# Faster Learned Sparse Retrieval with Guided Traversal

Antonio Mallia, Joel Mackenzie, Torsten Suel, and Nicola Tonellotto





### **Introduction and Preliminaries**

**Context:** Recently, a number of neural-based techniques have been explored for improving inverted index based retrieval.

- Term Expansion: Both offline (applied to documents) and online (applied to queries).
- Term (Re)weighting: Both offline (reweighted term frequencies, or learned impacts) and online (query term weighting).

search is cool search is fun search is fun for everyone



#### Postings Lists

search	•	<b></b>	0	1	1	1	2	1
cool	•		0	1				
fun	•		1	1	2	1		
everyone	•		2	1				



Score documents with a bag-of-words ranker -BM25, Language Model, DPH, etc



#### Postings Lists

search	•	<b></b>	0	1	1	1	2	1
cool	•		0	1				
fun	•		1	1	2	1		
everyone	•		2	1				



**Idea:** Add questions to documents to avoid vocabulary mismatch!

## search is cool

# search is cool + why is search cool? + who thinks search is cool?

+...



Idea: Add questions to documents to avoid vocabulary mismatch!

- Doc2Query, DocTTTTTQuery, TILDE



Idea: Add questions to documents to avoid vocabulary mismatch!

- Doc2Query, DocTTTTTQuery, TILDE

Expanded documents  $\rightarrow$  longer postings lists



#### Postings Lists

search	•	<b></b>	0	1	1	1	2	1
cool	•		0	1				
fun	•		1	1	2	1		
everyone	•		2	1				





- DeepCT, HDCT
- Use contextual language models to re-weight tf's
- Still score with traditional models!











#### Postings Lists

search	•	<b></b>	0	1	1	1	2	1
cool	•		0	1				
fun	•		1	1	2	1		
everyone	•		2	1				





- DeepImpact
- TILDE(v2)



- DeepImpact
- TILDE(v2)



- DeepImpact
- TILDE(v2)



- DeepImpact
- TILDE(v2)





+ Learn weights for each query



- + Learn weights for each query
- uniCOIL, SPLADEv2

#### Default

999391 where is bulli creek queensland australia

#### uniCOIL-TILDE

999391 where where

#### [where 42][is 44][bull 194][##i 116][creek 165] [queensland 149][australia 70]

#### SPLADEv2 [Adds query expansion]

[where 268] [bull 237] [creek 219] [queensland 211] [australia 183] [##i 132] [headquarters 92] [australian 86] [habitat 72] [stream 71] [brisbane 64] [him 40] [lake 37] [scotland 31] [it 29] [are 7] [from 6] [river 5] [italy 3] **Problem:** Learned sparse models negatively impact processing latency.

**Purpose:** Make learned sparse models as fast as traditional models.

Approach: Heuristics.

## Why are learned sparse models slower?

#### **Dynamic Pruning at a glance**

Assumption: Ranking is *additive*.

**Requirement:** We have pre-computed the maximum impact for each term (postings list) and stored it.

**Requirement:** During processing, we have access to the *k* th highest score "seen so far" (a threshold  $\theta$ ).

Score [and bound]

**Document Space**




br	onx
<i>ti</i>	he







**Document Space** 







**Document Space** 



**Document Space** 



**Document Space** 



**Document Space** 



**Document Space** 



**Document Space** 









**Document Space** 

















## But what if the distributions change?



**Document Space** 

Learned Models violate the embedded assumption of IDF: The more rare a term, the more highly it should be weighted with respect to the other terms.







## **Our Heuristic**

**Observation:** DeepImpact is built over a DocT5Query Index... Postings lists are (almost) one-to-one.

## DocT5 12 5 20 2 22 1 29 4 DeepImpact 12 8 20 7 22 1 29 9





## **Hypothesis:** BM25 visits the right documents; it just doesn't rank them well.
## **Guided Traversal (GT)**



**Interpolation** is easy/fast (we call this GTI). In our experiments, we just do an unweighted linear interpolation.





# **Limitations and Future Work**

Packing method is naïve and difficult to compress. Specialized solution; works only when postings align.

Only works in contexts where BM25 (over an expanded index) *visits* the "good" documents.

Concluding Notes: Benchmarking is hard!

### Improvements That Don't Add Up: Ad-Hoc Retrieval Results Since 1998

Timothy G. Armstrong, Alistair Moffat, William Webber, Justin Zobel

Computer Science and Software Engineering The University of Melbourne Victoria 3010, Australia {tgar,alistair,wew,jz}@csse.unimelb.edu.au

**CIKM 2009:** "How confident are we that a technique that yields an improvement over a weak baseline would also give an improvement over a strong one?"



Figure 5: MAP as a function of number of options turned on, for Indri running against the TREC-5 Ad-Hoc test collection.



Figure 5: MAP as a function of number of options turned on, for Indri running against the TREC-5 Ad-Hoc test collection.



Figure 5: MAP as a function of number of options turned on, for Indri running against the TREC-5 Ad-Hoc test collection.



Figure 5: MAP as a function of number of options turned on, for Indri running against the TREC-5 Ad-Hoc test collection.



Improvements are broadly additive

#### Examining Additivity and Weak Baselines

SADEGH KHARAZMI, RMIT University & NICTA FALK SCHOLER, RMIT University DAVID VALLET, Google MARK SANDERSON, RMIT University

#### On the Additivity and Weak Baselines for Search Result Diversification Research

Mehmet Akcay MiAdla East Technical University & ASELSAN Ankara, Turkey

#### Critically Examining the "Neural Hype": Weak Baselines and the Additivity of Effectiveness Gains from Neural Ranking Models

Wei Yang,1 Kuang Lu,2 Peilin Yang, and Jimmy Lin1 1 David R. Cheriton School of Computer Science, University of Waterloo <sup>2</sup> Department of Electrical and Computer Engineering, University of Delaware meakcav@aselsan.com.tr

Craig Macdonald University of Glasgow Glasgow, Scotland, UK craig.macdonald@glasgow.ac.uk

Ismail Sengor Altingovde Middle East Technical University Ankara, Turkey altingovde@ceng.metu.edu.tr

Iadh Ounis University of Glasrow Glasgow, Scotland, UK iadh.ounis@glasgow.ac.uk

#### Are We Really Making Much Progress? A Worrying Analysis of **Recent Neural Recommendation Approaches**

Maurizio Ferrari Dacrema Politecnico di Milano. Italy maurizio.ferrari@polimi.it

Paolo Cremonesi Politecnico di Milano. Italy paolo.cremonesi@polimi.it

Dietmar Iannach University of Klagenfurt, Austria dietmar.iannach@aau.at

### Examining the Additivity of Top-k Query Processing Innovations

Joel Mackenzie The University of Melbourne Melbourne, Australia Alistair Moffat The University of Melbourne Melbourne, Australia

## This also applies to efficiency studies!

Benchmarking is difficult...

CPU vs GPU? Which CPU and which GPU??

Optimized codebase? Research code? How much do we know (and need to know) about the underlying codebase?

Context of measurement; Scale of collections, hardware, etc... tension between latency, throughput, space occupancy, scalability, computing power and cost  $\rightarrow$  CO2