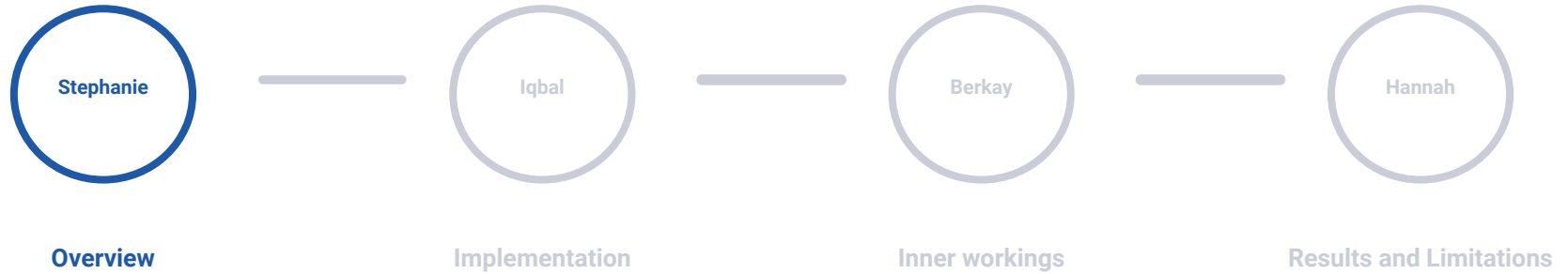# Search Engine Design with BM25 and DESM

## ECS735P Information Retrieval - Group 18

Stephanie Nicole Garibay Lim
Berkay Dur
Hannah Melkemaryam Claus
Iqbal Singh

Queen Mary
University of London

# Presentation Structure

| Stephanie | Iqbal | Berkay | Hannah |
|:---:|:---:|:---:|:---:|
| **Overview** | Implementation | Inner workings | Results and Limitations |

# Overview

- Importance of identification of the most relevant and accurate news information
  - Incapability to do so results to
    - Inaccurate or incomplete news information
    - Negative consequences such as misinformation and misinterpretation

- An IR system is capable to:
  - Accurate
  - Reliable
  - Timely

- Project scope:
  - Investigate and evaluate the design of news article search engine with 2 different models
  - Overview of the future design implementation of the system.

(Saracevic, 2010)

# Problem statement

## Dataset

- BBC News articles, scraped using Python.
- Contains all the recent news articles within the 20-day window, from the day it was scraped.
- Total data used for this project is **8,333**

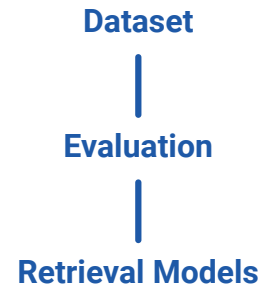| | id bigint | url text | content text | header text |
|---|---|---|---|---|
| 1 | 6965 | https://www.bbc.com/news/world-europe-17028059 | andorra media guide andorran media scene partly shaped proximity france sp... | Andorra media guide |
| 2 | 6966 | https://www.bbc.com/news/world-latin-america-20271246 | martinique media guide tv radio services provided french public overseas bro... | Martinique media guide |
| 3 | 6967 | https://www.bbc.com/news/world-latin-america-20219640 | cayman islands media guide four tv stations air caymans two run religious or... | Cayman Islands media guide |
| 4 | 6968 | https://www.bbc.com/news/world-africa-20274845 | guadeloupe profile facts known carib indian population karukera island beauti... | Guadeloupe profile - Facts |
| 5 | 6979 | https://www.bbc.com/news/world-latin-america-20413716 | french guiana media guide commercial broadcasters operate alongside servi... | French Guiana media guide |
| 6 | 6969 | https://www.bbc.com/news/world-africa-20274424 | guadeloupe media guide commercial broadcasters operate alongside service... | Guadeloupe media guide |
| 7 | 6970 | https://www.bbc.com/news/world-europe-17219246 | cyprus media guide cypriot media mirror island political division zone north o... | Cyprus media guide |
| 8 | 6971 | https://www.bbc.com/news/world-europe-18023383 | nato finland joining joining nato finland ending seven decades country finland... | What is Nato and why is Finland joinin... |
| 9 | 6972 | https://www.bbc.com/news/world-europe-17205118 | bulgaria media profile television internet media main sources information prin... | Bulgaria media profile |
| 10 | 6973 | https://www.bbc.com/sport/american-football | american football super bowl winner tom brady agrees become wnba champi... | American Football |
| 11 | 6974 | https://www.bbc.com/news/world-latin-america-19596910 | grenada media guide grenada free media guaranteed law country daily newsp... | Grenada media guide |
| 12 | 6988 | https://www.bbc.com/news/world-africa-14094381 | sierra leone media guide media freedom sierra leone limits media rights moni... | Sierra Leone media guide |
| 13 | 6989 | https://www.bbc.com/news/uk-politics-40031087 | terrorism threat levels work terrorism threat level northern ireland raised uk te... | How do terrorism threat levels work? |
| 14 | 6975 | https://www.bbc.com/news/world-latin-america-18425060 | falkland islands media guide coverage local affairs provided radio station terr... | Falkland Islands media guide |
| 15 | 6976 | https://www.bbc.com/news/world-europe-18249814 | greenland media guide kalaallit nunaata radio knr greenland broadcasting co... | Greenland media guide |
| 16 | 6977 | https://www.bbc.com/news/world-us-canada-17140680 | puerto rico media guide broadcasting regulated us federal communications c... | Puerto Rico media guide |

# Inverted Index

# Problem statement

**Evaluation**

- Limitations with dataset chosen:
  - Unlabelled data
- Manual scoring
  - Query and result relevance based on:
    - News category
    - News article content

# Problem statement

**Retrieval Models**

- BM25
- Dual Embedding Space Model (DESM)

# Retrieval Models

- **BM25**
  - It ranks the relevance of the document by weighing the similarity of the query terms in the document
    - Counting repetition of query terms in the document
  - Assumption:
    - Query terms are more useful for document ranking
  - Uses only original query terms and any additional query will be linked to the original query via relevance

- **DESM**
  - Two embeddings
    - Query words, Q,
    - Document words, D.
  - Ranking function is simply the mean cosine similarity of Q and D bar
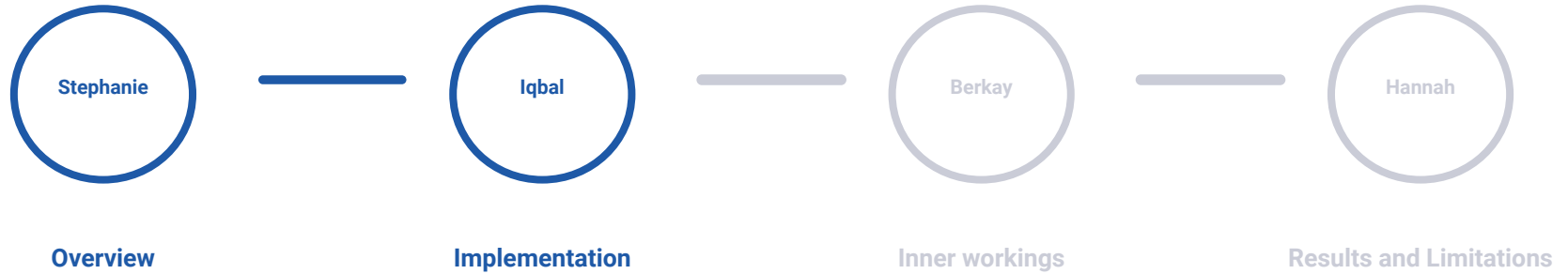  - Takes into account distributional semantics, which incorporates the relationship between words

$$DESM(Q, D) = \frac{1}{|Q|} \sum_{q_i \in Q} \frac{\mathbf{q}_i^T \overline{\mathbf{D}}}{\|\mathbf{q}_i\| \|\overline{\mathbf{D}}\|}$$
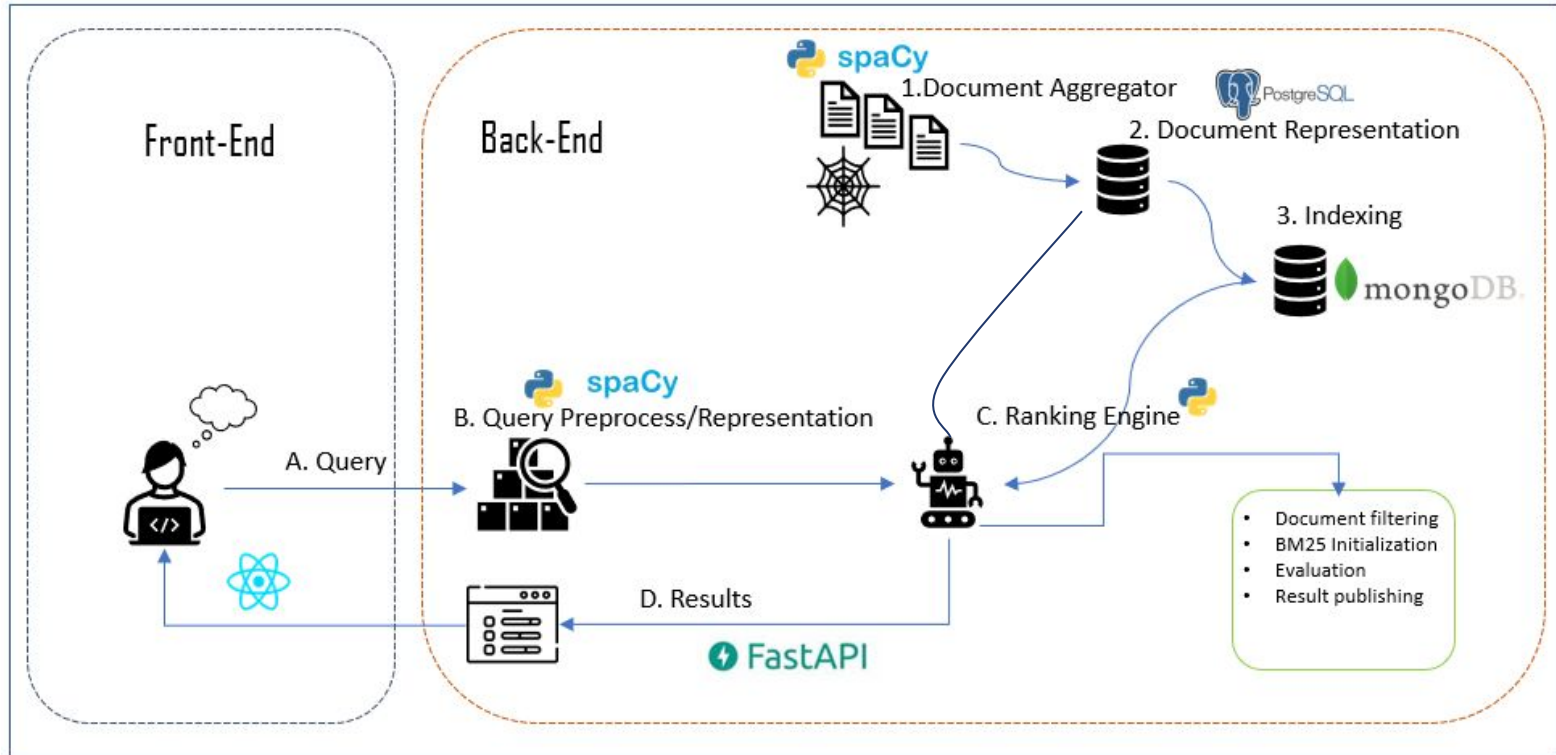
Where:

$$\overline{\mathbf{D}} = \frac{1}{|D|} \sum_{\mathbf{d}_j \in D} \frac{\mathbf{d}_j}{\|\mathbf{d}_j\|}$$

(Nalisnick *et al.,* 2016)

Queen Mary
University of London

# Presentation Structure



| Stephanie | Iqbal | Berkay | Hannah |
|-----------|-------|--------|--------|
| **Overview** | **Implementation** | Inner workings | Results and Limitations |

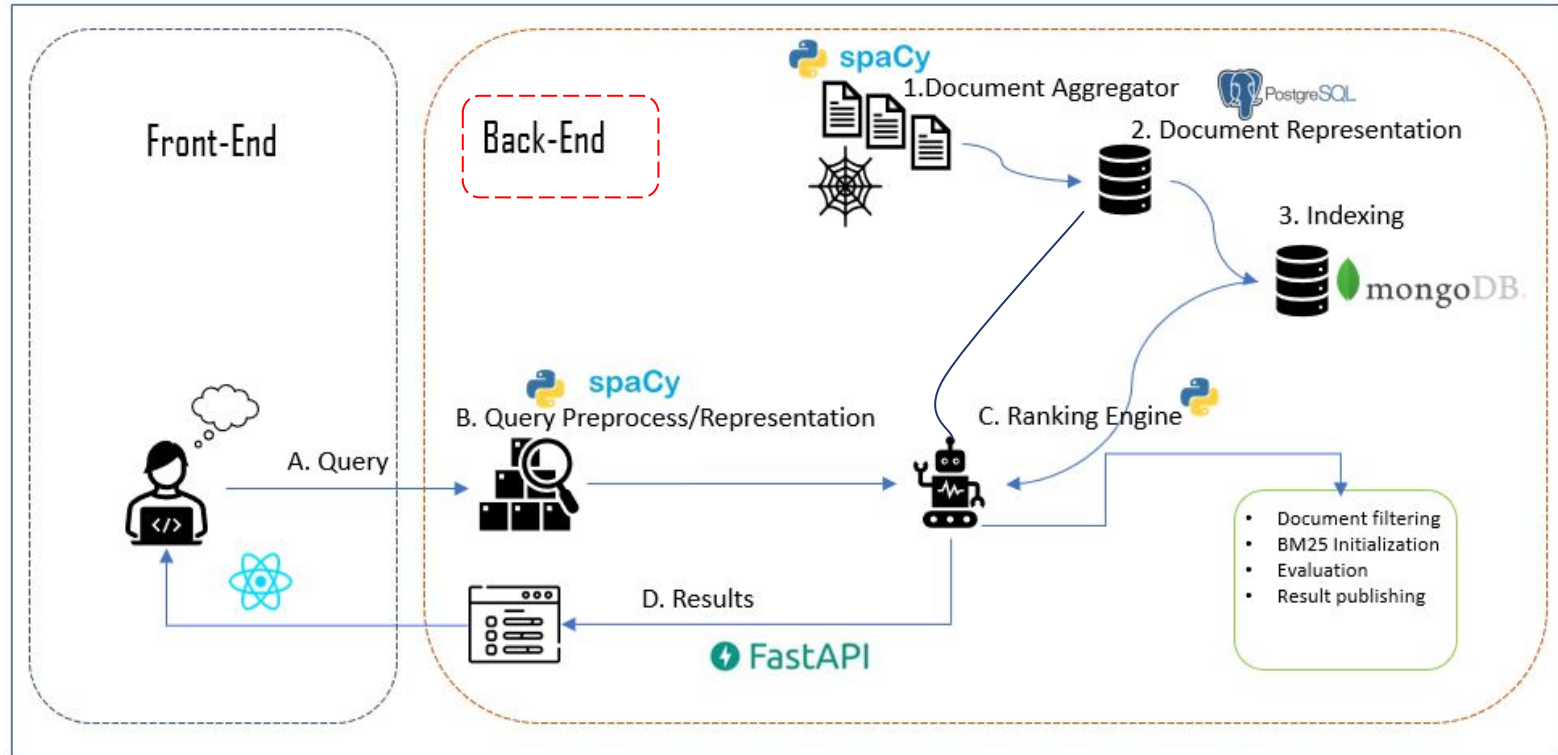# Implementation
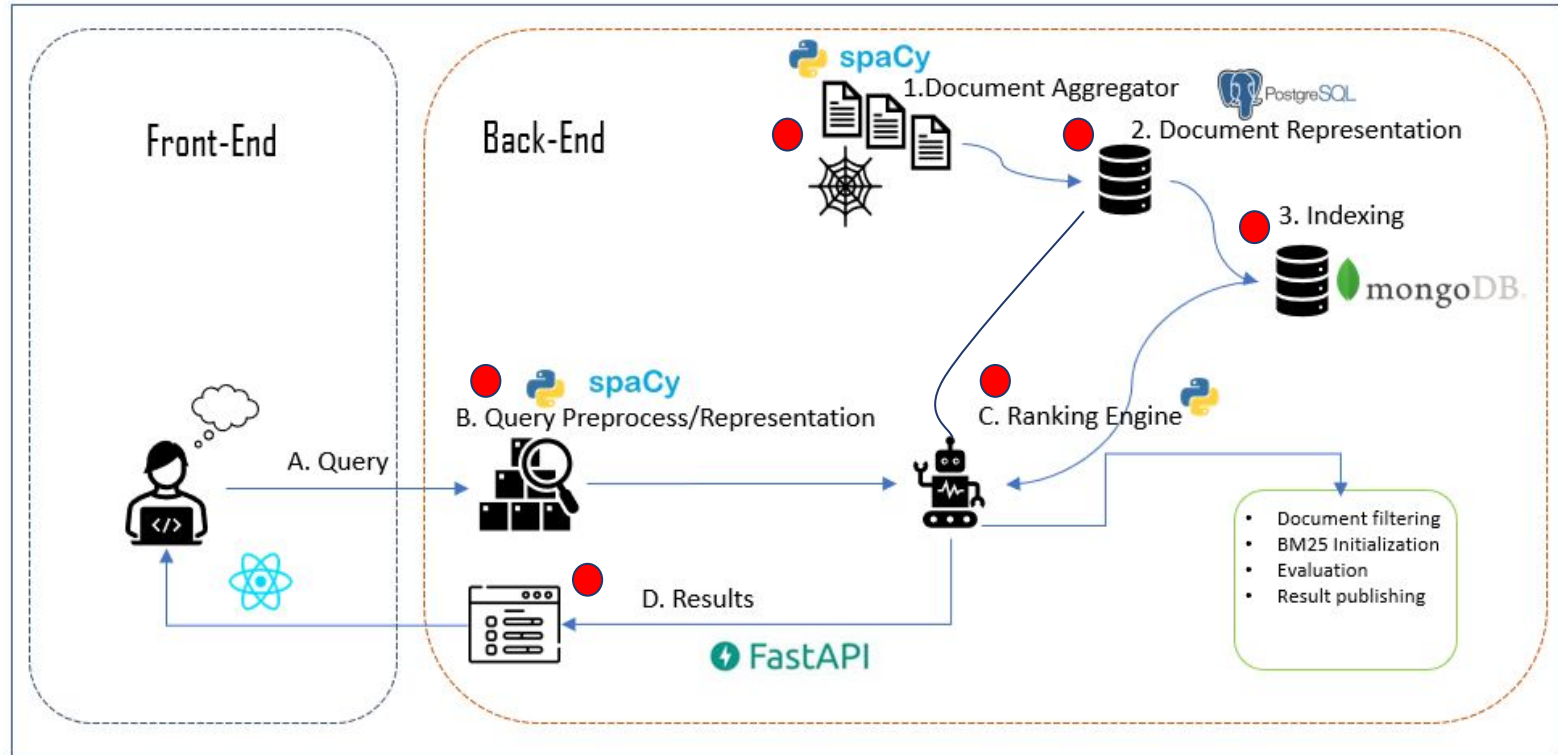
# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Implementation

# Presentation Structure

Stephanie — Iqbal — Berkay — Hannah

**Overview**　　　**Implementation**　　　**Inner workings**　　　Results and Limitations

# Inner Workings

**Frontend**

"Is Rishi Sunak Prime Minister?"

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

"Is Rishi Sunak Prime Minister?"

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
  "Is",
  "Rishi",
  "Sunak",
  "Prime",
  "Minister",
  "?"
]

# Inner Workings

**Frontend**          **Backend**

**GET**

"Is Rishi Sunak
Prime
Minister?"

**Preprocessing
Engine**
[
    "Is",
    "Rishi",
    "Sunak",
    "Prime",
    "Minister",
    "?"
]

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "Is",
    "Rishi",
    "Sunak",
    "Prime",
    "Minister"
]

# Inner Workings

**Frontend**  |  **Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

[
    "Is",
    "Rishi",
    "Sunak",
    "Prime",
    "Minister"
]

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
   "Is",
   "Rishi",
   "Sunak",
   "Prime",
   "Minister"
]

Queen Mary
University of London

# Inner Workings

**Frontend**  |  **Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

[
    "is",
    "rishi",
    "sunak",
    "prime",
    "minister"
]

# Inner Workings

**Frontend**                    **Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**<u>Preprocessing Engine</u>**

[
    "is",
    "rishi",
    "sunak",
    "prime",
    "minister"
]

# Inner Workings

**Frontend**    **Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "is",
    "rishi",
    "sunak",
    "prime",
    "minister"
]

Queen Mary
University of London

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[

   "rishi",
   "sunak",
   "prime",
   "minister"

]

Queen Mary
University of London

# Inner Workings

**Frontend**                    **Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

Queen Mary
University of London

# Inner Workings

**Frontend**

**Backend**

**Ranking Engine**

BM25

DESM

**Preprocessing Engine**

[
   "rishi",
   "sunak",
   "prime",
   "minister"
]

**GET**

"Is Rishi Sunak Prime Minister?"

Queen Mary
University of London

# Inner Workings

**Frontend**

**Backend**

**Ranking Engine**

DESM

**Preprocessing Engine**

[

    "rishi",

    "sunak",

    "prime",

    "minister"

]

**GET**

"Is Rishi Sunak Prime Minister?"

# Inner Workings

**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

### Preprocessing Engine

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

### Ranking Engine

DESM
```
{
 "arsenal" : [4, 22, 36, …],
 "first" : [5, 6, 8, 35, …],
 "rishi" : [1, 2, 5, 23, ...],
 "justice" : [2, 4, 6, 33, …],
 "sunak" : [1, 5, 7, 43, ...],
 "prime" : [2, 7, 23, 33, ...],
 "tennis" : [64, 79, 81, …],
 "fired" : [33, 64, 99, …],
 "ukraine" : [33, 98, …],
 "train" : [14, 85, 108, …],
 "minister" : [2, 5, 7, ...],
 …
}
```

# Inner Workings

**Frontend**

**Backend**

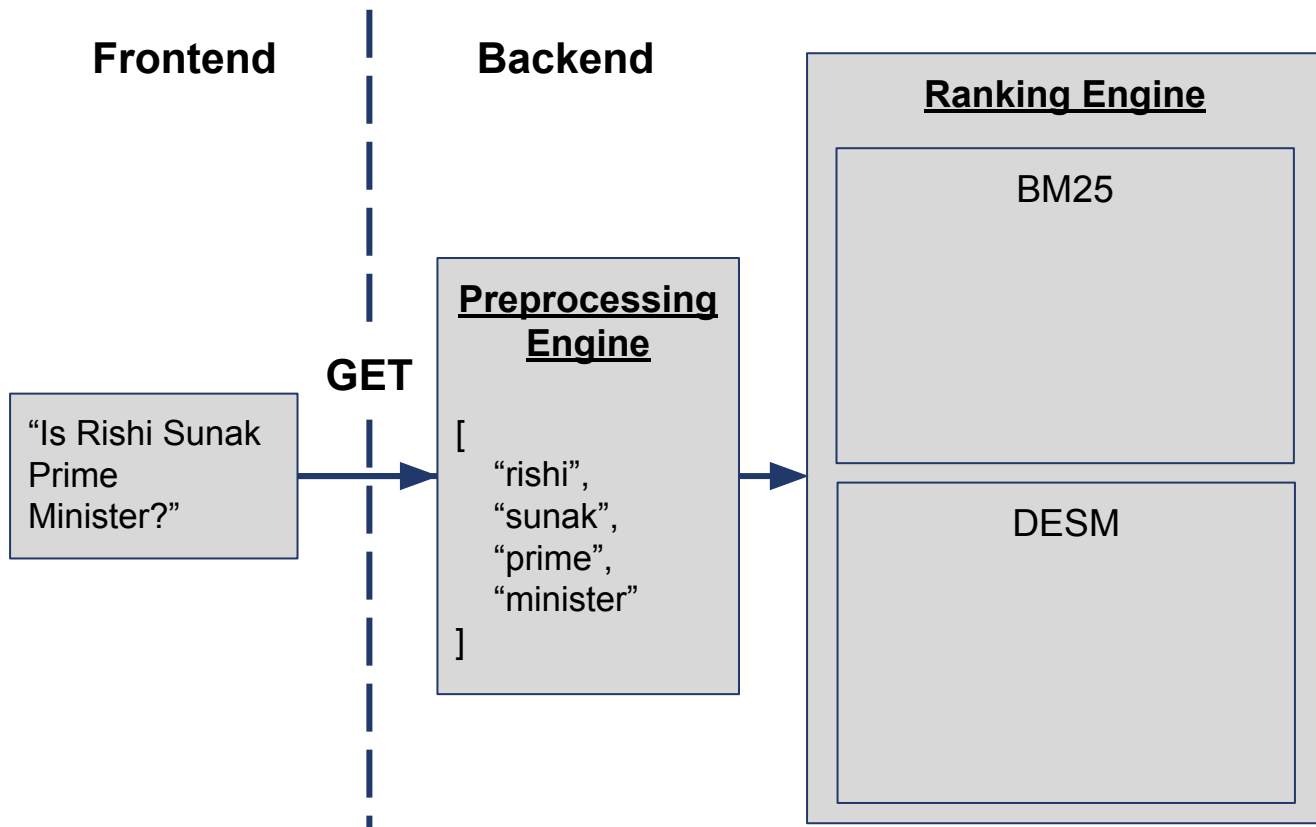**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM

{
 "arsenal" : [4, 22, 36, …],
 "first" : [5, 6, 8, 35, …],
 "rishi" : [1, 2, 5, 23, ...],
 "justice" : [2, 4, 6, 33, …],
 "sunak" : [1, 5, 7, 43, ...],
 "prime" : [2, 7, 23, 33, ...],
 "tennis" : [64, 79, 81, …],
 "fired" : [33, 64, 99, …],
 "ukraine" : [33, 98, …],
 "train" : [14, 85, 108, …],
 "minister" : [2, 5, 7, ...],
 …
}

Queen Mary
University of London

# Inner Workings

**Frontend**

**Backend**

**Ranking Engine**

DESM

```
{
 "rishi" : [1, 2, 5, 23, ...],
 "sunak" : [1, 5, 7, 43, ...],
 "prime" : [2, 7, 23, 33, ...],
 "minister" : [2, 5, 7, ...]
}
```

**GET**

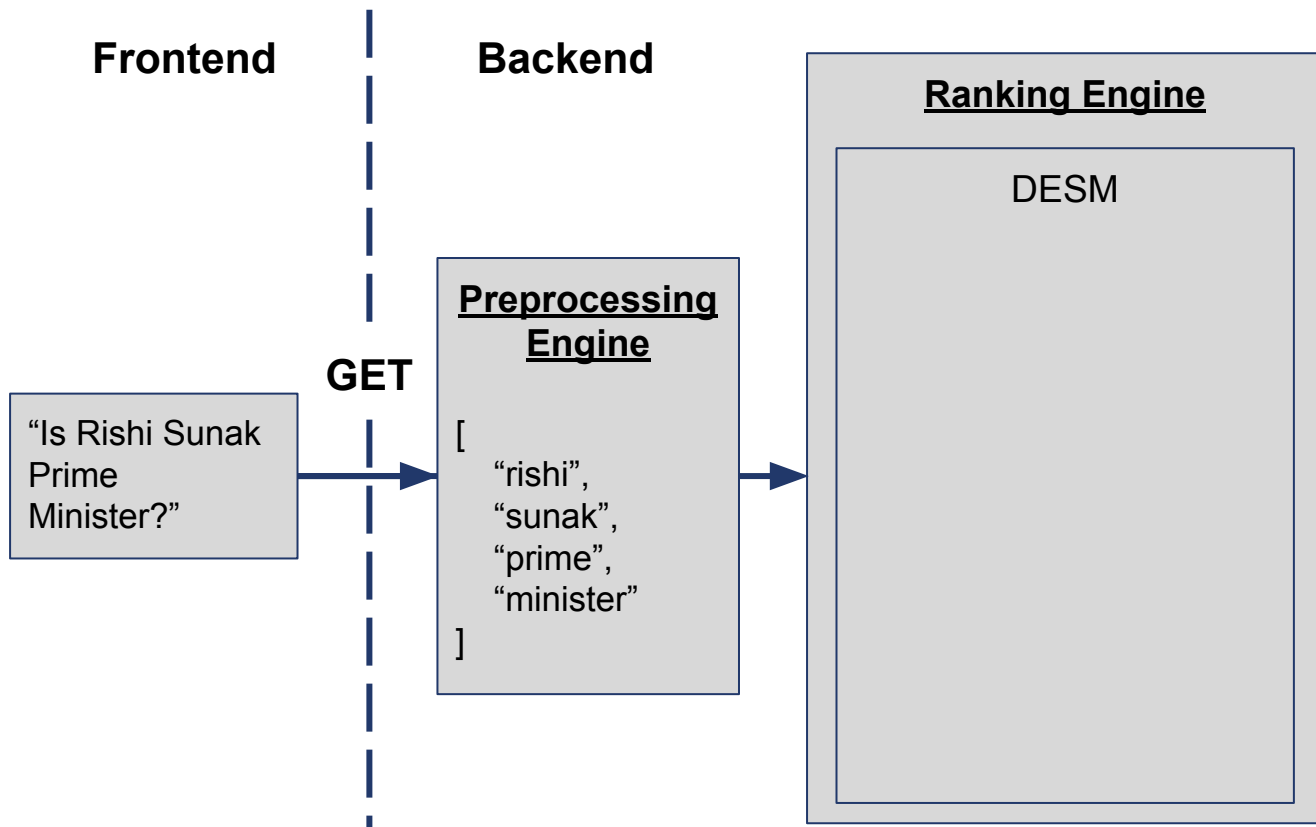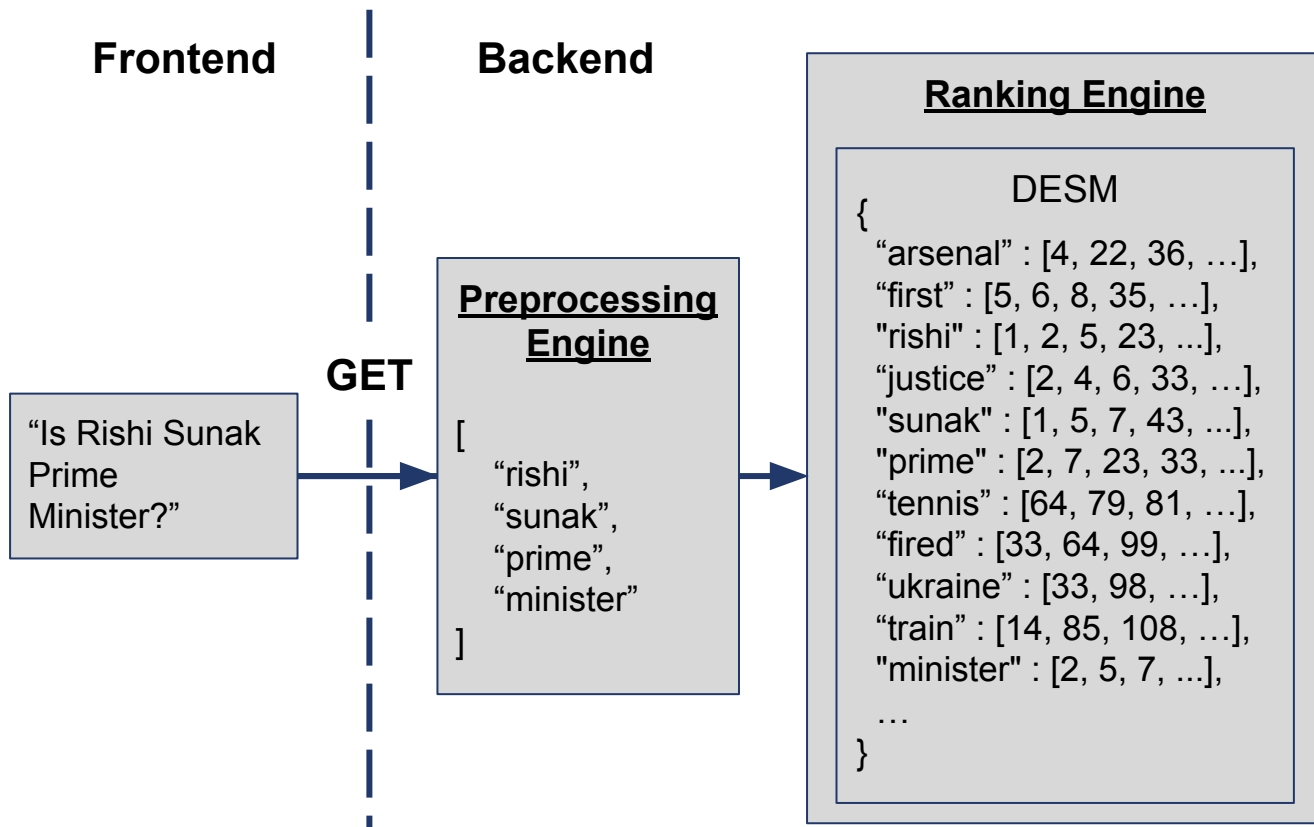**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

"Is Rishi Sunak Prime Minister?"

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM
```
{
  "rishi" : [1, 2, 5, 23, ...],
  "sunak" : [1, 5, 7, 43, ...],
  "prime" : [2, 7, 23, 33, ...],
  "minister" : [2, 5, 7, ...]
}
```

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM
{
 "rishi" : [1, 2, 5, 23, ...],
 "sunak" : [1, 5, 7, 43, ...],
 "prime" : [2, 7, 23, 33, ...],
 "minister" : [2, 5, 7, ...]
}

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM
```
{
  doc_ids : [
    1, 2, 5, 7, 23, 33, 43, …
  ]
}
```

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM
{
 doc_ids : [
    1, 2, 5, 7, 23, 33, 43, …
  ]
}

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM

Document Centroids: {
  1 : [0.002, 0.23, 0.03, ...],
  2 : [0.224, 4.324, …],
  5 : [0.3, 0.34, 0.009, …],
  …
}

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
   "rishi",
   "sunak",
   "prime",
   "minister"
]

**Ranking Engine**

DESM

Queen Mary
University of London

# Inner Workings

**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

## Preprocessing Engine

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

## Ranking Engine

### DESM

```
{
    "arsenal" : [0.23, 0.01 …],
    "first" : [2.9, 0.28, …],
    "rishi" : [0.02, 0.43, ...],
    "justice" : [1.32, 0.91, …],
    "sunak" : [0.64, 0.22, ...],
    "prime" : [0.01, 0.09, ...],
    "tennis" : [0.11, 0.83, …],
    "fired" : [9.1, 0.01 …],
    "ukraine" : [5.3, 0.22, …],
    "train" : [0.54, 0.63, …],
    "minister" : [0.98, 0.9, ...],
    …
}
```

Queen Mary
University of London

# Inner Workings



**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM
```
{
  "arsenal" : [0.23, 0.01 …],
  "first" : [2.9, 0.28, …],
  "rishi" : [0.02, 0.43, ...],
  "justice" : [1.32, 0.91, …],
  "sunak" : [0.64, 0.22, ...],
  "prime" : [0.01, 0.09, ...],
  "tennis" : [0.11, 0.83, …],
  "fired" : [9.1, 0.01 …],
  "ukraine" : [5.3, 0.22, …],
  "train" : [0.54, 0.63, …],
  "minister" : [0.98, 0.9, ...],
  …
}
```

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM

```
{
  "rishi" : [0.02, 0.43, ...],
  "sunak" : [0.64, 0.22, ...],
  "prime" : [0.01, 0.09, ...],
  "minister" : [0.98, 0.9, ...]
}
```

# Inner Workings

**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM
Query Embeddings: {
  "rishi" : [0.02, 0.43, ...],
  "sunak" : [0.64, 0.22, ...],
  "prime" : [0.01, 0.09, ...],
  "minister" : [0.98, 0.9, ...]
}

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM

Document Centroids: {
 1 : [0.002, 0.23, 0.03, ...],
 2 : [0.224, 4.324, …],
 5 : [0.3, 0.34, 0.009, …],
 …
}

Query Embeddings: {
 "rishi" : [0.02, 0.43, ...],
 "sunak" : [0.64, 0.22, ...],
 "prime" : [0.01, 0.09, ...],
 "minister" : [0.98, 0.9, ...]
}

# Inner Workings

**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

**Ranking Engine**

DESM

Document Centroids: {
 1 : [0.002, 0.23, 0.03, ...],
 2 : [0.224, 4.324, …],
 5 : [0.3, 0.34, 0.009, …],
 …
}

Query Embeddings: {
 "rishi" : [0.02, 0.43, ...],
 "sunak" : [0.64, 0.22, ...],
 "prime" : [0.01, 0.09, ...],
 "minister" : [0.98, 0.9, ...]
}

# Inner Workings

**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
   "rishi",
   "sunak",
   "prime",
   "minister"
]

**Ranking Engine**

DESM

Document Centroids: {
  1 : 0.0023,
  2 : [0.224, 4.324, …],
  5 : [0.3, 0.34, 0.009, …],
  …
}

Query Embeddings: {
 "rishi" : [0.02, 0.43, ...],
 "sunak" : [0.64, 0.22, ...],
 "prime" : [0.01, 0.09, ...],
 "minister" : [0.98, 0.9, ...]
}

# Inner Workings

**Frontend**

**Backend**

### Ranking Engine

DESM

Document Centroids: {
  1 : 0.0023,
  2 : 0.04323,
  5 : 0.002395,
  …
}

Query Embeddings: {
  "rishi" : [0.02, 0.43, ...],
  "sunak" : [0.64, 0.22, ...],
  "prime" : [0.01, 0.09, ...],
  "minister" : [0.98, 0.9, ...]
}

**GET**

"Is Rishi Sunak Prime Minister?"

### Preprocessing Engine

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

# Inner Workings

**Frontend**

**Backend**

"Is Rishi Sunak Prime Minister?"

**GET**

**Ranking Engine**

DESM

```
{
 1 : 0.0023,
 2 : 0.04323,
 5 : 0.002395,
 …
}
```

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

# Inner Workings



Frontend | Backend

"Is Rishi Sunak Prime Minister?"

**GET**

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM

```
{
  1 : 0.0023,
  2 : 0.04323,
  5 : 0.002395,
  …
}
```

**Results**

# Inner Workings



**Frontend**

**Backend**

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

```
[
    "rishi",
    "sunak",
    "prime",
    "minister"
]
```

**Ranking Engine**

DESM

```
{
  1 : 0.0023,
  2 : 0.04323,
  5 : 0.002395,
  …
}
```

**Results**

```
{
1 : {score: 0.043, doc_id: 2},
2 : {score: 0.0024, doc_id: 5},
...,
10: {score: 0.001, doc_id: 73}
}
```

# Inner Workings

**Frontend**

**Backend**

**Ranking Engine**

DESM

**GET**

"Is Rishi Sunak Prime Minister?"

**Preprocessing Engine**

[
    "rishi",
    "sunak",
    "prime",
    "minister"
]

{
  1 : 0.0023,
  2 : 0.04323,
  5 : 0.002395,
  ...
}

**Results**

{
1 : {score: 0.043, doc_id: 2},
2 : {score: 0.0024, doc_id: 5},
...,
10: {score: 0.001, doc_id: 73}
}

Queen Mary
University of London

# Inner Workings

# Presentation Structure

| Stephanie | Iqbal | Berkay | Hannah |
|-----------|-------|--------|--------|
| **Overview** | **Implementation** | **Inner workings** | **Results and Limitations** |

# Results



**Legend:**

Retrieval method with the best search results for the particular query

Average score of the 10 shown articles

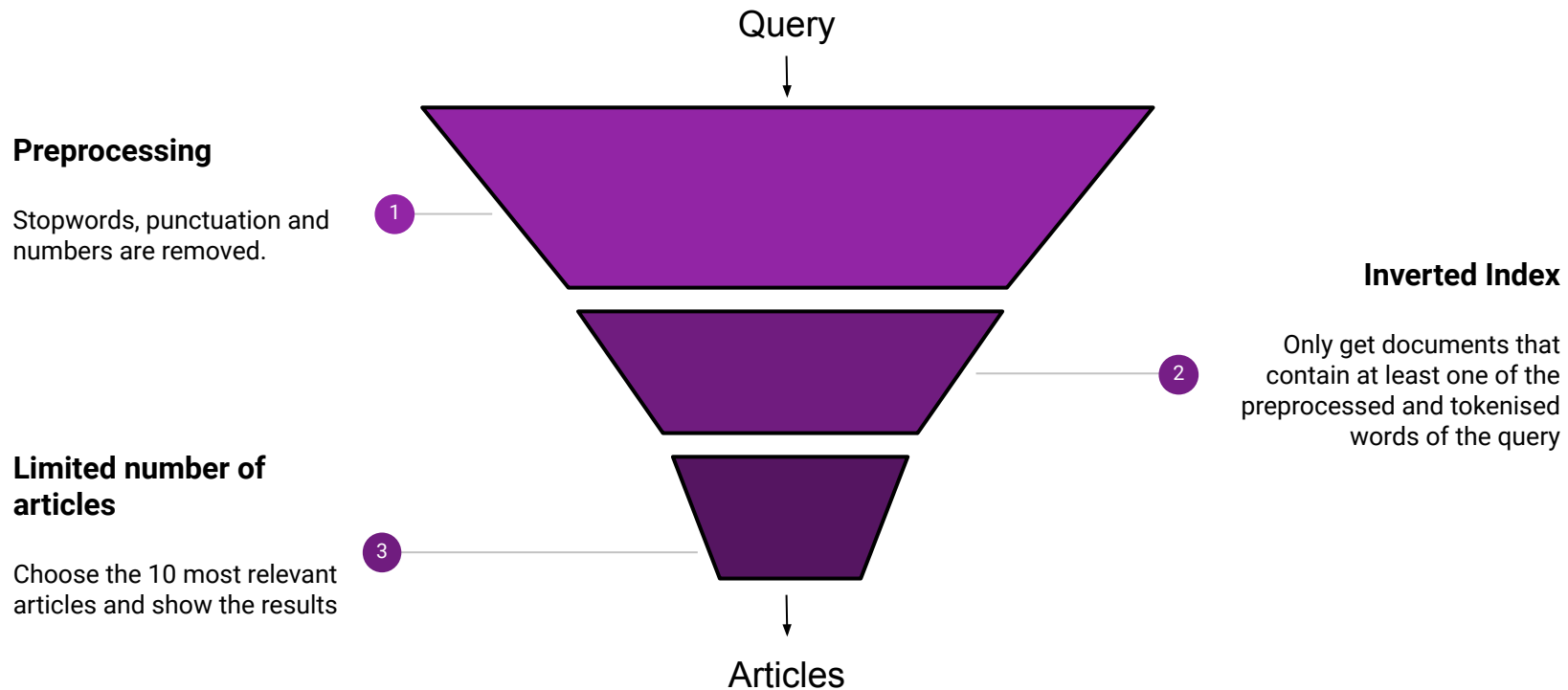Average score of the 10 shown articles after the relevance feedback has been submitted (Evaluation of the performance)

Green: Improvement of the average score
Red: Decrease of the average score

# Results

## Comparing implemented retrieval methods



DESM
48.0%

BM25
44.0%

No relevant results
8.0%

# Results



Query

**Preprocessing**

Stopwords, punctuation and numbers are removed.

**Inverted Index**

Only get documents that contain at least one of the preprocessed and tokenised words of the query

**Limited number of articles**

Choose the 10 most relevant articles and show the results

Articles

# Limitations

**Indexing**

Search engines cannot index all web pages due to lack of access or the great amount of data available

**Accuracy**

Search results cannot be fully accurate all the time, due to certain factors such as spam, clickbait, or ads

**05**

**01**

**02**

**03**

**04**

**Language**

Some language might not be as well represented as others, i.e. English, French, etc.The search engine might also prioritise results from specific regions

**Bias**

Search results can be influenced by personalisation, user data tracking or ads

**Context**

Search engines might not always understand the context of the query, i.e. apple as a fruit or the company

# Limitations



**Query: Chess**

Article with first index: BBC News CEO Secrets

Intended article: Chess gets a risqué makeover

# Limitations



**Query: Table**

Results: BBC Sport News

Intended article: Something about tables (furniture)

# Outlook

Include the order of words in the query, i.e. n-grams

Extend the human evaluation

Use a bigger and more diverse dataset, exceeding the 20 days window
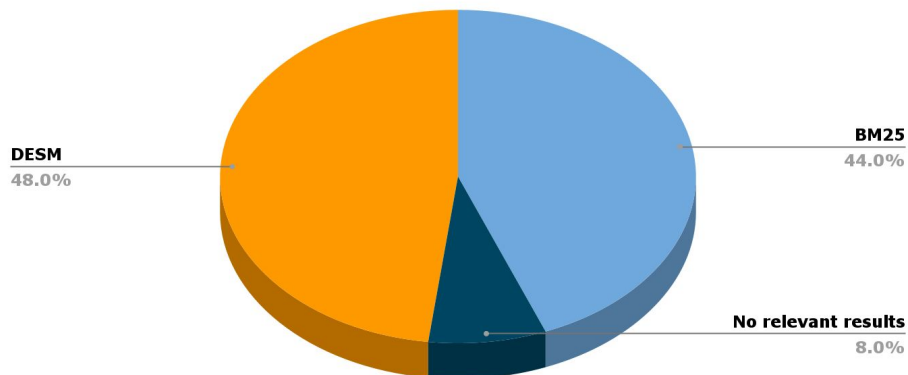
Use IN-OUT embeddings in the DESM

# Conclusion

DESM performs better than BM25

- DESM is more capable of showcasing relevant news according to the given query.

**Comparing implemented retrieval methods**



**DESM 48.0%**

**BM25 44.0%**

**No relevant results 8.0%**

## DESM > BM25

# Appendix

**Presentation's Journal Paper:**

Nalisnick, E. *et al.* (2016) "Improving document ranking with dual word embeddings," *Proceedings of the 25th International Conference Companion on World Wide Web - WWW '16 Companion*, pp. 83–84. Available at: https://doi.org/10.1145/2872518.2889361.

**Notable citations:**
- Guo, J., Cai, Y., Fan, Y., Sun, F., Zhang, R. and Cheng, X. (2022). Semantic Models for the First-Stage Retrieval: A Comprehensive Review. ACM Transactions on Information Systems, [online] 40(4), pp.1–42. doi:https://doi.org/10.1145/3486250.
- Muhammad, I., Bollegala, D., Coenen, F., Gamble, C., Kearney, A. and Williamson, P. (2021). Document Ranking for Curated Document Databases Using BERT and Knowledge Graph Embeddings: Introducing GRAB-Rank. Big Data Analytics and Knowledge Discovery, pp.116–127. doi:https://doi.org/10.1007/978-3-030-86534-4_10.
- Chy, A.N., Ullah, M.Z. and Aono, M. (2019). Query Expansion for Microblog Retrieval Focusing on an Ensemble of Features. Journal of Information Processing, [online] 27, pp.61–76. doi:https://doi.org/10.2197/ipsjjip.27.61.
- Khattab, O., Hammoud, M. and Elsayed, T. (2020). Finding the Best of Both Worlds. Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. doi:https://doi.org/10.1145/3397271.3401076.
- Järvelin, K. and Kekäläinen, J. (2017). IR evaluation methods for retrieving highly relevant documents. ACM SIGIR Forum, 51(2), pp.243–250. doi:https://doi.org/10.1145/3130348.3130374.
- Mitra, B., Nalisnick, E., Craswell, N., Caruana, R. and Redmond, M. (2016). *A Dual Embedding Space Model for Document Ranking*. [online] Available at: https://arxiv.org/pdf/1602.01137.pd
- Weng, L. (2017) *Learning word embedding*, *Lil'Log (Alt + H)*. Available at: https://lilianweng.github.io/posts/2017-10-15-word-embedding/ (Accessed: February 20, 2023).