
DIFFERENTIABLE SEARCH INDEXING *

Alessio Borgi, Eugenio Bugli, Damiano Imola

1952442, 1934824, 2109063

Sapienza Università di Roma

Rome

{borgi.1952442, bugli.1934824, imola.2109063}@studenti.uniroma1.it

ABSTRACT

Keywords First keyword · Second keyword · More

1 Introduction

Trainable Information Retrieval (IR) systems are characterized by two phases:

- **Indexing:** Indexing of a corpus, which means to associate the content of each document with its corresponding docid.
- **Retrieval:** Learn how to retrieve efficiently from the index, which means to

Instead of using Contrastive Learning based Dual Encoders, the paper proposes an architecture which directly map a query \mathbf{q} to a relevant docid \mathbf{j} . This architecture, called DSI, it's implemented with a pre-trained Transformer and all the information of the corpus are encoded within the parameters of the language model. When we are doing inference, the give to the trained model a text query as input and we expect to obtain a docid as output. If we are interested in a ranked list of relevant documents we can also use Beam Search. Our DSI system uses standard model inference to map from encodings to docids, instead of learning internal representations that optimize a search procedure. DSI can be extended in different ways:

- **Document representation:** there are several ways to represent documents (e.g. full text, bag-of-words representations, ...)
- **Docid representation:** (e.g. unique tokens, structured semantic docids, text strings, ...)

1.1 Indexing Methods

Given a sequence of document tokens, the model is trained to predict the docids. There can be used different strategies:

- **Inputs2Target:**
- **Targets2Inputs:**
- **Bidirectional:**

2 Dataset Generation

Given the MS Marco Dataset, we need to perform some preprocessing in order to create an association between queries and docids. This Dataset is composed by 100k passages text which are partitioned into Training, Validation and Test. Each partition is composed by lists of passages texts and related queries. We create a set $\mathcal{U} = \mathcal{O} \cup \mathcal{P}$, where \mathcal{O} is the

**Citation:* Authors. Title. Pages.... DOI:000000/11111.

set composed by each segment (passage text) of the MS Marco Dataset, while \mathcal{P} is the set of pseudo-queries from each segment of \mathcal{O} , which are artificially generated by a model (docT5query). After this procedure, we need to use segments belonging to \mathcal{U} as queries with a ranker. This procedure gives us the first \parallel ranked documents that have relevant information with respect to the segment. By doing this, we check if a segment has enough information to represent the document and can be used as a suitable query.

3 Filtering by Damiano

You have dense model which are good to maintain the knowledge of the inputs. Filtering: starting from a fragment, you use the dense model to perform ranking in order to obtain a list of k docid related to your document. You have obtained relevant fragments which will be inserted inside your training data.

4 Document Representation Strategies

5 Docid Representation for Retrieval

6 Headings: first level

See Section 6.

6.1 Headings: second level

6.1.1 Headings: third level

Paragraph

7 Examples of citations, figures, tables, references

[1].

The documentation for natbib may be found at

<http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf>

Of note is the command `\citet`, which produces citations appropriate for use in inline text. For example,

```
\citet{hasselmo} investigated\dots
```

produces

Hasselmo, et al. (1995) investigated...

<https://www.ctan.org/pkg/booktabs>

7.1 Figures

See Figure 1. Here is how you add footnotes.²

7.2 Tables

See awesome Table 1.

8 Conclusion

Your conclusion here

²Sample of the first footnote.

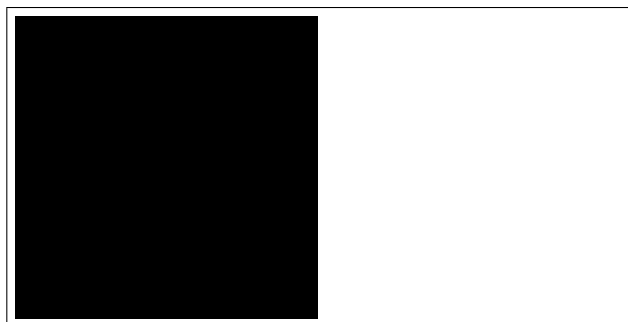


Figure 1: Sample figure caption.

Table 1: Sample table title

| Part | | |
|----------|-----------------|------------------------|
| Name | Description | Size (μm) |
| Dendrite | Input terminal | ~ 100 |
| Axon | Output terminal | ~ 10 |
| Soma | Cell body | up to 10^6 |

Acknowledgments

This was supported in part by.....

References

- [1] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations. In *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*, pages 2356–2362, 2021.