

# REAL: Retrieval-Augmented Language MOdel Pre-Training

Shubham Gupta

February 20, 2020

## 1 Introduction

- REALM is a paper mentioned in the T5 paper titled: **How Much Knowledge Can You Pack Into The Parameters of a Language Model?**
- TLDR: This paper retrieves documents that have the information present while solving Question-Answer type problems.
- Introduced a latent *knowledge retriever*, which can attend and retrieve documents over large corpus and can be trained in unsupervised manner using masked language modelling technique and backprop through retriever which considers lots of docs.

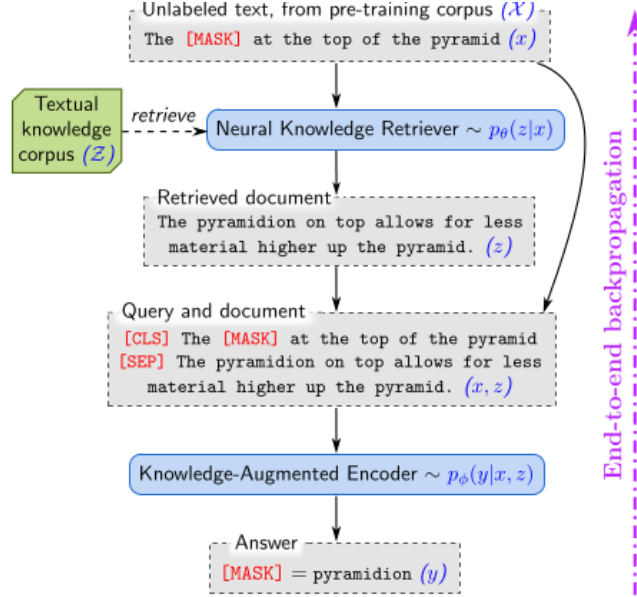


Figure 1. REALM augments language model pre-training with a **neural knowledge retriever** that retrieves knowledge from a **textual knowledge corpus**,  $Z$  (e.g., all of Wikipedia). Signal from the language modeling objective backpropagates all the way through the retriever, which must consider millions of documents in  $Z$ —a significant computational challenge that we address.

Figure 1: Training process for REALM

- Key point: Train retriever using a performance-based signal from unsupervised text.
- Retrieval based LM = More computational resources = More money
  - Solution: Computation performed for each doc is cached and can be used again. Best doc selected using *Maximum Inner Product Search (MIPS)*. Read the paper here.
- REALM retriever can be used on downstream tasks via transfer learning.
- REALM is SOTA on NQ-Open, WQ and CuratedTrec.

## 2 Approach

### 2.1 Retrieve-then-predict generative process

- Training: Masked-LM. Fine-tuning: Open QA task
- $p(y|x)$  decomposed into two steps:

- Given  $x$ , retrieve documents  $z$  from corpus  $Z$ . Modelled as:  $p(z|x)$
- Condition of both  $z$  and  $x$  to generate output  $y$  i.e  $p(y|z, x)$
- Overall likelihood  $y$  is generated by treating  $z$  as latent variable and marginalizing over all documents  $z$

$$p(y|x) = \sum_{z \in Z} p(y|z, x) * p(z|x) \quad (1)$$

## 2.2 Architecture

- **Neural Knowledge Retriever** i.e models  $p(z|x)$
- **Knowledge Augmented Encoder** i.e models  $p(y|z, x)$

## 2.3 Neural Knowledge Retriever

- Dense inner product model.

$$p(z|x) = \frac{\exp(f(x, z))}{\sum_{z'} \exp(f(x, z'))} \quad (2)$$

$$f(x, z) = \text{Embed}_{input}(x)^T \text{Embed}_{doc}(z)$$

- $\text{Embed}_{input}$  and  $\text{Embed}_{doc}$  are embedding functions
- $f(x, z)$  is called **relevance score**. It is inner product of vector embeddings.
- Relevant Distribution is softmax over all relevance scores
- Embedding implement using BERT-style transformers. Join using  $\text{SEP}_i$ , prefix using  $\text{CLS}_i$  and append  $\text{SEP}_i$  as the end.

$$\text{join}_{BERT}(x) = [\text{CLS}]x[\text{SEP}] \quad (3)$$

$$\text{join}_{BERT}(x_1, x_2) = [\text{CLS}]x_1[\text{SEP}]x_2[\text{SEP}]$$

- Pass above into transformer, which gives over vector for each token. Perform linear projection to reduce dimensionality of vector

$$\begin{aligned} \text{Embed}_{input}(x) &= W_{input} \text{BERT}_{CLS}(\text{join}_{BERT}(x)) \\ \text{Embed}_{doc}(z) &= W_{doc} \text{BERT}_{CLS}(\text{join}_{BERT}(z_{title}, z_{body})) \end{aligned} \quad (4)$$

## 2.4 Knowledge-Augmented Encoder

- Given input  $x$  and relevant doc  $z$ , this defines  $p(y|z, x)$
- Join  $x$  and  $z$  into single sequence and feed into transformer
- Here, training is different for pre-training vs fine-tuning
  - For pre-training, predict [MASK] token. Use same Masked LM (MLM) loss as in Transformer (Devlin)
  - For Open-QA, we need to produce string  $y$ . Assumption:  $y$  occurs as sequence of tokens in some document in the corpus. Skipping the math bit for now.

## 2.5 Training

- Compute gradients in  $\theta$  and  $\phi$  and optimize using SGD.
- Challenge: Computing  $p(y|x)$
- Approx by summing over top  $k$  documents with highest prob under  $p(z|x)$
- Question: How to find top  $k$  docs? Answer: Use MIPS
- Need to precompute  $Embed_{doc}(x)$  for all docs. Problems? It changes with each step of SGD.
- *Solution*: Async refresh  $Embed_{doc}$  every 500 steps
- Use MIPS to select top  $k$  docs. For these docs, recompute  $p(z|x)$  using new  $\theta$ .

### 2.5.1 Implementing async MIPS refreshes

- Two jobs running in parallel:
  - *Primary trainer*: Perform gradient updates on parameters
  - *Secondary index builder*: Embeds and indexes the docs

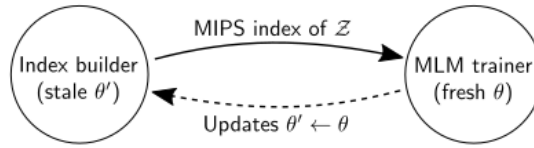


Figure 3. REALM pre-training with asynchronous MIPS refreshes.

Figure 2: Async MIPS implementation

- Async refresh used only for pre-training
- For fine tuning, build index once from pre-trained  $\theta$  and use it.

### 2.5.2 What does retriever learn?

- Retriever promotes docs that improve accuracy
- This can be analyzed by analyzing gradient wrt the parameters

### 2.5.3 Injecting inductive biases into pre-training

- **Salient span masking**: Some questions require only local context. Select named entities and dates and mask one of them. Performs better.
- **Null document**: Add null document to top  $k$  documents to allow answers even when no context is required

- **Prohibiting trivial retrievals:** If knowledge corpus  $Z$  is the same as pre-training corpus  $X$ , it can predict  $y$  by looking at  $x$  in  $z$ . Exclude trivial candidate
- **Initialization:** Warm up  $Embed_{input}$  and  $Embed_{doc}$  using Inverse Cloze Task(ICT) i.e model trained to retrieve the doc where the sentence came from.

### 3 Experiments

- REALM outperforms all approaches by a big margin.
- Not adding results here for now.

### 4 Future Work

- Structured knowledge where we learn entities which are informative
- Multi lingual setting. Retrieving knowledge in high resource language to better represent text in low resource language
- Multi model setting. Retrieve images or videos that can provide knowledge not present in text