplicability, and Generalizability of Results (RIGOR). *SIGIR Forum* 49, 2 (jan 2016), 107–116.
- [4] BAJAJ, P., CAMPOS, D., CRASWELL, N., DENG, L., GAO, J., LIU, X., MAJUMDER, R., ANDREW MCNAMARA, B. M., NGUYEN, T., ROSENBERG, M., SONG, X., STOICA, A., TIWARY, S., AND WANG, T. MS MARCO: A Human Generated MACHine Reading COMprehension Dataset. In *InCoCo@NIPS* (2016).
- [5] BALOG, K. *Entity-oriented search*. Springer Nature, Gewerbestrasse 11, 6330 Cham, Switzerland, 2018.
- [6] BALOG, K., RAMAMPIARO, H., TAKHIROV, N., AND NØRVÅG, K. Multi-Step Classification Approaches to Cumulative Citation Recommendation. In *Proceedings of the 10th Conference on Open Research Areas in Information Retrieval* (2013), OAIR '13, p. 121–128.

- [7] CAO, N. D., IZACARD, G., RIEDEL, S., AND PETRONI, F. Autoregressive Entity Retrieval. In *International Conference on Learning Representations* (2021).
- [8] CHATTERJEE, S., AND DIETZ, L. BERT-ER: Query-specific BERT Entity Representations for Entity Ranking. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2022), SIGIR '22, p. 1466–1477.
- [9] CHURCH, K. W., AND GALE, W. A. Poisson mixtures. *Natural Language Engineering* 1, 2 (1995), 163–190.
- [10] CLANCY, R., FERRO, N., HAUFF, C., LIN, J., SAKAI, T., AND WU, Z. Z. The SIGIR 2019 Open-Source IR Replicability Challenge (OSIRRC 2019). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2019), SIGIR'19, Association for Computing Machinery, p. 1432–1434.
- [11] DALTON, J., DIETZ, L., AND ALLAN, J. Entity Query Feature Expansion Using Knowledge Base Links. In *Proceedings of the 37th International ACM SIGIR Conference on Research; Development in Information Retrieval* (New York, NY, USA, 2014), SIGIR '14, Association for Computing Machinery, p. 365–374.
- [12] DEVEAUD, R., ALBAKOUR, M.-D., MACDONALD, C., AND OUNIS, I. On the Importance of Venue-Dependent Features for Learning to Rank Contextual Suggestions. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management* (New York, NY, USA, 2014), CIKM '14, Association for Computing Machinery, p. 1827–1830.
- [13] FERRAGINA, P., AND SCAIELLA, U. TAGME: On-the-Fly Annotation of Short Text Fragments (by Wikipedia Entities). In *Proceedings of the 19th ACM International Conference on Information and Knowledge Management* (2010), CIKM '10, p. 1625–1628.
- [14] FERRUCCI, D. Introduction to “This is Watson”. *IBM Journal of Research and Development* 56 (05 2012), 1:1–1:15.

- [15] FRANCIS, N., GREEN, A., GUAGLIARDO, P., LIBKIN, L., LINDAAKER, T., MARSAULT, V., PLANTIKOW, S., RYDBERG, M., SELMER, P., AND TAYLOR, A. Cypher: An Evolving Query Language for Property Graphs. In *Proceedings of the 2018 International Conference on Management of Data* (New York, NY, USA, 2018), SIGMOD '18, Association for Computing Machinery, p. 1433–1445.
- [16] FUHR, N. Models for Integrated Information Retrieval and Database Systems. *IEEE Data Engineering Bulletin* 19, 1 (1996), 3–13.
- [17] FUHR, N., AND RÖLLEKE, T. A Probabilistic Relational Algebra for the Integration of Information Retrieval and Database Systems. *ACM Trans. Inf. Syst.* 15, 1 (jan 1997), 32–66.
- [18] GAO, L., DAI, Z., CHEN, T., FAN, Z., VAN DURME, B., AND CALLAN, J. Complement Lexical Retrieval Model with Semantic Residual Embeddings. In *Advances in Information Retrieval* (Cham, 2021), ECIR '21, Springer International Publishing, pp. 146–160.
- [19] GERRITSE, E. J., HASIBI, F., AND DE VRIES, A. P. Graph-Embedding Empowered Entity Retrieval. In *Proceedings of the 42nd European Conference on Information Retrieval* (2020), pp. 97–110.
- [20] GERRITSE, E. J., HASIBI, F., AND DE VRIES, A. P. Entity-Aware Transformers for Entity Search. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2022), SIGIR '22, pp. 1455–1465.
- [21] HASIBI, F., BALOG, K., AND BRATSBERG, S. E. Exploiting Entity Linking in Queries for Entity Retrieval. In *Proceedings of the 2016 ACM International Conference on the Theory of Information Retrieval* (New York, NY, USA, 2016), ICTIR '16, Association for Computing Machinery, p. 209–218.
- [22] HASIBI, F., BALOG, K., AND BRATSBERG, S. E. Exploiting Entity Linking in Queries for Entity Retrieval. In *Proceedings*

- of the 2016 ACM International Conference on the Theory of Information Retrieval* (2016), ICTIR '16, p. 209–218.
- [23] HASIBI, F., BALOG, K., GARIGLIOTTI, D., AND ZHANG, S. Nordlys: A Toolkit for Entity-Oriented and Semantic Search. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2017), SIGIR '17, p. 1289–1292.
 - [24] KAMPHUIS, C., AND DE VRIES, A. P. Reproducible IR needs an (IR) (graph) query language. In *Proceedings of the Open-Source IR Replicability Challenge co-located with 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, OSIRRC@SIGIR 2019, Paris, France, July 25, 2019* (Aachen, 2019), CEUR-WS.org, pp. 17–20.
 - [25] KAMPHUIS, C., AND DE VRIES, A. P. The OldDog Docker Image for OSIRRC at SIGIR 2019. In *Proceedings of the Open-Source IR Replicability Challenge co-located with 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, OSIRRC@SIGIR 2019, Paris, France, July 25, 2019* (Aachen, 2019), CEUR-WS.org, pp. 47–49.
 - [26] KAMPHUIS, C., AND DE VRIES, A. P. GeeseDB: A Python Graph Engine for Exploration and Search. In *Proceedings of the 2nd International Conference on Design of Experimental Search & Information REtrieval Systems* (Aachen, 2021), DESIRES '21, CEUR-WS.org, pp. 10–18.
 - [27] KAMPHUIS, C., DE VRIES, A. P., BOYTSOV, L., AND LIN, J. Which BM25 Do You Mean? A Large-Scale Reproducibility Study of Scoring Variants. In *Advances in Information Retrieval* (Cham, 2020), J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins, Eds., ECIR '20, Springer International Publishing, pp. 28–34.
 - [28] KAMPHUIS, C., HASIBI, F., DE VRIES, A. P., AND CRIJNS, T. Radboud University at TREC 2019. In *Proceedings of The Twenty-Eight Text REtrieval Conference* (Gaithersburg, Maryland, USA, 2019), TREC '19, National Institute for Standards and Technology (NIST).

- [29] KAMPHUIS, C., HASIBI, F., LIN, J., AND DE VRIES, A. P. REBL: Entity Linking at Scale. In *Proceedings of the 3rd International Conference on Design of Experimental Search & Information REtrieval Systems* (Aachen, 2022), DESIRES '22, CEUR-WS.org.
- [30] LE, P., AND TITOV, I. Improving Entity Linking by Modeling Latent Relations between Mentions. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (July 2018), pp. 1595–1604.
- [31] LIN, J., MA, X., LIN, S.-C., YANG, J.-H., PRADEEP, R., AND NOGUEIRA, R. Pyserini: An easy-to-use Python toolkit to support replicable IR research with sparse and dense representations. *arXiv preprint arXiv:2102.10073* (2021).
- [32] LIN, J., MACKENZIE, J., KAMPHUIS, C., MACDONALD, C., MALLIA, A., SIEDLACZEK, M., TROTMAN, A., AND DE VRIES, A. Supporting Interoperability Between Open-Source Search Engines with the Common Index File Format. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2020), SIGIR '20, Association for Computing Machinery, p. 2149–2152.
- [33] LIN, S., YANG, J., AND LIN, J. Distilling Dense Representations for Ranking using Tightly-Coupled Teachers. *CoRR abs/2010.11386* (2020).
- [34] LIN, T., MAUSAM, AND ETZIONI, O. Entity Linking at Web Scale. In *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX)* (June 2012), pp. 84–88.
- [35] LUAN, Y., EISENSTEIN, J., TOUTANOVA, K., AND COLLINS, M. Sparse, Dense, and Attentional Representations for Text Retrieval. *Transactions of the Association for Computational Linguistics 9* (2021), 329–345.
- [36] LV, Y., AND ZHAI, C. Adaptive Term Frequency Normalization for BM25. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management* (New York, NY,

- USA, 2011), CIKM '11, Association for Computing Machinery, p. 1985–1988.
- [37] LV, Y., AND ZHAI, C. Lower-Bounding Term Frequency Normalization. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management* (New York, NY, USA, 2011), CIKM '11, Association for Computing Machinery, p. 7–16.
- [38] LV, Y., AND ZHAI, C. When Documents Are Very Long, BM25 Fails! In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2011), SIGIR '11, Association for Computing Machinery, p. 1103–1104.
- [39] MACDONALD, C., SANTOS, R. L., AND OUNIS, I. On the Usefulness of Query Features for Learning to Rank. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management* (New York, NY, USA, 2012), CIKM '12, Association for Computing Machinery, p. 2559–2562.
- [40] MACDONALD, C., TONELLOTO, N., MACAVANEY, S., AND OUNIS, I. PyTerrier: Declarative Experimentation in Python from BM25 to Dense Retrieval. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management* (New York, NY, USA, 2021), CIKM '21, Association for Computing Machinery, p. 4526–4533.
- [41] MACLEOD, I. A. Text Retrieval and the Relational Model. *Journal of the American Society for Information Science* 42, 3 (1991), 155–165.
- [42] MENDES, P. N., JAKOB, M., GARCÍA-SILVA, A., AND BIZER, C. DBpedia Spotlight: Shedding Light on the Web of Documents. In *Proceedings of the 7th International Conference on Semantic Systems* (2011), I-Semantics '11, p. 1–8.
- [43] MÜHLEISEN, H., SAMAR, T., LIN, J., AND DE VRIES, A. Old Dogs Are Great at New Tricks: Column Stores for IR Prototyping. In *Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval* (New York,

- NY, USA, 2014), SIGIR '14, Association for Computing Machinery, p. 863–866.
- [44] OUNIS, I., AMATI, G., PLACHOURAS, V., HE, B., MACDONALD, C., AND JOHNSON, D. Terrier Information Retrieval Platform. In *Advances in Information Retrieval* (Berlin, Heidelberg, 2005), D. E. Losada and J. M. Fernández-Luna, Eds., Springer Berlin Heidelberg, pp. 517–519.
- [45] RAASVELDT, M., AND MÜHLEISEN, H. DuckDB: An Embeddable Analytical Database. In *Proceedings of the 2019 International Conference on Management of Data* (New York, NY, USA, 2019), SIGMOD '19, Association for Computing Machinery, p. 1981–1984.
- [46] REINANDA, R., MELJ, E., AND DE RIJKE, M. Mining, Ranking and Recommending Entity Aspects. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2015), SIGIR '15, p. 263–272.
- [47] ROBERTSON, S., AND ZARAGOZA, H. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [48] ROBERTSON, S. E., AND JONES, K. S. Relevance weighting of search terms. *Journal of the American Society for Information Science* 27, 3 (1976), 129–146.
- [49] ROBERTSON, S. E., WALKER, S., JONES, S., HANCOCK-BEAULIEU, M., AND GATFORD, M. Okapi at TREC-3. In *TREC* (1994).
- [50] ROUSSEAU, F., AND VAZIRGIANNIS, M. Composition of TF Normalizations: New Insights on Scoring Functions for Ad Hoc IR. In *Proceedings of the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2013), SIGIR '13, Association for Computing Machinery, p. 917–920.
- [51] SCHEK, H.-J., AND PISTOR, P. Data Structures for an Integrated Data Base Management and Information Retrieval System. In

- Proceedings of the 8th International Conference on Very Large Data Bases* (San Francisco, CA, USA, 1982), VLDB '82, Morgan Kaufmann Publishers Inc., p. 197–207.
- [52] SINGHAL, A., BUCKLEY, C., AND MITRA, M. Pivoted Document Length Normalization. In *Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 1996), SIGIR '96, Association for Computing Machinery, p. 21–29.
- [53] SOBOROFF, I., HUANG, S., AND HARMAN, D. TREC 2018 News Track Overview. In *Proceedings of The Twenty-Seventh Text REtrieval Conference* (Gaithersburg, Maryland, USA, 2018), TREC '18, National Institute for Standards and Technology (NIST).
- [54] SPITKOVSKY, V. I., AND CHANG, A. X. A Cross-Lingual Dictionary for English Wikipedia Concepts. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)* (2012), European Language Resources Association (ELRA), pp. 3168–3175.
- [55] TROTMAN, A., JIA, X., AND CRANE, M. Towards an Efficient and Effective Search Engine. In *OSIR@ SIGIR* (2012), pp. 40–47.
- [56] TROTMAN, A., PUURULA, A., AND BURGESS, B. Improvements to BM25 and Language Models Examined. In *Proceedings of the 2014 Australasian Document Computing Symposium* (New York, NY, USA, 2014), ADCS '14, Association for Computing Machinery, p. 58–65.
- [57] VAN HULST, J. M., HASIBI, F., DERCKSEN, K., BALOG, K., AND DE VRIES, A. P. REL: An Entity Linker Standing on the Shoulders of Giants. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2020), SIGIR '20, Association for Computing Machinery, p. 2197–2200.
- [58] XIONG, C., CALLAN, J., AND LIU, T.-Y. Word-Entity Duet Representations for Document Ranking. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2017), SIGIR '17, p. 763–772.

- [59] YANG, P., FANG, H., AND LIN, J. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval* (New York, NY, USA, 2017), SIGIR '17, Association for Computing Machinery, p. 1253–1256.
- [60] YANG, Y., IRSOY, O., AND RAHMAN, K. S. Collective Entity Disambiguation with Structured Gradient Tree Boosting. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)* (2018), pp. 777–786.

Summary

Samenvatting

Acknowledgements

Research Data Management

Curriculum Vitae