

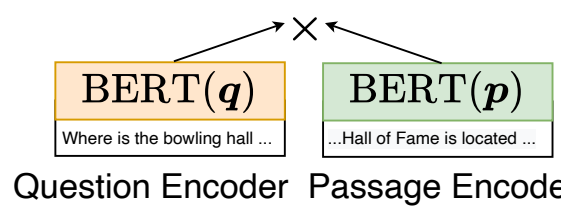
Questions Are All You Need to Train a Dense Passage Retriever

Devendra Singh Sachan^{1,2}, Mike Lewis³, Dani Yogatama⁴, Luke Zettlemoyer^{3,5}, Joelle Pineau^{1,2,3}, Manzil Zaheer⁶
¹McGill University, ²Mila-Quebec AI Institute, ³Meta AI, ⁴USC, ⁵University of Washington, ⁶Google DeepMind

sachande@mila.quebec

Introduction

Dense Passage Retriever (DPR)

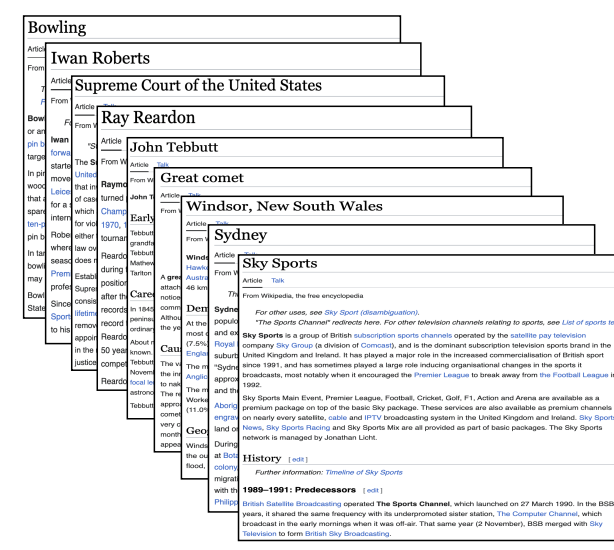


- Embeds question and passages in the same embedding space.
- Retrieved passages are very useful for tasks like *question answering*.

Our Approach to train DPR

Training Data =

When was the last time the detroit lions won a championship?
Who played general chang in star trek 6?
Where did the rule of 72 come from?



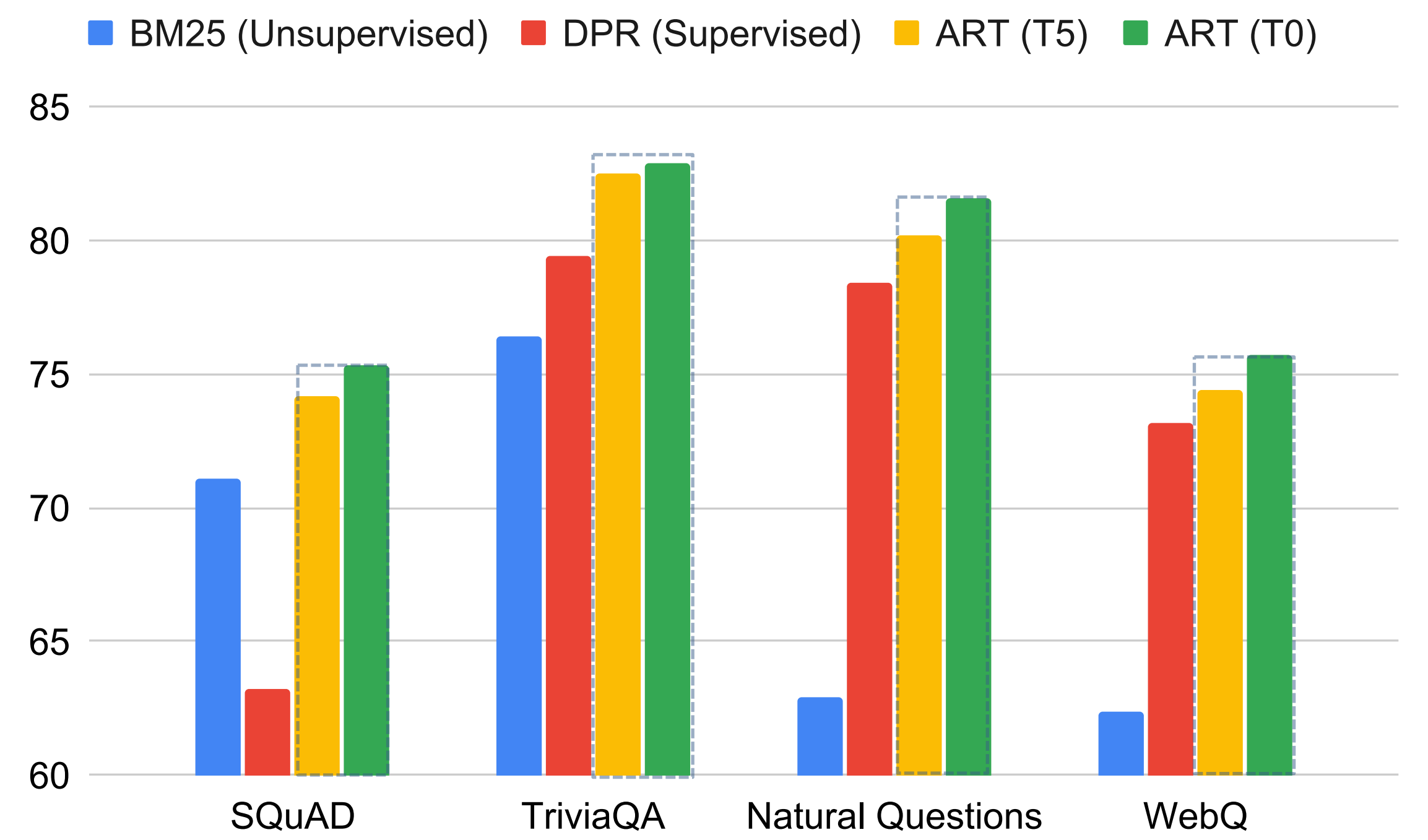
Collection of questions Database of documents (evidence)
Ex: English Wikipedia ~ 21M passages

Key Idea

Leverage knowledge contained in **Large Language Models** to score candidate passages.

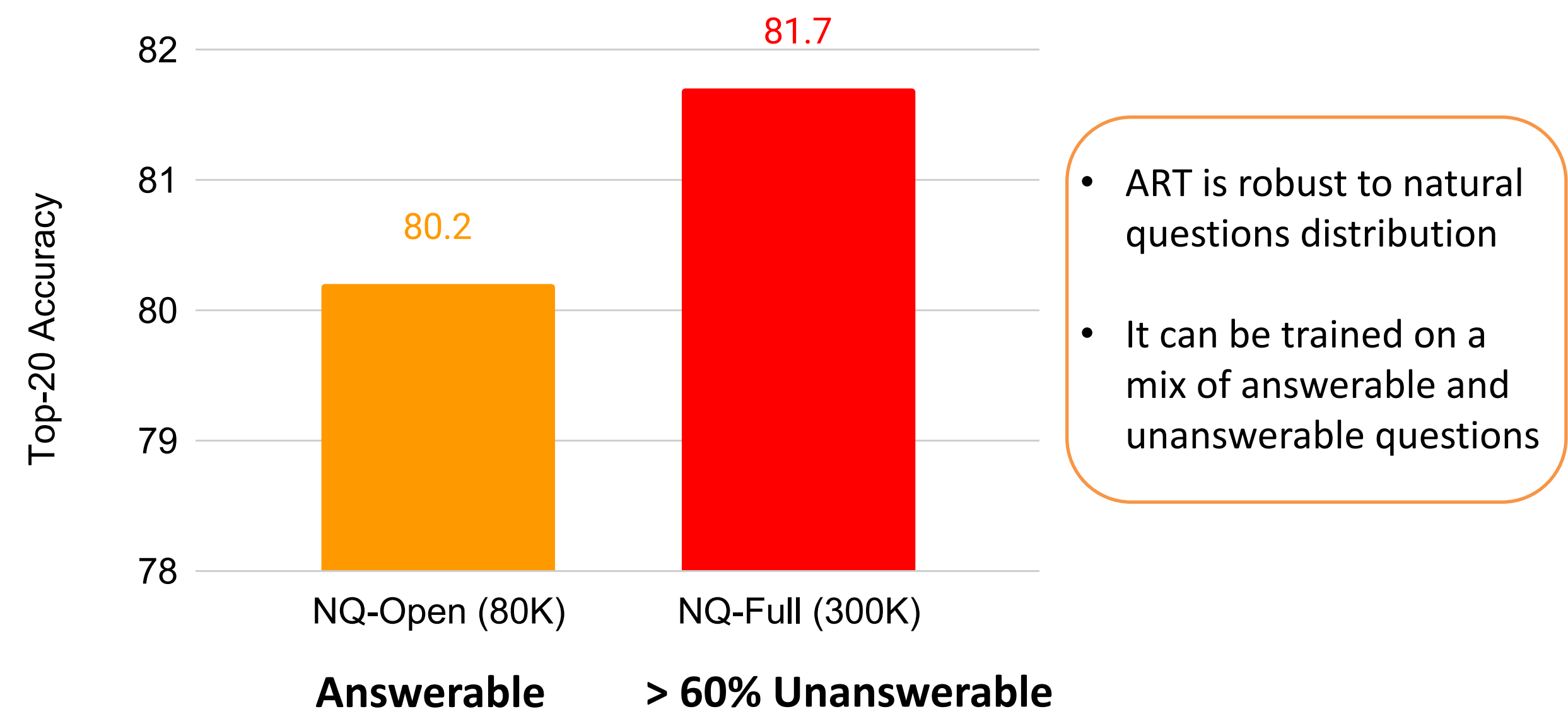
Results

Retrieval Accuracy (Top-20)



- ART outperforms unsupervised retrievers such as BM25.
- It matches or exceeds the performance of supervised models.

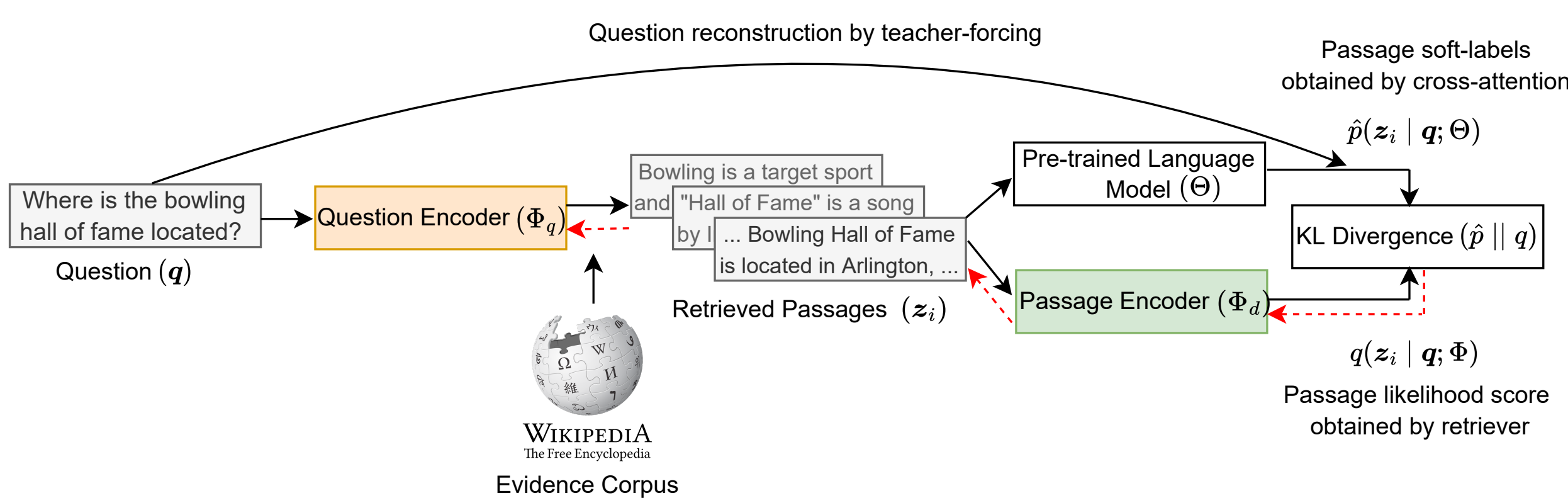
Robustness to Question Answerability



- ART is robust to natural questions distribution
- It can be trained on a mix of answerable and unanswerable questions

Method

ART: Autoencoding-based Retriever Training



Step 1: K-NN Search over Evidence

- Compute question similarity with all evidence passages
- Select top-K (such as 32) passages with highest scores $\mathcal{Z} = \{z_1, \dots, z_K\}$

Step 2: Retriever Likelihood Calculation

- Calculate scores using "current" passage encoder weights

$$s(\mathbf{q}, \mathbf{z}_i; \Phi) = f_q(\mathbf{q}; \Phi_q)^\top f_d(\mathbf{z}_i; \Phi_d)$$

- Define retriever distribution

$$q(\mathbf{z}_i | \mathbf{q}, \mathcal{Z}; \Phi) = \text{softmax } s(\mathbf{q}, \mathbf{z}_i), \forall \mathbf{z}_i$$

Step 3: Zero-Shot Relevance Score Estimation

- Use a LLM to score question tokens conditioned on each passage (teacher-forcing)

$$p(\mathbf{z}_i | \mathbf{q}; \Theta) \propto \frac{1}{|\mathbf{q}|} \sum_t \log p(q_t | \mathbf{q}_{<t}, \mathbf{z}_i; \Theta)$$

- Obtain soft relevance scores *i.e.*, teacher distribution as

$$\hat{p}(\mathbf{z}_i | \mathbf{q}, \mathcal{Z}) = \text{softmax } p(\mathbf{z}_i | \mathbf{q}; \Theta), \forall \mathbf{z}_i$$

Step 4: Loss Calculation and Backpropagation

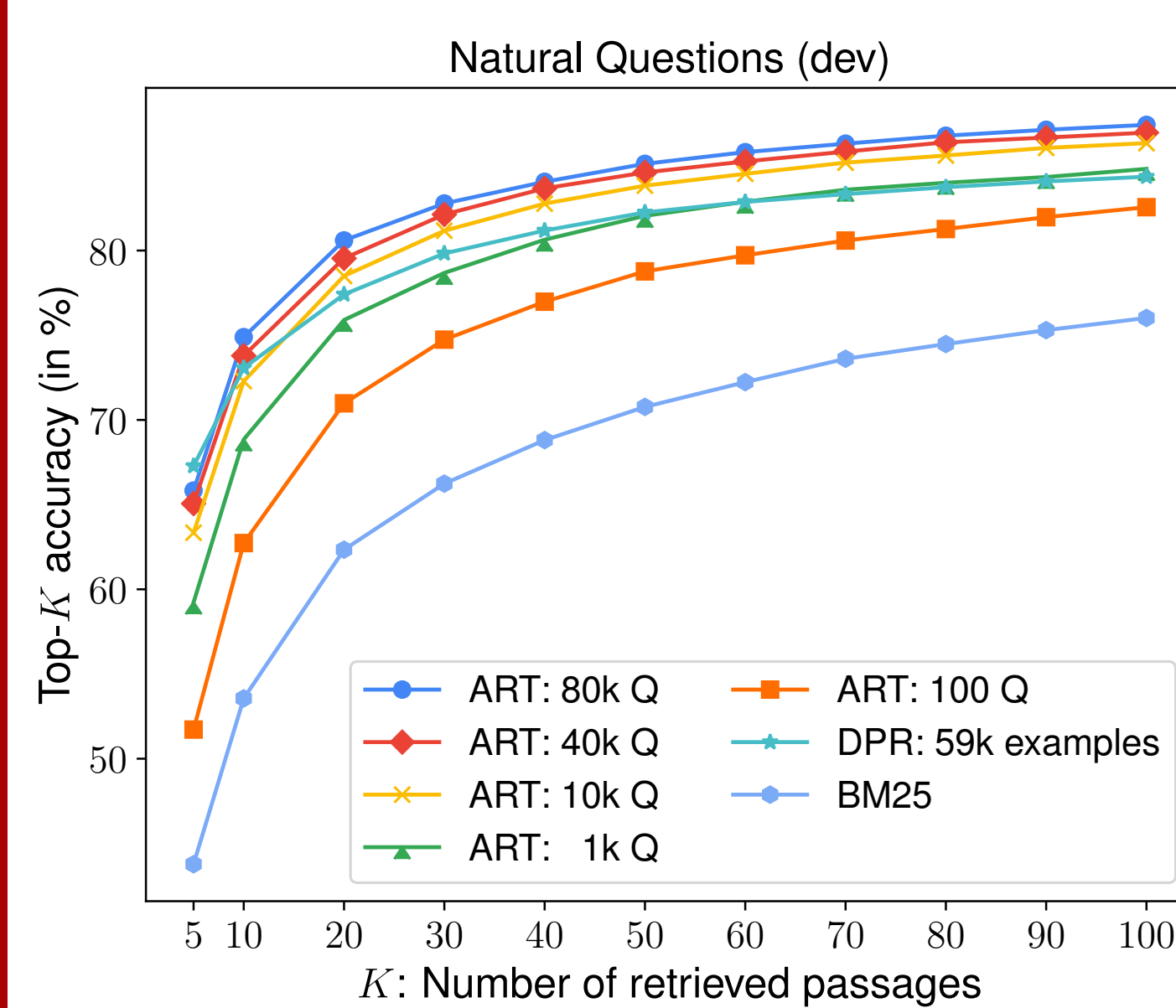
- Align retriever distribution with relevance score distribution

$$\mathcal{L}(\Phi) = \mathbb{KL}(\hat{p}(\mathbf{z}_i | \mathbf{q}, \mathcal{Z}) || \underbrace{q(\mathbf{z}_i | \mathbf{q}, \mathcal{Z}; \Phi)}_{\text{retriever training}})$$

Step 5: Periodically Update Evidence Embeddings

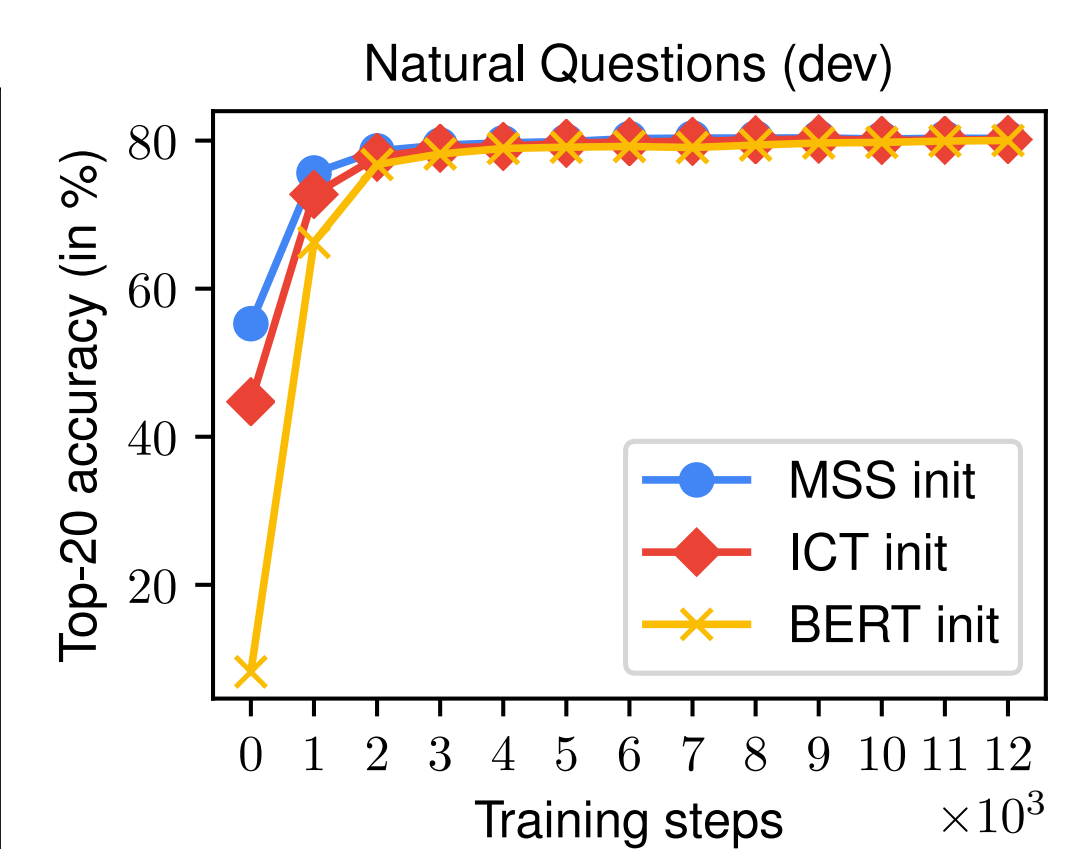
Analysis

Sample Efficiency



- Outperforms BM25 with just 100 examples

Retriever Initialization



- Different initializations converge to the same final state
- Useful for bootstrapping retriever from BERT

Conclusion

- **ART**: an approach to train dense retriever using unaligned pairs of questions and passages.
- Custom hard negative mining approaches are **not required**.
- Uses off-the-shelf large language models as a black-box (w/o finetuning).

Related Work

1. *Improving Passage Retrieval with Zero-Shot Question Generation*, Devendra Sachan et al., EMNLP 2022
2. *RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking*, Ruiyang Ren et al., EMNLP 2021