# Navigating Neural Search: Avoiding Common Pitfalls

Jo Kristian Bergum @jobergum

# This talk

- Highlight common pitfalls
- Pre-trained Language Models (PLM)
- Quick overview of neural search using PLM
    - Three neural models built on pretrained language models (PLM)
- Text embedding models and embedding retrieval

# Not in this talk

- Retrieval Augmented Generation (RAG)
- Generative Large Language Models (GPT, LLAMA)

# Pretrained Language Models (PLM)

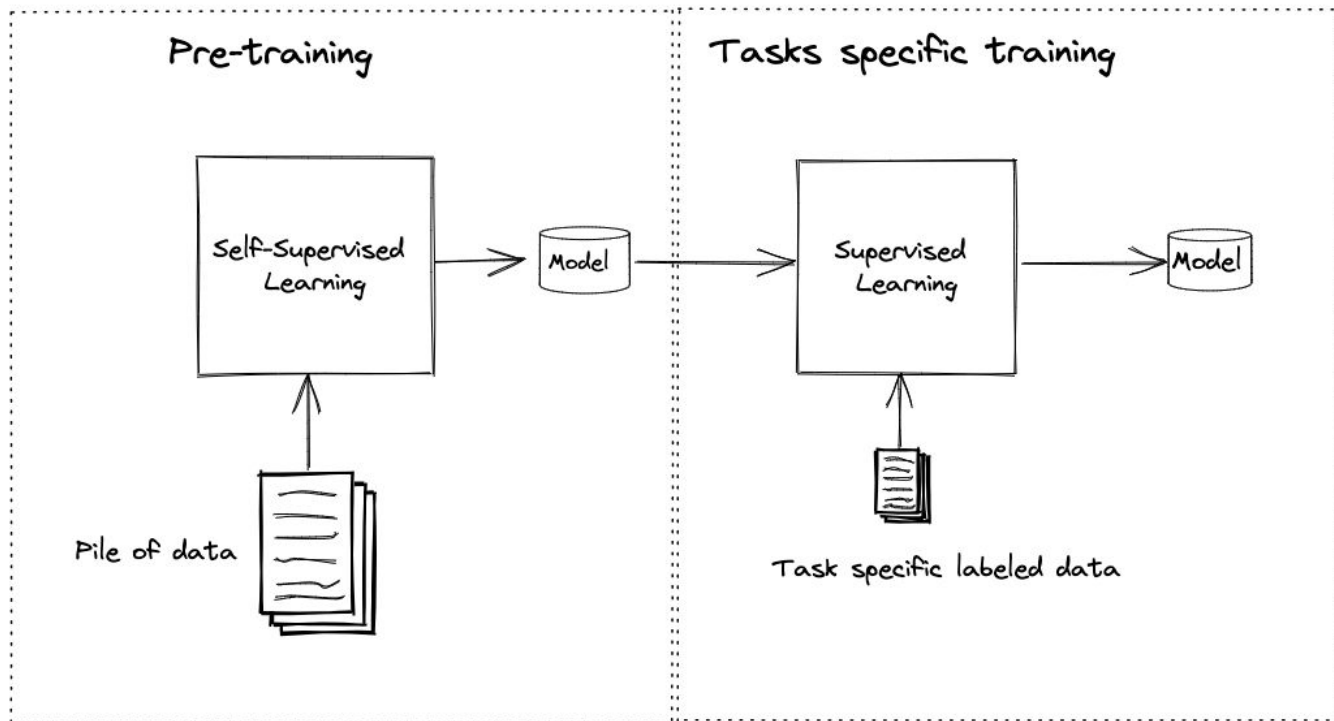- Attention is All you Need (Google 2017)
- BERT (Bidirectional Encoder Representations from Transformers)
- Trained using masked word language objective

*The cat sits on the [MASK] looking at the [MASK]*

- Masking objective is genius - Enables self-supervision with large corpuses of text
- Pre-trained model weights uses as starting weights for downstream tasks
  - Search
  - Classification
  - And more

# Transfer Learning 101

# Language Model

**A tokenizer + fixed vocabulary**

**A deep neural network architecture**

**Small, medium, large, xxx large?**

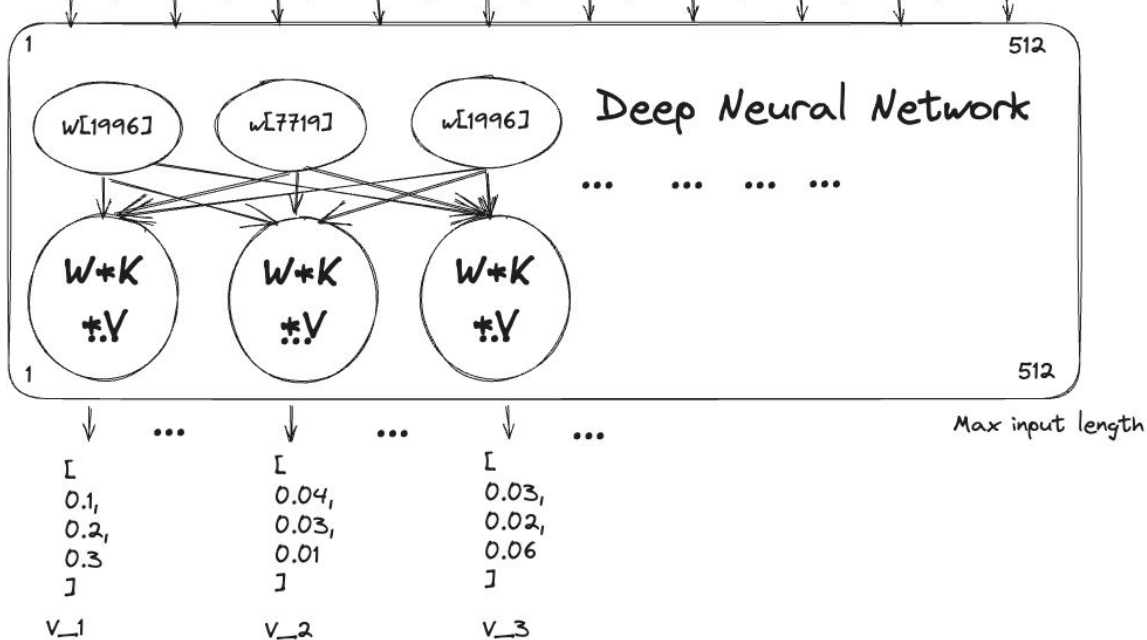"The cat sits on the table looking at the dog"

Tokenization

[1996, 4937, 7719, 2006, 1996, 2795, 2559, 2012, 1996, 3899]

1                512

W[1996]    W[7719]    W[1996]    Deep Neural Network

...    ...    ...    ...

W*K *.V     W*K *.V     W*K *.V

1                512

Max input length

[ 0.1, 0.2, 0.3 ]     [ 0.04, 0.03, 0.01 ]     [ 0.03, 0.02, 0.06 ]

$V\_1$       $V\_2$       $V\_3$

# LM Tokenization

**Happy Path Tokenization**

**10 words maps to 10 token ids**

"The cat sits on the table looking at the dog"

Tokenization

['the', 'cat', 'sits', 'on', 'the', 'table', 'looking', 'at', 'the', 'dog']
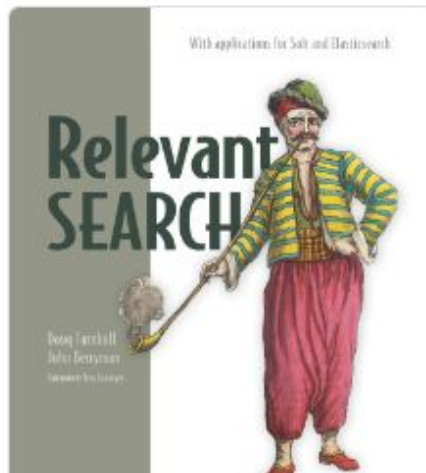
[1996, 4937, 7719, 2006, 1996, 2795, 2559, 2012, 1996, 3899]

# LM tokenization

- Different tokenizer implementations
- Tied to model
- Fixed vocabulary size
- Learned "word" embedding vectors per word in vocab
- Vocab fixed before pre-training of neural network weights

## Relevant Search ♡
### With applications for Solr and Elasticsearch

★★★★★ 7 reviews

Doug Turnbull and John Berryman
*Foreword by Trey Grainger*

June 2016 · ISBN 9781617292774 · 360 pages · printed in black & white

Data    Data Science

Tokenization

['isbn', '978', '##16', '##17', '##29', '##27', '##7', '##4']

# LM Tokenization

Are LM insensitive to spelling mistakes?

"Elasticsearch versus Vespa for vector search"

⬇

[21274, 17310, 11140, 6431, 2310, 13102, 2050, 2005, 9207, 3945]

"Elasticssearch versus Vespa for vector searrch"

⬇

[21274, 11393, 2906, 2818, 6431, 2310, 13102, 2050, 2005, 9207, 2712, 12171, 2818]

# LLM tokenization impact vector representation

## Variant and tokens

annoyance => annoyance

anoyance => ['an', '##oya', '##nce']

annyoance => ['ann', '##yo', '##ance']

## Top-3 retrieved words (vector search over WordNet)

frustration, anger, rage (0.91)

loyalty, consciousness, treasure (0.83)

anniversary, old age, tendency (0.84)

WordNet® is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept.

# LM tokenization

Linguistics and language matters !(?)

- Multilingual
- English

Not that many language specific LM models (except for English)

Don't know newer words

- 2023 (202, ##3)
- Covid-19 (co,##vid,-, 19)
- GPT (gp, ##t)

"die Katze sitzt auf dem tisch und schaut den hund an"

⬇ English Tokenizer (word piece)

['die','katz','##e','sit','##z','##t','auf','dem','tis','##ch','und','sc','##ha','##ut','den','hu','##nd','an']

"die Katze sitzt auf dem tisch und schaut den hund an"

⬇ Multilingual tokenizer (sentence piece)

['_die','_ka','tze','_si','tzt','_auf','_dem','_Tisch','_und','_schaut','_den','_Hund','_an']
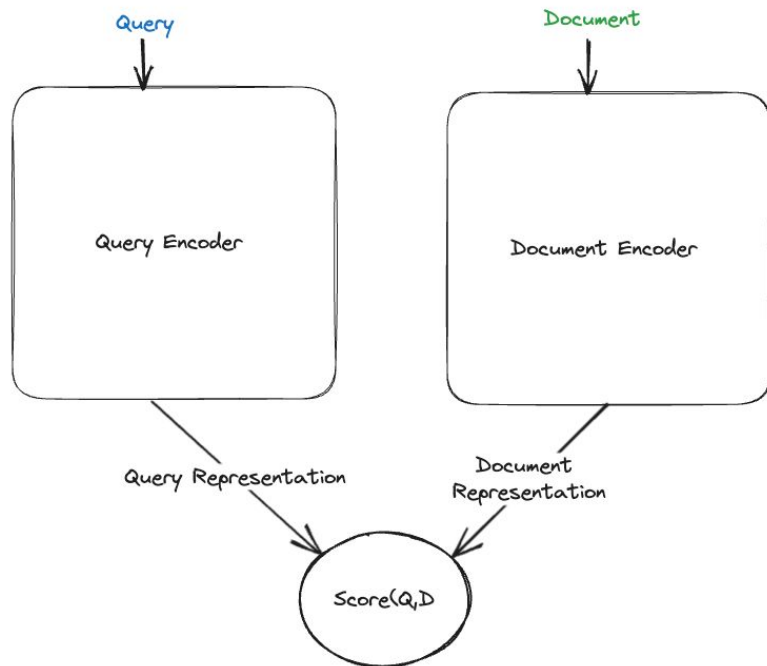
# Applying LMs to search

# Searching over data with sublinear complexity

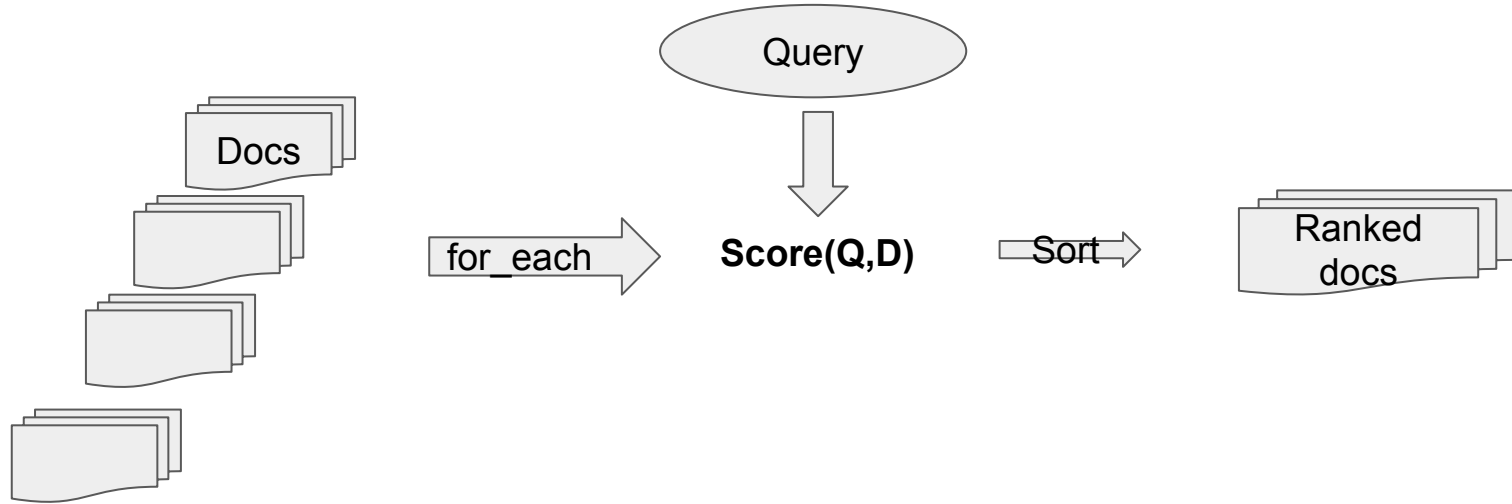*Conceptual representational model for retrieval*

- Representation of queries and documents
    - So that relevant documents for a query is scored higher than irrelevant documents
- Dense/Sparse/Mixed
- Score(Q,D) complexity constraints
- Supervised (learned) versus unsupervised

Bi-Encoder Architecture

Query

Document

Query Encoder

Document Encoder

Query Representation

Document Representation

Score(Q,D)

# Motivation for representational approach

Avoid scoring all documents D in collection for a query Q

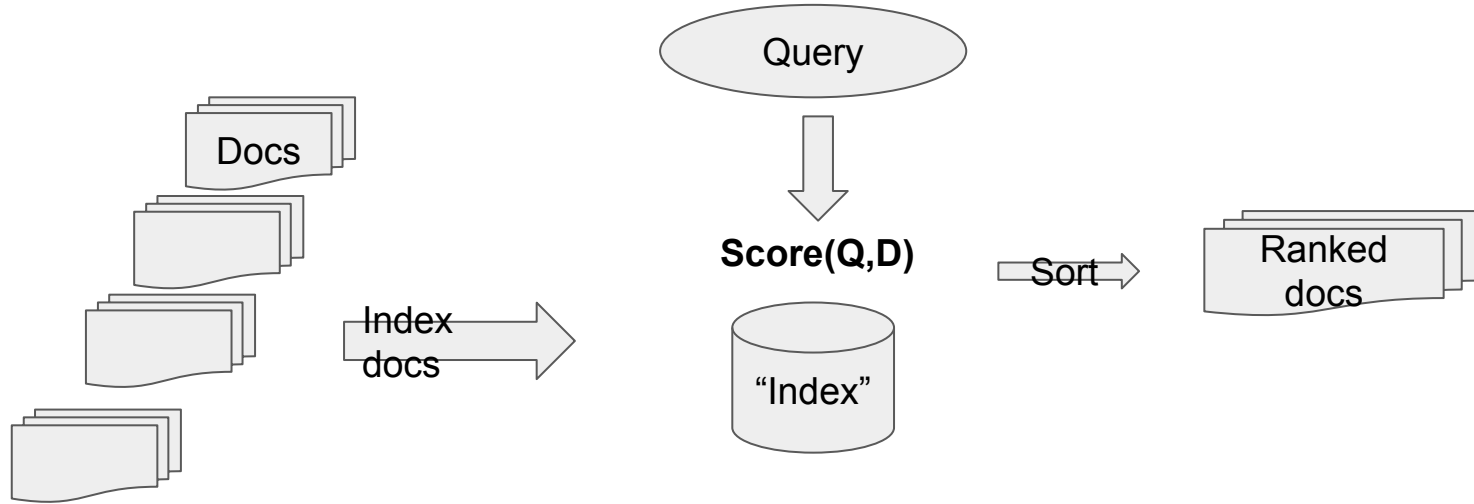Docs

Query

for_each

**Score(Q,D)**

Sort

Ranked docs

# Motivation for representational approach

Avoid scoring all documents D in collection for a query Q

# Make it more concrete

Logical representation versus physical implementation.

**Accelerating scoring over sparse representations**

- Build Inverted Index data structures
- Search accelerated with algorithms like WAND, MaxScore, BM-WAND++

**Accelerating scoring over dense representations**

- Build Vector Index (IVF, Quantization, HNSW, ++)
- Search accelerated with algorithms tied to vector index structure

# Also: Phased retrieval and ranking

100s — global-phase ranking

Thousands — Second phase ranking
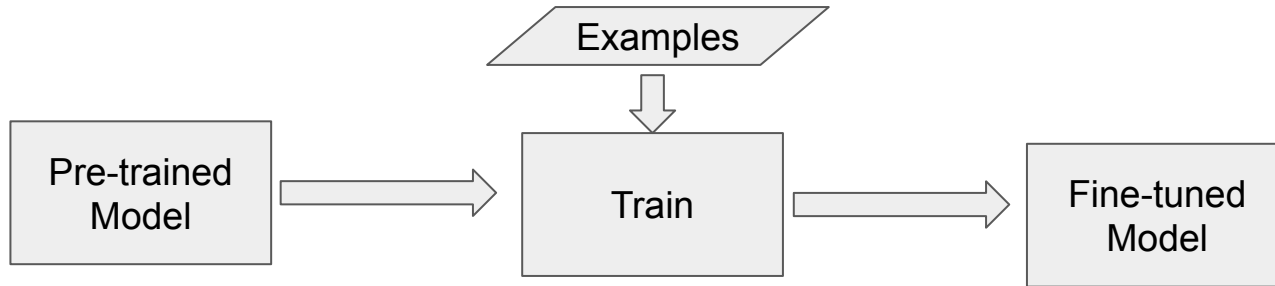
Millions — First phase ranking

Billions

# 3 Neural Methods for Search using LM

All methods require - **Labeled examples** - usually triplets

*<query, relevant document, irrelevant document>*

# Cross-Encoder

Encodes both query and document at the same time (cross)

all-to-all attention between all tokens in query and document

Most effective on IR benchmarks (nDCG)

High compute complexity (n^2)

No efficient way to "index"

Query

"who sits on tables"

Document

"The cat sits on the table looking at the dog"

101   2040, 7719, 2006, 7251   102   1996, 4937, 7719, 2006, 1996, 2795, 2559, 2012, 1996, 3899   102

Language Model
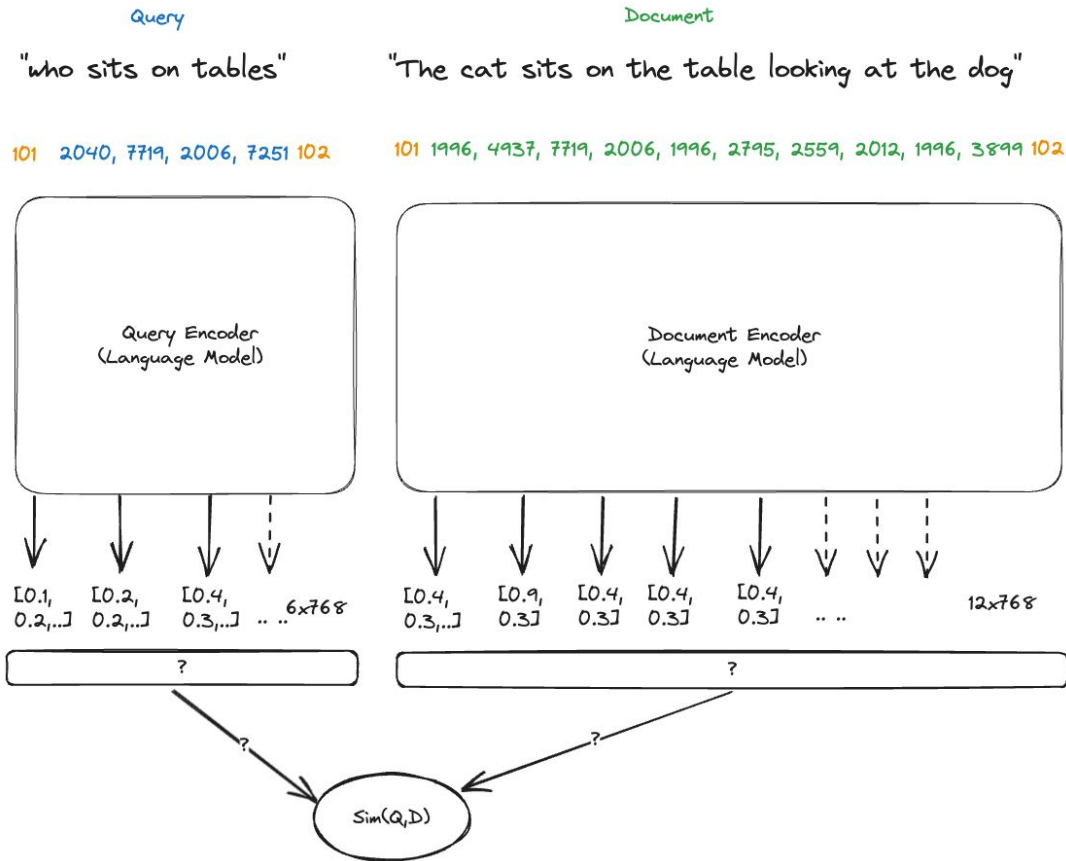
512

Classification Layer

Score

# Bi-Encoder

Encode queries and documents independently

No token level attention between query and document (no cross)

Enables indexing documents offline

Sim(Q,D):

- Dot product (sparse or dense)
- Cosine/Euclidean/Hamming/Many

Query

"who sits on tables"

101  2040, 7719, 2006, 7251 102

Query Encoder
(Language Model)

[0.1,    [0.2,    [0.4,
0.2,..]  0.2,..]  0.3,..] .. ..    6x768

?

Document

"The cat sits on the table looking at the dog"

101 1996, 4937, 7719, 2006, 1996, 2795, 2559, 2012, 1996, 3899 102

Document Encoder
(Language Model)

[0.4,    [0.9,   [0.4,   [0.4,   [0.4,
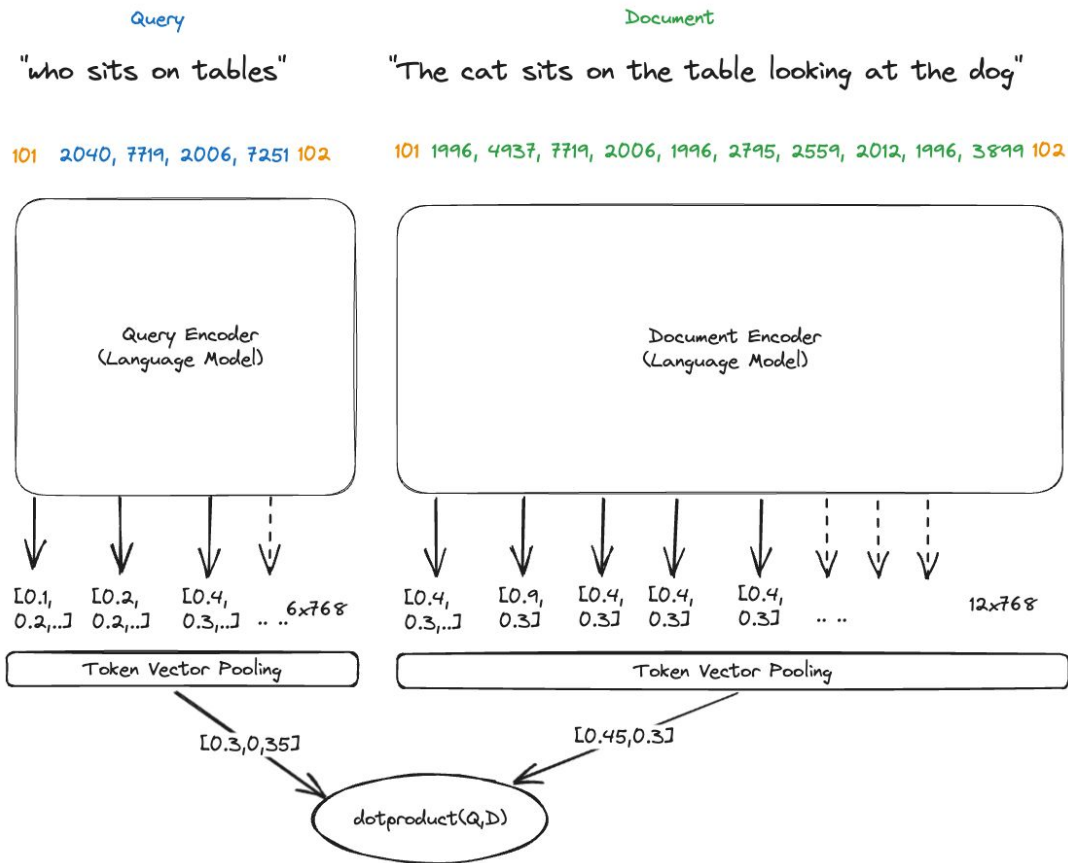0.3,..]  0.3]    0.3]    0.3]    0.3]    .. ..    12x768

?

Sim(Q,D)

# Bi-Encoder

Output Pooling

From a **token** vector representations to a vector representation of a sequence

- Average?
- 101/CLS token?

# Bi-Encoder adv

Learn token vectors instead of sequence vectors

Not pooled

Query

"who sits on tables"

Document

"The cat sits on the table looking at the dog"

101 2040, 7719, 2006, 7251 102

101 1996, 4937, 7719, 2006, 1996, 2795, 2559, 2012, 1996, 3899 102

Query Encoder
(Language Model)

Document Encoder
(Language Model)

[0.1, 0.2,..] [0.2, 0.2,..] [0.4, 0.3,..] .. .. 6x768

[0.4, 0.3,..] [0.9, 0.3] [0.4, 0.3] [0.4, 0.3] [0.4, 0.3] .. .. 12x768
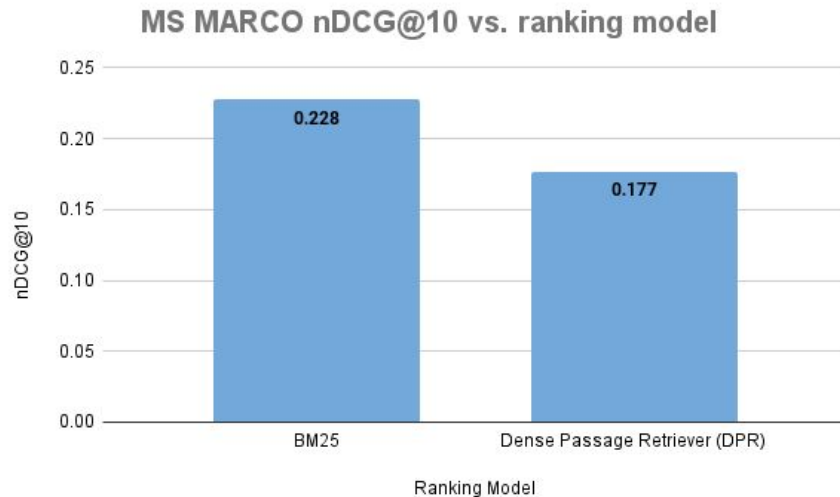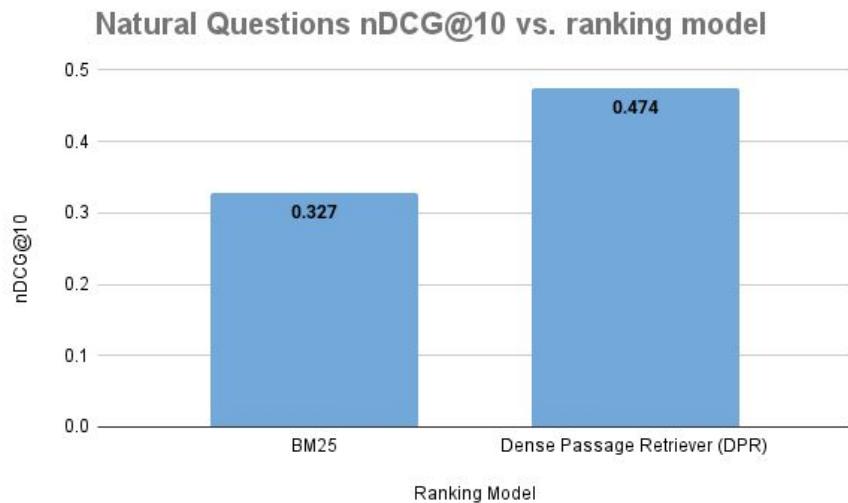
Per token vector rep (ColBERT)

Sum of max dot products

# Learned representations - No better than the examples?

Remember: The representation of queries and documents are learned

- Your data might not look like the examples



Natural Questions nDCG@10 vs. ranking model

| | |
|---|---|
| BM25 | 0.327 |
| Dense Passage Retriever (DPR) | 0.474 |

Ranking Model



MS MARCO nDCG@10 vs. ranking model

| | |
|---|---|
| BM25 | 0.228 |
| Dense Passage Retriever (DPR) | 0.177 |

Ranking Model

# Data the vector model was trained on

Photo by <u>Vidar Nordli-Mathisen</u> on
<u>Unsplash</u>

# Your data

# TLDR; Neural Methods for Retrieval & Ranking

Accuracy versus cost

Model not better than the examples it was trained on

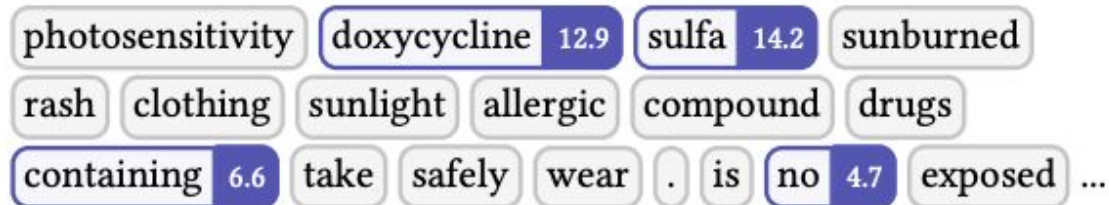Explain/score interpretability difficult with pooled representations

Introducing Neural Bag of Whole-Words with ColBERTer

https://arxiv.org/abs/2203.13088

**Q does doxycycline contain sulfa**

*BERT tokenized (9 subword-tokens): 'does', 'do', '##xy',*
*'##cy', '##cl', '##ine', 'contain', 'sul', '##fa'*

**ColBERTer BOW$^2$** *(30 saved vectors from 84 subword-tokens):*

| photosensitivity | doxycycline 12.9 | sulfa 14.2 | sunburned |
| rash | clothing | sunlight | allergic | compound | drugs |
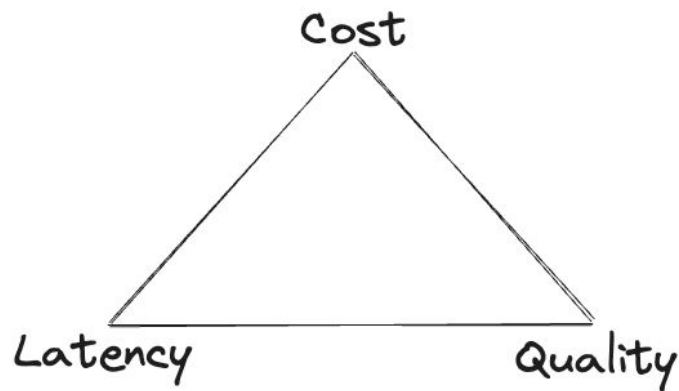| containing 6.6 | take | safely | wear | . | is | no 4.7 | exposed | ...

**Fulltext:** No doxycycline is not a sulfa containing compound, so you may take it safely if you are allergic to sulfa drugs. You should be aware, however, that doxycycline may cause photosensitivity, so you should wear appropriate clothing, or you may get easily sunburned or develop a rash if you are exposed to sunlight.
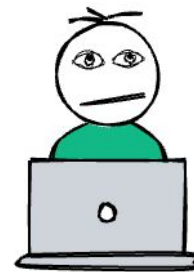
**Figure 1: Example of ColBERTer's BOW$^2$ (Bag Of Whole-Words): ColBERTer stores and matches unique whole-word representations. The words in BOW$^2$ are ordered by implicitly learned query-independent term importance. Matched words are highlighted in blue with whole-word scores displayed in a user-friendly way next to them.**

# Off-the-shelf text-embedding models

- Size of model
- Embedding dimensionality
- Sequence length
- Quality/Accuracy (for your use case)
- Language capabilities
- Licence/Commercial use

Cost

Latency

Quality

Which Embedding model?

# MTEB (massive text embedding benchmark)

Great guide

Many different tasks

https://huggingface.co/spaces/mteb/leaderboard

Benchmark hacks?

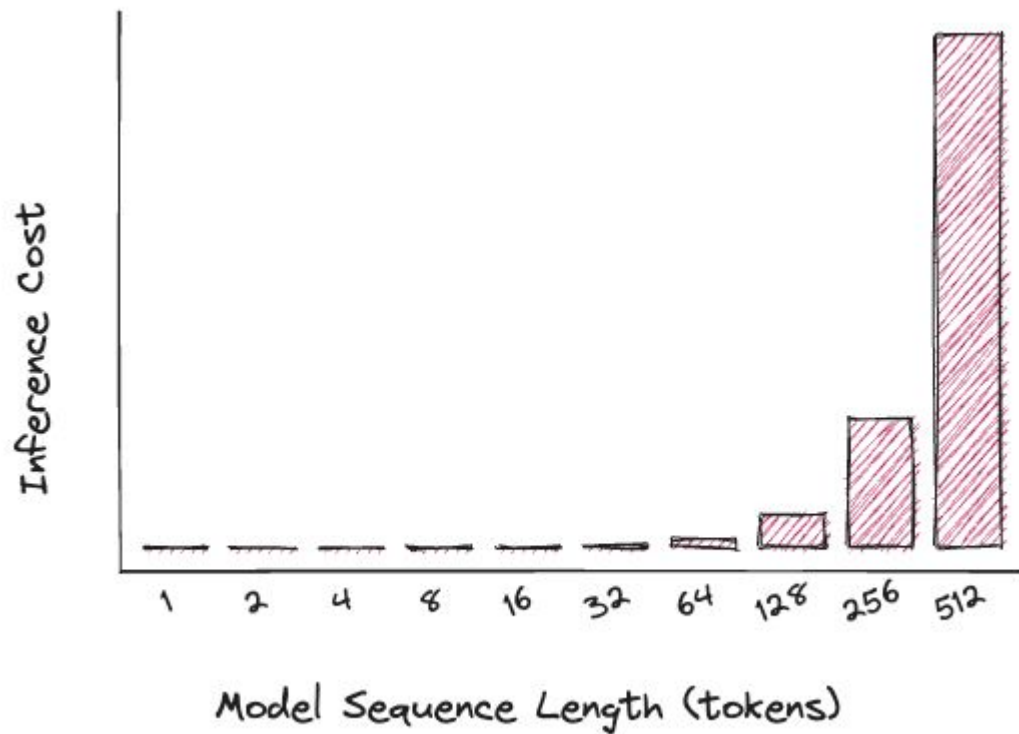| Rank | Model | Model Size (GB) | Embedding Dimensions | Sequence Length | Average (56 datasets) |
|---|---|---|---|---|---|
| 1 | bge-large-en-v1.5 | 1.34 | 1024 | 512 | 64.23 |
| 2 | bge-base-en-v1.5 | 0.44 | 768 | 512 | 63.55 |
| 3 | gte-large | 0.67 | 1024 | 512 | 63.13 |
| 4 | gte-base | 0.22 | 768 | 512 | 62.39 |
| 5 | e5-large-v2 | 1.34 | 1024 | 512 | 62.25 |
| 6 | bge-small-en-v1.5 | 0.13 | 384 | 512 | 62.17 |
| 7 | instructor-xl | 4.96 | 768 | 512 | 61.79 |
| 8 | instructor-large | 1.34 | 768 | 512 | 61.59 |

# Embedding Retrieval

Embedding inference + Retrieval

Model size (GPU needed?)

Sequence length scaling

Dimensionality

    1536 dims (4x cost of 384)

    Not 4x accuracy !



Inference Cost vs. Model Sequence Length (tokens): 1, 2, 4, 8, 16, 32, 64, 128, 256, 512

# Vector Search

Brute Force Search Might Be All You Need?

Assume 64GB/s memory bandwidth

1M vectors with1536 dimensions using float is approx 6GB

Quiz: How many QPS can one node support at max?

# The A in ANN

Approximate search instead of brute-force search

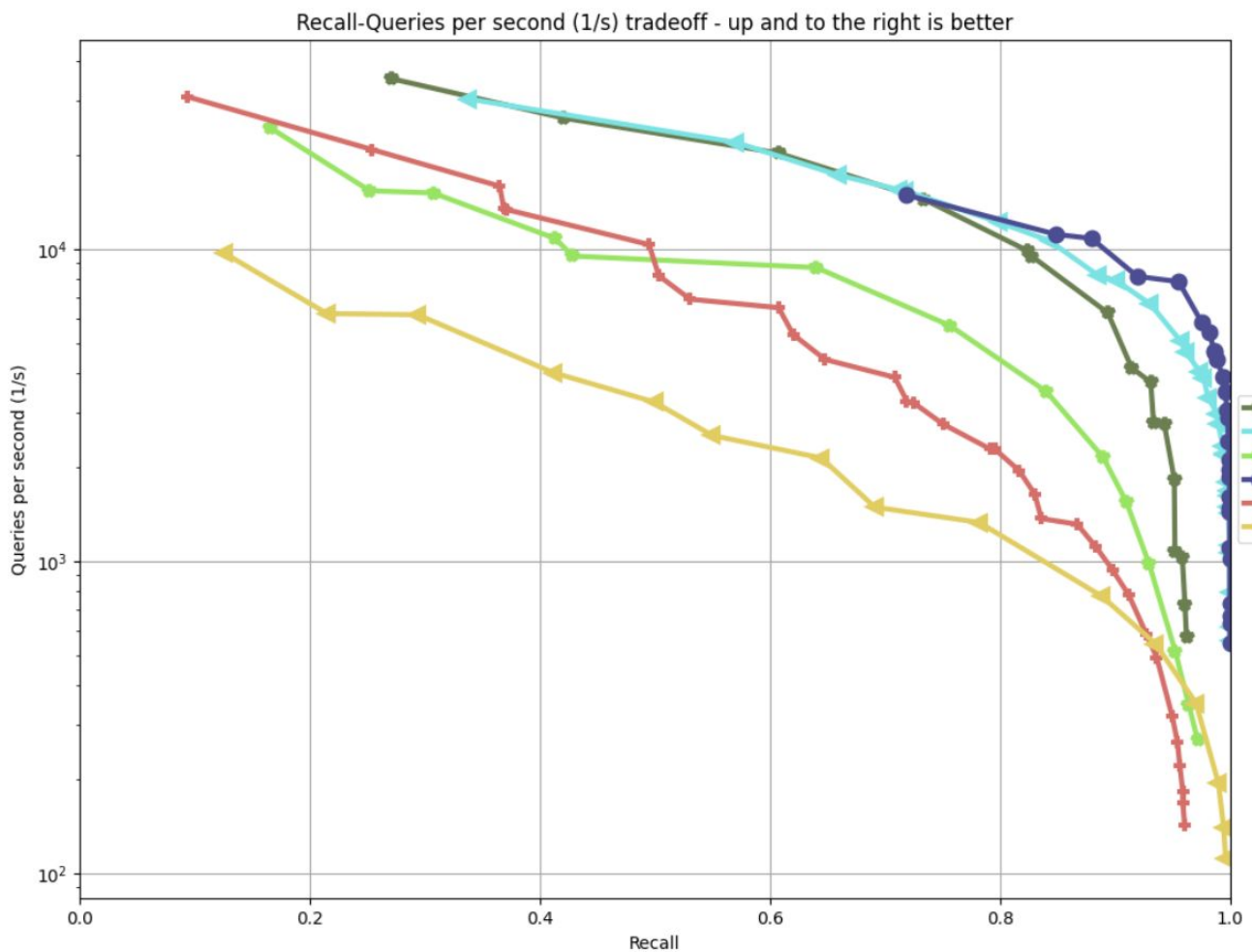Speed up retrieval, by building an index, **sounds familiar**?

Many different ANN algorithms and associated tradeoffs

- Query speedup
- Quality (What is the error introduced by approximate search)
- Real-time (Mutatable, grow from to zero to N)
- Resource footprint, index build time

Exact and
Approximate
(recall@k)

Overlap@k is a better
name for us working
with search metrics



Recall-Queries per second (1/s) tradeoff - up and to the right is better

Queries per second (1/s)

Recall

# Impact of ANN choice & parameters on search quality

Our search quality metrics

- Recall (Are we finding all the relevant hits)
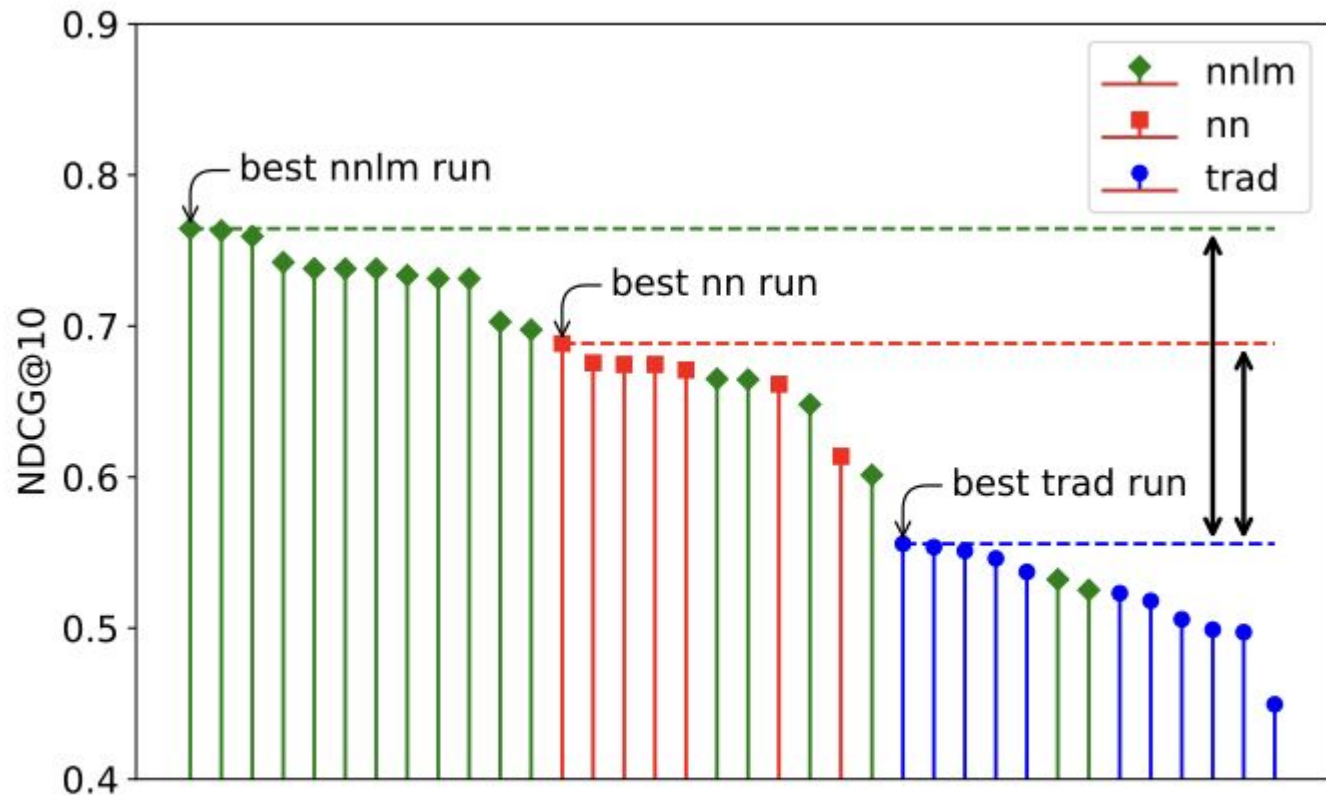- Precision (Are we finding nothing but relevant?)

# LADR

[https://arxiv.org/abs/2307.16779](https://arxiv.org/abs/2307.16779), BM25 on DL19 is about 0.55 NDCG@10

| Method | DL19 ~4ms | | DL19 ~8ms | | DL20 ~4ms | | DL20 ~8ms | |
|---|---|---|---|---|---|---|---|---|
| | nDCG | R@1k | nDCG | R@1k | nDCG | R@1k | nDCG | R@1k |
| **TAS-B** (Exh.) | 0.715 | 0.842 | 0.715 | 0.842 | 0.713 | 0.875 | 0.713 | 0.875 |
| IVF [$I$] | 0.374 | 0.414 | 0.474 | 0.536 | 0.503 | 0.559 | 0.579 | 0.677 |
| ScaNN [$S$] | 0.475 | 0.519 | 0.537 | 0.598 | 0.476 | 0.527 | 0.553 | 0.641 |
| HNSW [$H$] | - | - | 0.614 | 0.707 | - | - | 0.699 | 0.836 |
| GAR [$G$] | 0.543 | 0.540 | 0.688 | 0.755 | 0.568 | 0.594 | 0.684 | 0.796 |
| Re-Ranking [$R$] | 0.589 | 0.605 | 0.684 | 0.755 | 0.615 | 0.667 | 0.691 | 0.805 |
| Proactive LADR | $_{GR}^{IS}$**0.690** | $_{GR}^{IS}$**0.771** | $_{GR}^{ISH}$0.730 | $_{GR}^{ISH}$0.850 | $_{GR}^{IS}$**0.691** | $_{GR}^{IS}$**0.807** | $_{GR}^{ISH}$0.722 | $_{GR}^{IS}$0.857 |
| Adaptive LADR | - | - | $_{GR}^{ISH}$**0.738** | $_{GR}^{ISH}$**0.872** | - | - | $_{GR}^{ISH}$**0.739** | $_{GR}^{ISH}$**0.900** |

| max-links-per-node | neighbors-to-explore-at-insert | hnsw.exploreAdditionalHits | NDCG@10 |
|---|---|---|---|
| 16 | 100 | 0 | 0.5115 |
| 16 | 100 | 100 | 0.6415 |
| 16 | 100 | 300 | 0.6588 |
| 32 | 500 | 0 | 0.6038 |
| 32 | 500 | 100 | 0.6555 |
| 32 | 500 | 300 | 0.6609 |

*Summarization of the HNSW parameters and the impact on NDCG@10.*

As the table above demonstrates, we can reach the same NDCG@10 as the exact search by using `max-links-per-node` 32, `neighbors-to-explore-at-insert` 500, and `hnsw.exploreAdditionalHits` 300. The high `hnsw.exploreAdditionalHits` setting indicates that we could alter the index time settings upward, but we did not experiment further. Note the initial HNSW setting in row 1 and the significant negative impact on retrieval quality.

(b) Passage retrieval task

# TLDR;

- Tokenization and vocabulary matters
- Language matters
- **Representations, representations, representations**
- Your data (queries and documents) might not match training examples

- Embedding inference
    - Sequence length
    - Dimensionality
- Embedding retrieval (vector search)
    - Brute force versus approximate

- Approximate Search Does Introduce Errors..

# Resources

Lots on Blog.vespa.ai, for example

https://blog.vespa.ai/improving-zero-shot-ranking-with-vespa-part-two/

https://blog.vespa.ai/accelerating-transformer-based-embedding-retrieval-with-vespa/

# Hated it? Tweet me

Jo Kristian Bergum  🐦 jobergum