Bridging the Training-Inference Gap for Dense Phrase Retrieval



Gyuwan Kim¹ Jinhuk Lee² Barlas Oğuz³ Wenhan Xiong³ Yizhe Zhang³ Yashar Mehdad³ William Yang Wang¹ ¹University of California, Santa Barbara ²Korea University ³Meta AI



10⁹

Motivation

- Components to build a dense retrieval system
 - ➤ Training a dual encoder
 - > Selecting the best model with validation
 - Constructing an index for efficient search are loosely connected each other
 - e.g., model training does not directly optimize the retrieval performance from the full corpus
- Goal: minimize the training-inference gap of dense retrievers to achieve better retrieval

Efficient Validation

- To expedite modeling innovation correctly, we measure retrieval accuracy on an index from a smaller subset of the full corpus (C)
- C₀: gold passages from the development set (minimal set ensuring to contain answers)
- ✤ Random Subcorpus (R_r): C_0 + random passages, $|R_r| = r|C|$
- Hard Subcorpus (H_k): C₀ + all context passages from top-k retrieval results using a

performance (focusing on *phrase retrieval*)

(b) **Ours**

Checkpoint

Retriever

Index

n = 6M

Metric

 $\mathcal{L}_{ ext{train}}$

Pseudo Corpus 🔶

 $Q_{\rm dev}$

(a) Lee et al., 2021

Full Corpus ---- Retriever

 $Q_{\rm dev}$

 $= \mathcal{L}_{\text{train}}$

pre-trained dense retriever



Comparison of (a) original and (b) proposed procedure of DensePhrases training (top) and validation (bottom)

Index

n = 3B

Metric

Checkpoint

Validation results with different size of random ($r \in \{0, 1/100, 1/10\}$) and hard ($k \in \{1, 2, 4, 8, 10, 16, 32, 64\}$) subcorpora, \flat : before query-side fine-tuning

Optimized Training of DensePhrases

Unified loss (UL)

Training

Validation

Checkpoint

- We should find an answer phrase among all possible candidates at once in test time
- > Put all negatives together into contrastive targets with different λ coefficients
- Hard negatives (HN)
 - > Fix mistakes from the first round model
 - Mining: extract model-based hard negatives from top-k retrieval results for questions in the training set
- \succ Use all tokens in context passages
 - # of negatives: in-passage (L-1), in-batch
 - (B-1 \rightarrow B*L-1), pre-batch (B*T \rightarrow B*T*L)
- Training: fine-tune a dual encoder by appending sampled hard negatives as negative targets for each training step

Experiments

- The relative order of accuracy between models on hard subcorpus converges quickly
- Both UL and HN are shown to be effective

 We improves phrase retrieval by 2-3% in top-1 accuracy and passage retrieval by 2-4% in top-20 accuracy from DensePhrases

Model	NaturalQ	TriviaQA EM@1	Model	N ₂ translo					TrinicOA				
	EM@1			NaturalQ				IriviaQA					
$DPR^{\diamond} + BERT$ reader	41.5	56.8		Top-1	Top-5	Top-20	MRR@20	P@20	Top-1	Top-5	Top-20	MRR@20	P@20
DPR [♠] + BERT reader	41.5	56.8	DPR♦	46.0	68.1	79.8	55.7	16.5	54.4	-	79.4	_	-
$RePAO^{\diamondsuit}$ (retrieval-only)	41.2	38.8	DPR [•]	44.2	66.8	79.2	54.2	17.7	54.6	70.8	79.5	61.7	30.3
RePAQ ^(retrieval-only)	41.7	41.3	DensePhrases	50.1	69.5	79.8	58.7	20.5	-	-			-
DensePhrases [♡]	40.9	50.7	DensePhrases [•]	51.1	69.9	78.7	59.3	22.7	62.7	75.0	80.9	68.2	38.4
DensePhrases	41.3	53.5	DensePhrases [©] -III	57 1	757	837	65.2	22.0	62.0	74.6	80.6	67.6	333
DensePhrases [♥] -UL DensePhrases [♥] -UL-HN DensePhrases [♠] -UI	43.5 44.0 42.4	51.3 47.0	DensePhrases [♥] -UL-HN DensePhrases [♠] -UL	58.6 56.7	75.7 75.9	83.4 83.8	66.1 65.2	22.0 21.9 23.7	60.3 65.0	73.3 76.6	79.6 82.7	66.1 70.2	32.3 39.0
Densel mases -OL	ч∠.т	55.5											

◊: trained on each dataset independently, ★: trained on multiple datasets, ♡: trained on Natural Questions datasets