

Dense Neural Retrieval for Scientific Documents at Zeta Alpha

ECIR 2022 Industry Day, Stavanger - April 14th 2022

Jakub Zavrel, Marzieh Fadaee, Artem Grotov and Rodrigo Nogueira



🖾 Zeta Alpha

Made in Amsterdam, since 2019.

Leverage Neural Language Understanding to help people make better AI supported decisions. A smarter way to discover and organize knowledge for AI and Data Science teams

Zeta Alpha helps you win in the fast moving world of innovation and research in Al. Our user friendly, focused, state-ofthe-art discovery platform becomes your research assistant that streamlines how you organize, share, and stay up-to-date.

Start free 14-day trial



Zeta Alpha Al discovery platform

① Discover content using neural search, find similar, visualization, code popularity and social media.

② Organize your work in personalized topic tags

③ Receive timely need to know recommendations tailored to your interests

(4) Share and re-use knowledge with your team

🖾 Zeta Alpha



Al focused technical content from arXiv, conferences, companies, blogs, news, github code, twitter

+ import private data

Why Neural Search?



- Semantic understanding of data as opposed to surface keywords: bridge *the lexical gap*
- Context and relationships crucial in interpreting meaning: *handles complex and relational queries*
- Unstructured data accessible without classification and taxonomies, even *multi-lingual* and *cross-lingual*
- *Multi-modal* capabilities: potential to combine text, audio, images and video







expanding the horizon of exploration

A well tuned BM25 is hard to beat...





Dense Neural Retrieval, it's different...







BM25 struggles on more complex queries

≡	🖾 Zeta Alpha	low cost training BERT X Q	4 ⁹⁹⁺ 🖬 🗢 JZ
0	Discover	Any time 🗸 Sources 🗸 Code 🗸 Countries 🗸 Organizations 🗸 Owner 🗸	
•	Recommendations 🔹		
	People	260,382 results	Transformer Powered Search 👔 🖉 🔳 Relevance 🖛
*	Favorites	1	
۵	My documents	A statistic stat	Visualization () Explore more
≡	Notes	na karan kar	
~ III	Tags	Fine-tuning pre-trained language models improves the quality of commercial reply suggestion systems, but at the cost of unsustainable training times Popular training ti more	VOSviewer I < I I I I I I I I I I I I I I I I I
Q	My tags Following Shared	ar la characterization and a characterization and characterization and a characterization a	
\sim M	ly tags	Set and the set of	Large-Scale Differentially Pri How to Train BERT
1	approximate k-NN	\odot \odot \odot Ξ Find similar \equiv Notes \square Tag \Leftrightarrow \lt :	How to Train BERT with an Acad
8	Question Answering	II	Optimizing small BERTs trained
	long form qa 🚢	Networks	iBERT: Enhancing BERT with Li
#	recommender syste	13 Jan 2021 Max W. Y. Lam, Jun Wang, Dan Su & et al. (1)	ROSITA: Refined BERT cOmpreSsi
8	contrastive learning	recurrent (GALR) network On one hand, with only 1.5M parameters, it has ac more	den Costs of Low Qualit
8	ai index 2021	🛅 1 🖹 42 🖾 PDF Reader 💭 🛪 3	Cimple 4:4724 for distilling Documents 50 Similarity links: 560 Groups: 5
	mlops 🚉		
=	transformer new vari	\bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc Find similar \equiv Notes \square Tag \bigcirc \checkmark :	
11	expert search	II	What are the low cost training BERT?
=	citation analysis ≗	arXiv AT-BERT: Adversarial Training BERT for Acronym Identification Winning Solution for SDU@AAAI- 21	Layer Dropping and Layer Freezing
8	BERT	11 Jan 2021 Danqing Zhu, Wangli Lin, Yang Zhang & et al. (4)	
::	financial time series	In this paper, we present an Adversarial Training BERT method named AT-BERT, our winning solution to	

Dense retrieval "gets it"



≡ [🗟 Zeta Alpha	low cost training BERT	× ९ 🖬		. <mark>99+</mark> +	٥		JZ
Ø Disc	cover	Any time 👻 Sour	ces 🗸 Code 🖌 Countries 🗸 Organizations 🗸 Owner 🗸					
Receiption	ommendations 🌣	11,980 results		Transformer Powered Search 🚯 🔲	E	Releva	nce F	-
Peo	ple							
Favo	orites documents		TDS Training BERT at a University	Visualization 👔		Explore	more	
≡ Note	es		If you're reading this post, you've probably heard about the remarkable performance of new machine					-
🗸 📕 Tags	s		learning models like BERT, GPT-2/3, and other deep learning models for language, image, audio, and	NOSviewer		2 13	U	e.
Q My	tags Following Shared			Dic Modeling with BERT				
∨ My ta	gs <u>=</u> +			How to Train BERT				
∷ nei ∷api	ural retrieval 🚓 proximate k-NN 🚜		\bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \updownarrow Find similar \equiv Notes \square Tag \Rightarrow $<$:	How to Train BERT with an Acad	A Benchr	mark for	+	wit
∷ Qu	estion Answering			Fine-tuning BERT for Low-Resou	-		-	
∷ Ion	ng form qa 🏔	BERT	13 May 2021 Bansidhar Mangalwedhekar	Optimizing small BERTs trained			Ŷ	J.
∷ rec	commender syste	Rev 6 6 - 6	This paper studies the efficiency of transferring BERT learnings to low complexity models like BiLSTM,	Distilling BERT for low comple RØSITA: Refined BERT compreSsi				
# Cor	ntrastive learning	Bingle Dentarca	more	Iden Costs of Low Qualit				
II ai i	index 2021		聞 0 🖹 4 📖 PDF Reader	Cimple totatel for distilling		\bigcirc	- 4.9	
# <u>m</u> l	ops 🚉							Ł
:: tra	nert search							
:: ex;	ation analysis	Barris Sant Management		What are the low cost training BERT?		C	eta 🚹	2
II BE	RT .	 Source and the second se	15 Apr 2021 Peter Izsak, Moshe Berchansky & Omer Levy	Paper microscopes, public data repositories, and hosted noteboo	ks		0	6
∷ fin	ancial time series	 And Parkets (± 4.0), since a print 	While large language models a la BERT are used ubiquitously in NLP, pretraining them is considered a luxury that only a few well-funded industry labs can afford. How can one train such models with a m				Ç	

Zeta Alpha domain specific evaluation



We benchmark on ~100 queries in AI domain from 5 types: short phrases, knowledge graph questions, quora questions, freq. user queries, and paper titles.

Model	P@10	R@10	F1@10	MRR@10
ZA keyword retrieval	0.71	0.17	0.27	0.91
ZA dense retrieval	0.84	0.21	0.34	0.94

Challenges dense retrieval:

- **Calibration** (always gets an answer)
- Explainability
- q2doc vs doc2doc seem to require different embeddings.

Combining dense and keyword search



User does not ask for two separate retrieval mechanisms, but for relevant results.

Dense en keyword retrieval top-10 overlap only ~ 15%.

How can we combine the rankings?

Sum:

 $score(q, d) = sim(q, d)/sim_{max}(q, d) + BM25(q, d)/BM25_{max}(q, d)$

Mul:

```
score(q,d) = \sin(q,d) \times \mathrm{BM25}(q,d)
```

RRF: (k=60)

$$RBF(q,d,M) = \sum_{m \in M} \frac{1}{k + \pi^m(q,d)}$$

	P@10	MRR@10	nDCG@10
Keyword	0.806	0.930	0.858
Dense	0.754	0.873	0.863
Sum	0.830	0.947	0.883
Mul	0.846	0.940	0.881
RRF	0.833	0.957	0.877

Just released: PDF Reader







X

66

3

L(N, D).

Approach 1

Approach 2 - Approach 3 --- Kaplan et al (2020

Just released: PDF Reader + margin notes



training tokens and model parameters, we are interested in minimizing L under the constraint

🗛 🕦 🚯 31 💵 PDE Reader 🔰 🕅 🕵 287

V

A Chinchilla (708) gopher (280B) # GPT-3 (1758) Megatron-Turing NLG (530B) We have definitely entered the exaFLOPs era of 66 Chinchilla outperforms Gopher and the other In fact it also outperforms GPT-3

argmin

?? really ?? Documents: 50 \rightarrow



Just released: PDF Reader + search in and with notes



How about using a neural reranker?



Add monoT5 based cross-encoder (Pradeep, Nogueira and Lin, 2021) on top of BM25

First stage	Second stage	Prec@10	Recall@10	F1@10	MRR@10
keyword	-	0.69	0.15	0.24	0.87
keyword	rerank 30	0.71	0.15	0.25	0.91
keyword	rerank 100	0.53	0.15	0.19	0.88

Marginal improvements or even degradation of results...



We lose our boosting effects





Dense Retrieval + Reranker



Add monoT5 based cross-encoder (Pradeep, Nogueira and Lin, 2021) on top of ZA Dense (knn)

First stage	Second stage	Prec@10	Recall@10	F1@10	MRR@10
knn	-	0.83	0.19	0.31	0.93
knn	rerank 30	0.82	0.19	0.31	0.96
knn	rerank 50	0.82	0.19	0.31	0.97
knn	rerank 100	0.85	0.19	0.31	0.97

Marginal improvements, but at what cost?



Better top results





Cross encoder reranking: Ouch, Performance!

top N	api call	Direct model call 3.0s 3.1s 3.0s				
10	0.9s	3.0s				
30	1.9s	3.1s				
50	4.2s	3.0s				
100	4.7s	3.1s				
1000	44.2s	3.4s				



InPars: Data Augmentation for Unsupervised IR



Ranking models are finetuned on a synthetic dataset built by augmenting documents with queries using generative LLMs like GPT-3.

This represents our recipe for **unsupervised domain adaptation**. With very good results on the BEIR benchmark...



Unsupervised Dataset Generation for Information Retrieval

The Information Retrieval (IR) community has recently witnessed a revolution due to large pretrained transformer models. Another key ingredient for this revolution was the MS MARCO dataset, whose scale and diversity has enabled zero-shot transfer learning to various tasks. However, not all IR tasks and domains can benefit from one single dataset equally. Extensive research in various NLP tasks has shown that using domain-specific training data, as opposed to a general-purpose one, improves the performance of neural models [43, 54]. In this work, we harness the few-shot canabilities of large pretrained language models as synthetic data generators for IR tasks. We show that models finetuned solely on our unsupervised dataset outperform strong baselines such as BM25 as well as recently proposed self-supervised dense retrieval methods. Code, models, and data are available at <hidden url>.

CCS CONCEPTS

 Information systems → Novelty in information retrieval: Computing methodologies \rightarrow Neural networks.

KEYWORDS

Few-shot Models, Large Language Models, Generative Models, Question Generation, Synthetic Datasets, Multi-stage Ranking

ACM Reference Format

. 2018. Unsupervised Dataset Generation for Information Retrieval. In Wood-

billions of parameters. As of February 2022, they charge 0.06 USD per 1000 tokens for their largest model. If each candidate document contains 250 tokens, naively using this API for a reranking task would cost approximately 15 USD per query.

Dense retrievers [13, 14] avoid this expensive reranking step by precomputing vector representations for each document in the collection prior to retrieval. When a query comes in, only its vector representations are computed, and a fast vector search framework can be used to retrieve the nearest document vectors to the vector representation of the query [12]. Despite being computationally cheaper at inference time, dense retrievers need one inference pass to compute the vector representation of each document in the collection, which also makes billion-parameter neural models impracticable to be used as dense retrievers.¹ Another challenge in developing neural models for IR is the lack of domain-specific training data. Manually constructing high-quality datasets is difficult as it requires queries from real users. While there are a few general-purpose labeled data available [17, 28], they are not always effective in generalizing to out-of-domain datasets [26, 44]. For this goal, zero-shot and few-shot learning models are in particular promising. However, a cost-effective manner of using large LMs in IR tasks is still an open question.

In this work, we propose a simple yet effective approach towards efficiently using large LMs in retrieval and obtain improvements across several IR datasets. Rather than using large LMs directly in the retrieval process, we harness them to generate labeled data



InPars: Data Augmentation for IR using LLM's

		MARCO	TREC-DL 2020		Rol	Robust04		TRECC
		MRR@10	MAP	nDCG@10	MAP	nDCG@20	nDCG@10	nDCG@10
	Unsupervised							
(1)	BM25	0.1874	0.2876	0.4876	0.2531	0.4240	0.3290	0.6880
(2)	Contriever (Izacard et al., 2021)	-	-	-	-	-	0.2580	0.2740
(3)	cpt-text (Neelakantan et al., 2022)	0.2270	-	-	-	-	-	0.4270
	OpenAI Search reranking 100 docs from	n BM25						
(4)	Ada (300M)	\$	0.3141	0.5161	0.2691	0.4847	0.4092	0.6757
(5)	Curie (6B)	\$	0.3296	0.5422	0.2785	0.5053	0.4171	0.7251
(6)	Davinci (175B)	\$	0.3163	0.5366	0.2790	0.5103	\$	0.6918
	InPars (ours)							
(7)	monoT5-220M	0.2585	0.3599	0.5764	0.2490	0.4268	0.3354	0.6666
(8)	monoT5-3B	0.2967	0.4334	0.6612	0.3180	0.5181	0.5133	0.7835
	Supervised [> MARCO]							
(9)	Contriever (Izacard et al., 2021)	-	-	-	-	8	0.4980	0.5960
(10)	cpt-text (Neelakantan et al., 2022)	-	-	-	-	-	-	0.6490
(11)	ColBERT-v2 (Santhanam et al., 2021)	0.3970	-	-	-	8	0.5620	0.7380
(12)	GPL (Wang et al., 2021)	-	-	-	-	-	-	0.7400
(13)	miniLM reranker	$^{\dagger}0.3901$	-	-	-	-	$^{\ddagger}0.5330$	$^{\ddagger}0.7570$
(14)	monoT5-220M (Nogueira et al., 2020)	0.3810	0.4909	0.7141	0.3279	0.5298	0.5674	0.7775
(15)	monoT5-3B (Nogueira et al., 2020)	0.3980	0.5281	0.7508	0.3876	0.6091	0.6334	0.7948
	InPars (ours) $[\triangleright$ MARCO \triangleright unsup in-d	omain]						
(16)	monoT5-3B	0.3894	0.5087	0.7439	0.3967	0.6227	0.6297	0.8471



InPars: out of domain data augmentation

Example 1:

Document: We don't know a lot about the effects of caffeine during pregnancy on you and your baby. So it's best to limit the amount you get each day. If you are pregnant, limit caffeine to 200 milligrams each day. This is about the amount in 1½ 8-ounce cups of coffee or one 12-ounce cup of coffee.

Relevant Query: Is a little caffeine ok during pregnancy?

Example 2:

Document: Passiflora herbertiana. A rare passion fruit native to Australia. Fruits are green-skinned, white fleshed, with an unknown edible rating. Some sources list the fruit as edible, sweet and tasty, while others list the fruits as being bitter and inedible. **Relevant Query:** What fruit is native to Australia?

Example 3:

Document: The Canadian Armed Forces. 1 The first large-scale Canadian peacekeeping mission started in Egypt on November 24, 1956. 2 There are approximately 65,000 Regular Force and 25,000 reservist members in the Canadian military. 3 In Canada, August 9 is designated as National Peacekeepers' Day. **Relevant Ouery:** How large is the Canadian military?

Example 4: Document: {document_text} Relevant Query:

Example 1:

Document: We don't know a lot about the effects of caffeine during pregnancy on you and your baby. So it's best to limit the amount you get each day. If you are pregnant, limit caffeine to 200 milligrams each day. This is about the amount in 1½ 8-ounce cups of coffee or one 12-ounce cup of coffee.

Good Question: How much caffeine is ok for a pregnant woman to have?

Bad Question: Is a little caffeine ok during pregnancy?

Example 2:

Document: Passiflora herbertiana. A rare passion fruit native to Australia. Fruits are green-skinned, white fleshed, with an unknown edible rating. Some sources list the fruit as edible, sweet and tasty, while others list the fruits as being bitter and inedible. **Good Question:** What is Passiflora herbertiana (a rare passion fruit) and how does it taste like?

Bad Question: What fruit is native to Australia?

Example 3:

Document: The Canadian Armed Forces. 1 The first large-scale Canadian peacekeeping mission started in Egypt on November 24, 1956. 2 There are approximately 65,000 Regular Force and 25,000 reservist members in the Canadian military. 3 In Canada, August 9 is designated as National Peacekeepers' Day.

Good Question: Information on the Canadian Armed Forces size and history.

Bad Question: How large is the Canadian military?

Example 4: Document: {document_text} Good Question:



Figure 3: MRR@10 on the MS MARCO development set achieved by InPars using monoT5-220M reranker trained on synthetic questions generated by GPT-3 models of different sizes. Figures for cpt-text are from (Neelakantan et al., 2022). Note the log scale for the x-axis.

















Any Questions?

લ વ

Publications

InPars: Data Augmentation for Information Retrieval using Large Language Models
 2022 | Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee & Rodrigo Nogueira

🖾 Zeta Alpha

See paper

• Building a Platform for Ensemble-based Personalized Research Literature Recommendations for AI and Data Science at Zeta Alpha

2021 | Jakub Zavrel, Artem Grotov, & Jonathan Mitnik

See paper

 mMARCO: A Multilingual Version of the MS MARCO Passage Ranking Dataset
 2021 | Luiz Bonifacio, Vitor Jeronymo, Hugo Queiroz Abonizio, Israel Campiotti, Marzieh Fadaee, Roberto Lotufo & Rodrigo Nogueira

See paper

• Pretrained Transformers for Text Ranking: BERT and Beyond 2021 | Jimmy Lin, Rodrigo Noqueira, & Andrew Yates

Access book

A New Neural Search and Insights Platform for Navigating and Organizing AI Research
 2020 | Marzieh Fadaee, Olga Gureenkova, Fernando Rejon Barrera, Carsten Schnober,
 Wouter Weerkamp, Jakub Zavrel

See paper

• Effective Distributed Representations for Academic Expert Search 2020 | Mark Berger, Jakub Zavrel, & Paul Groth Get more information, sign up to use the platform:

www.zeta-alpha.com

Interested in our mission? Join the team!

See open positions