

Reading Wikipedia to Answer Open-Domain Questions

Authors - Danqi Chen

Open-domain QA

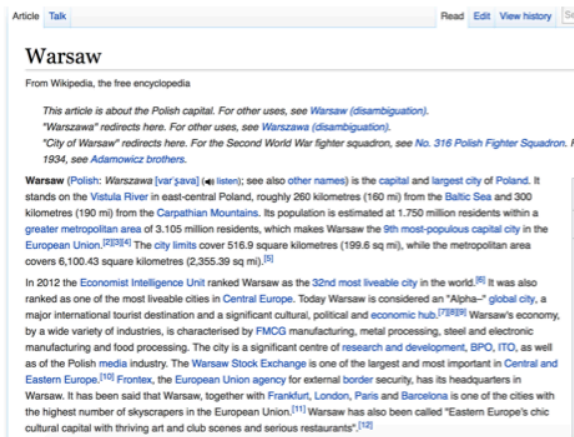
SQuAD, TREC, WebQuestions, WikiMovies

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



WIKIPEDIA
The Free Encyclopedia

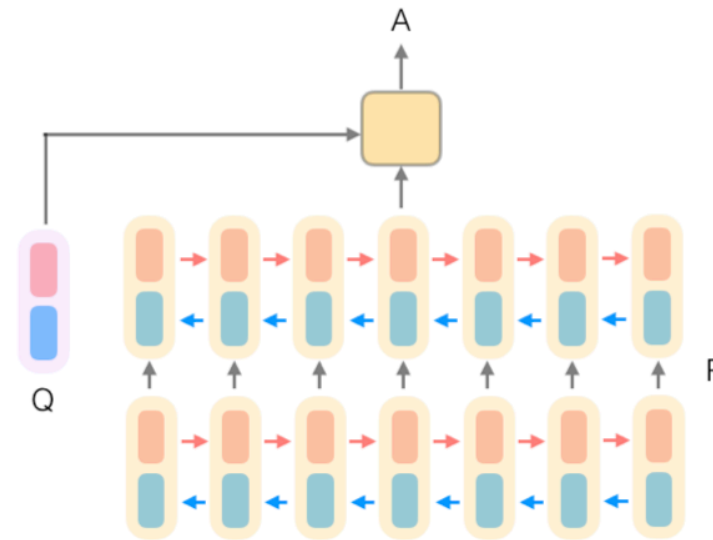
Document
Retriever



Document
Reader



833,500



Introduction

- Answering factoid questions in an open-domain setting
- Using Wikipedia as the unique knowledge source

Document Retriever

- Articles and questions are compared as TF-IDF weighted bag-of-word vectors. Additionally uses bigram counts for retrieval.
- Return 5 Wikipedia articles given any question

Document Reader

- A question q of l tokens , a paragraph p of m tokens
- Paragraph encoding
- Question encoding
- Prediction

Paragraph encoding

$$\{p_1, \dots, p_m\} = \text{RNN}(\{\tilde{p}_1, \dots, \tilde{p}_m\}), \quad \tilde{p}_i \in \mathbb{R}^d$$

$$f_{emb}(p_i) = E(p_i) \quad \text{300-dimensional}$$

$$f_{exact_match}(p_i) = \mathbb{I}(p_i \in q) \quad \text{3-dimensional} \quad \text{original, lowercase or lemma form}$$

$$f_{token}(p_i) = (POS(p_i), NER(p_i), TF(p_i)) \quad \text{term frequency (TF)}$$

$$f_{align}(p_i) = \sum_j a_{i,j} E(q_j) \quad a_{i,j} = \frac{\exp(\alpha(E(p_i)) \cdot \alpha(E(q_j)))}{\sum_{j'} \exp(\alpha(E(p_i)) \cdot \alpha(E(q_{j'})))}$$

$\alpha(\cdot)$ is single dense layer with ReLU nonlinearity. $a_{i,j}$ captures the similarity between p_i and each question words q_j

Question encoding

- Recurrent layer on top of word embeddings of questions words.
- Attention.

Prediction

- predict the two ends of the span that is most likely the correct answer
- input : paragraph vectors $\{p_1, \dots, p_m\}$ and the question vector \mathbf{q}
- two classifiers
$$P_{start}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_s \mathbf{q})$$
$$P_{end}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_e \mathbf{q})$$

the best span from token i to token i' $\max P_{start}(i) \times P_{end}(i') \quad i \leq i' \leq i + 15$

Dataset	Example	Article / Paragraph
SQuAD	<p>Q: How many provinces did the Ottoman empire contain in the 17th century?</p> <p>A: 32</p>	<p>Article: Ottoman Empire</p> <p>Paragraph: ... At the beginning of the 17th century the empire contained 32 provinces and numerous vassal states. Some of these were later absorbed into the Ottoman Empire, while others were granted various types of autonomy during the course of centuries.</p>
CuratedTREC	<p>Q: What U.S. state's motto is "Live free or Die"?</p> <p>A: New Hampshire</p>	<p>Article: Live Free or Die</p> <p>Paragraph: "Live Free or Die" is the official motto of the U.S. state of New Hampshire, adopted by the state in 1945. It is possibly the best-known of all state mottos, partly because it conveys an assertive independence historically found in American political philosophy and partly because of its contrast to the milder sentiments found in other state mottos.</p>
WebQuestions	<p>Q: What part of the atom did Chadwick discover?[†]</p> <p>A: neutron</p>	<p>Article: Atom</p> <p>Paragraph: ... The atomic mass of these isotopes varied by integer amounts, called the whole number rule. The explanation for these different isotopes awaited the discovery of the neutron, an uncharged particle with a mass similar to the proton, by the physicist James Chadwick in 1932. ...</p>
WikiMovies	<p>Q: Who wrote the film Gigli?</p> <p>A: Martin Brest</p>	<p>Article: Gigli</p> <p>Paragraph: Gigli is a 2003 American romantic comedy film written and directed by Martin Brest and starring Ben Affleck, Jennifer Lopez, Justin Bartha, Al Pacino, Christopher Walken, and Lainie Kazan.</p>

Wikipedia as knowledge source, curated Trec , webquestions and wiki movies doesn't contain Training paragraphs, so distant supervision is used to create training data.

Dataset	Train		Test
	Plain	DS	
SQuAD	87,599	71,231	10,570 [†]
CuratedTREC	1,486*	3,464	694
WebQuestions	3,778*	4,602	2,032
WikiMovies	96,185*	36,301	9,952

Dataset	Wiki Search	Doc. Retriever	
		plain	+bigrams
SQuAD	62.7	76.1	77.8
CuratedTREC	81.0	85.2	86.0
WebQuestions	73.7	75.5	74.4
WikiMovies	61.7	54.4	70.3

Method	Dev		Test	
	EM	F1	EM	F1
Dynamic Coattention Networks (Xiong et al., 2016)	65.4	75.6	66.2	75.9
Multi-Perspective Matching (Wang et al., 2016) [†]	66.1	75.8	65.5	75.1
BiDAF (Seo et al., 2016)	67.7	77.3	68.0	77.3
R-net [†]	n/a	n/a	71.3	79.7
DrQA (Our model, Document Reader Only)	69.5	78.8	70.0	79.0

Dataset	YodaQA	DrQA		
		SQuAD	+Fine-tune (DS)	+Multitask (DS)
SQuAD (<i>All Wikipedia</i>)	n/a	27.1	28.4	29.8
CuratedTREC	31.3	19.7	25.7	25.4
WebQuestions	39.8	11.8	19.5	20.7
WikiMovies	n/a	24.5	34.3	36.5

REALM: Retrieval-Augmented Language Model Pre-Training

Kelvin Guu * Kenton Lee *

Motivation

- Pre-trained model in like BERT and T5 contains large amount of world knowledge implicitly in their network parameters.
- Larger models for storing more world knowledge.
- Capture knowledge in more interpretable and modular way

Background

- Language model pre-training - Bert (Masked LM)
- Open domain question - answering
 - Retrieve top k document and predict the answer from them.

Approach

- For both pre-training and fine-training, REALM learns- $P(y|x)$
 - For pre-training , x is masked sentence , y is missing token
 - For fine-tuning , task -OpenQA, x is question- y is answer

$$p(y|x) = \sum_{z \in \mathcal{Z}} p(y|z, x) p(z|x).$$

- \mathcal{Z} - helpful documents

Knowledge Retriever

- Learn distribution of documents again question.

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

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

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

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

$$\text{Embed}_{\text{input}}(x) = \mathbf{W}_{\text{input}} \text{BERT}_{\text{CLS}}(\text{join}_{\text{BERT}}(x))$$

$$\text{Embed}_{\text{doc}}(z) = \mathbf{W}_{\text{doc}} \text{BERT}_{\text{CLS}}(\text{join}_{\text{BERT}}(z_{\text{title}}, z_{\text{body}}))$$

Z_{title} - document's title

Z_{body} - document's body

Knowledge Augmented Encoder

- Given input x , retrieved document z .
 - KAE defines $p(y|z,x)$
 - x and z are joined into single sentence and feed into transformer.
- Different architectures for pre-training and fine-tuning.

Pre-training

- Masked Language model
 - Predict original value of missing token in x .

$$p(y | z, x) = \prod_{j=1}^{J_x} p(y_j | z, x)$$

$$p(y_j | z, x) \propto \exp(w_j^\top \text{BERT}_{\text{MASK}(j)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})))$$

- J_x is the total number of [MASK] tokens in x ,
- W are learnable parameters.

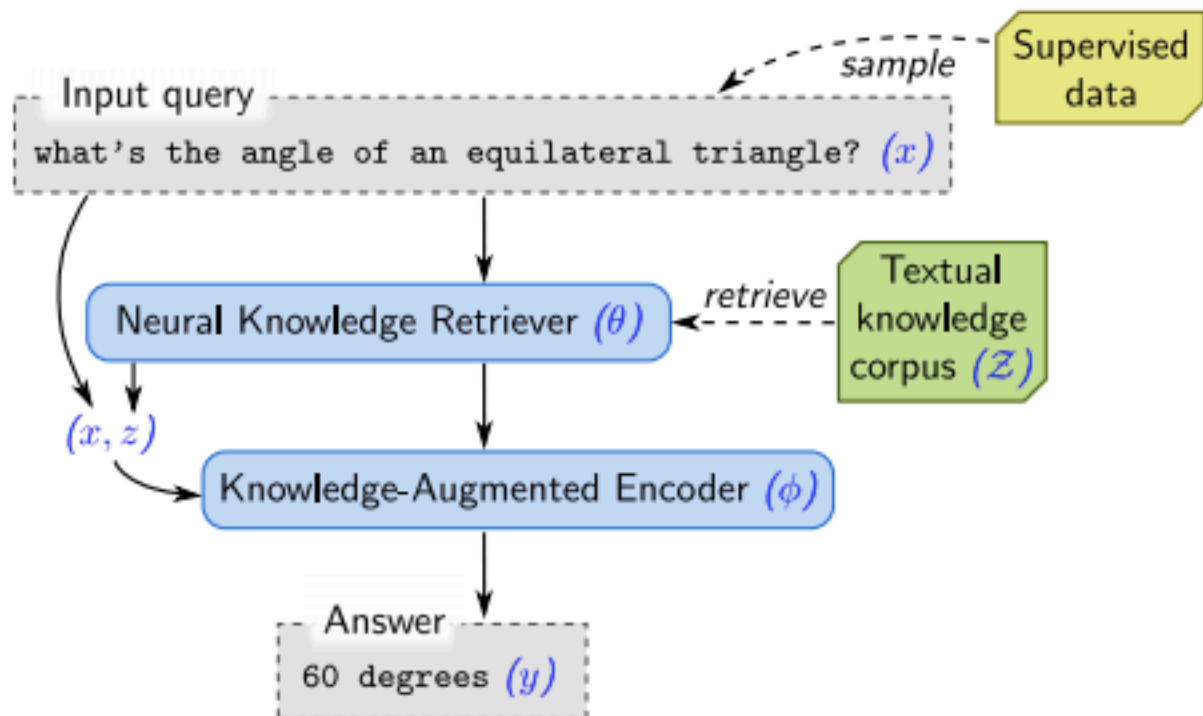
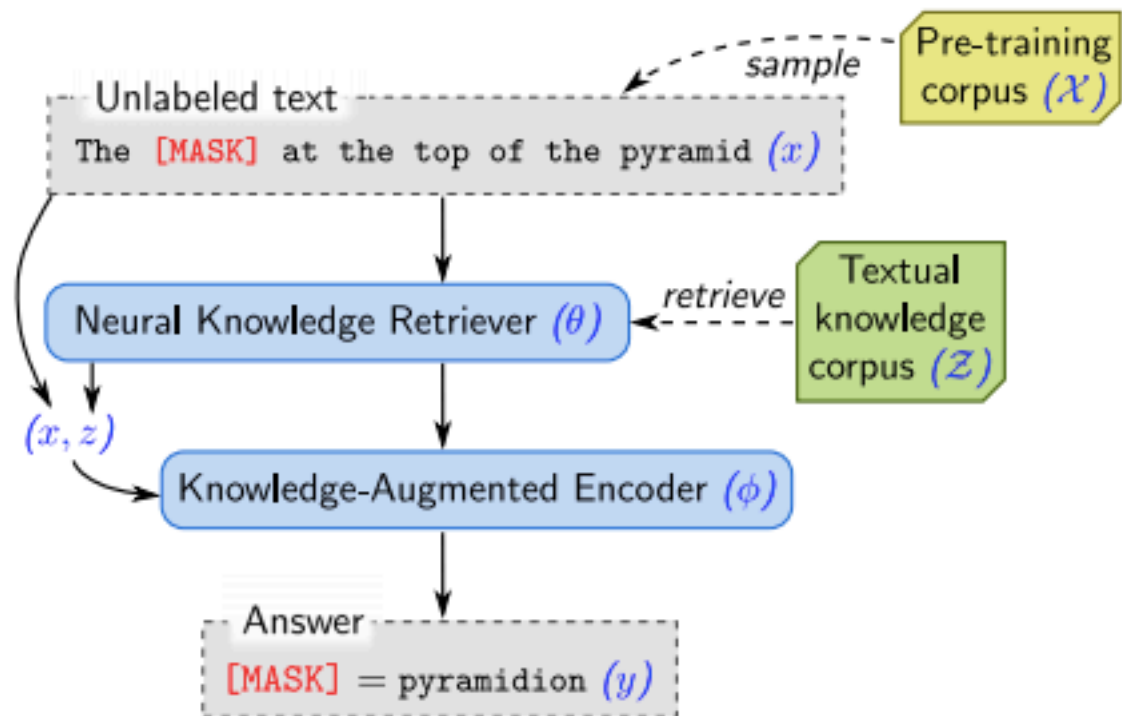
Fine - tuning

- Task - OpenQA
 - y is answer string.
 - Assumption - y is contiguous sequence of tokens in z
 - $S(z, y)$ be the set of spans matching y in z .

$$p(y | z, x) \propto \sum_{s \in S(z, y)} \exp(\text{MLP}([h_{\text{START}(s)}; h_{\text{END}(s)}]))$$

$$h_{\text{START}(s)} = \text{BERT}_{\text{START}(s)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})),$$

$$h_{\text{END}(s)} = \text{BERT}_{\text{END}(s)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})),$$



Training

- By maximizing the log-likelihood $\log p(y|x)$ of the correct output wrt to model parameters.
- Key challenge-
 - Marginal probability - $p(y | x) = \sum_{z \in \mathcal{Z}} p(y | x, z) p(z | x)$
 - Involves summation over all documents in knowledge source.
 - Approximate this by selecting top k under highest $p(z|x)$.
 - Reasonable since most of documents will have 0 probability.

Training

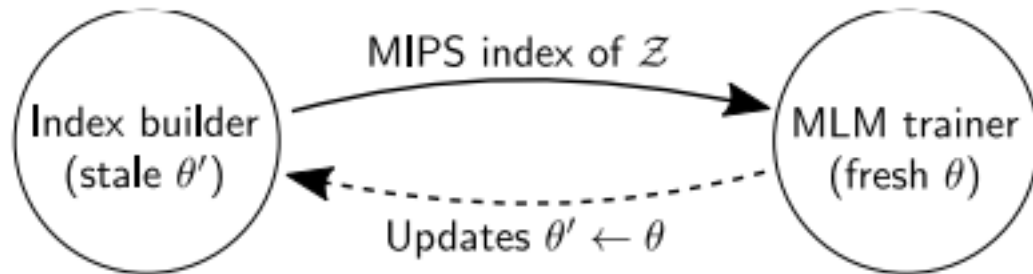
- $p(z|x)$ is equal to $f(x,z)$.

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

- Employ maximum inner product search to find approx top k documents.
- Need to precompute $\text{Embed}_{\text{doc}}(z)$ for every document.
- Construct an efficient search index.

Training

- But this index will become stale after update in parameters.
 - It only used to compute top-k documents.
 - Assuming no drastic change in parameters, index will slightly be stale.
 - Update the index asynchronously and train the MLM model.



Training

- MIPS index is refreshed after every few hundred training epochs for pre-training.
- Fine-tuning: index is built once and parameters of $Embed_{doc}(z)$ are not re-trained. $Embed_{input}$ is still fine-tuned to update retrieval function from query side.

What does retriever learn?

- Gradient of knowledge retriever wrt to parameters -

$$\nabla \log p(\mathbf{y} | \mathbf{x}) = \sum_{z \in \mathcal{Z}} r(z) \nabla f(\mathbf{x}, z)$$
$$r(z) = \left[\frac{p(\mathbf{y} | z, \mathbf{x})}{p(\mathbf{y} | \mathbf{x})} - 1 \right] p(z | \mathbf{x}).$$

- $p(\mathbf{y} | z, \mathbf{x})$ - probability of predicting the correct output \mathbf{y} given z .
- $p(\mathbf{y} | \mathbf{x})$ - is the expected value of $p(\mathbf{y} | \mathbf{x}, z)$.

Training strategies -

- Salient span masking
 - Some tokens only requires only local context.
 - Mask tokens which requires world knowledge.
 - “United Kingdom” or “July 1969”.
 - Identify such entities using NER and dates to mask them during pre training.

Training strategies

- Null document
 - Add empty document at top of k retrieved document.
 - This allows for cases where no-retrieval is necessary.
- Prohibiting trivial retrievals during pre-training -
 - If pre-training corpus and knowledge source are same,
 - KAE can trivially predict y by looking at unmasked version of x in z (which contains x).
 - This might result in KAE looking for string matches of x .
 - Remove such documents z during training.

Training strategies

- Initialization -
 - If not initialised,
 - Retriever doesn't retrieve relevant documents.
 - KAE starts ignoring documents by retriever.
 - Retriever will not receive any meaningful gradients.
 - Retriever can't improve
 - Vicious cycle.

Training strategies

- Initialization
 - Train the retriever using inverse cloze task.
 - Given a sentence, figure out from which document it came from.
 - Warm-start KAE using pre-trained BERT.

Experiments

- Open QA datasets -
 - Focus on datasets where authors didn't know the answer.
 - Avoid issues when questions is formulated with answer in mind.
 - Natural questions-Open(NQ) - google queries and their answers.
 - WebQuestions (WQ) - google suggest api and their answer from amazon mechanical turk.
 - Curated Trec (CT)- collection of question answer pair from sites like MSNSearch and AskJeeves.

Experiments

- Approaches compared -
 - Retrieval based OpenQA - like DrQA
 - Generation based OpenQA -
 - Text-to-text, encode question and predict answer token by token.
 - fine-tuned T5 for openQA
- Pre-training-
 - 200k steps on 64 TPUs, batch size 512, lr $3e-5$ and Bert's default optimizer.
- For each candidate , retrieve 8 candidate using MIPS including null document.

Results

Name	Architectures	Pre-training	NQ (79k/4k)	WQ (3k/2k)	CT (1k /1k)	# params
BERT-Baseline (Lee et al., 2019)	Sparse Retr.+Transformer	BERT	26.5	17.7	21.3	110m
T5 (base) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	27.0	29.1	-	223m
T5 (large) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	29.8	32.2	-	738m
T5 (11b) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	34.5	37.4	-	11318m
DrQA (Chen et al., 2017)	Sparse Retr.+DocReader	N/A	-	20.7	25.7	34m
HardEM (Min et al., 2019a)	Sparse Retr.+Transformer	BERT	28.1	-	-	110m
GraphRetriever (Min et al., 2019b)	GraphRetriever+Transformer	BERT	31.8	31.6	-	110m
PathRetriever (Asai et al., 2019)	PathRetriever+Transformer	MLM	32.6	-	-	110m
ORQA (Lee et al., 2019)	Dense Retr.+Transformer	ICT+BERT	33.3	36.4	30.1	330m
Ours (\mathcal{X} = Wikipedia, \mathcal{Z} = Wikipedia)	Dense Retr.+Transformer	REALM	39.2	40.2	46.8	330m
Ours (\mathcal{X} = CC-News, \mathcal{Z} = Wikipedia)	Dense Retr.+Transformer	REALM	40.4	40.7	42.9	330m

Results

- REALM outperforms T5 when approx 30 times lower in size.
- T5 has access to Squad data during pre-training.

Ablation	Exact Match	Zero-shot Retrieval Recall@5
REALM	38.2	38.5
REALM retriever+Baseline encoder	37.4	38.5
Baseline retriever+REALM encoder	35.3	13.9
Baseline (ORQA)	31.3	13.9
REALM with random uniform masks	32.3	24.2
REALM with random span masks	35.3	26.1
30× stale MIPS	28.7	15.1

Reviews (Pros)

- Thorough comparisons, experiments, training strategy,(Atishya, Jigyasa,Rajas, Lovish,Vipul)
- Dot product to retrieve documents, this allow for use of MIPS(Soumya,
- Improve SOTA(Soumya, Rajas,Saransh,Makkunda)
- Pre-training in retrieval phase(keshav)
- Provide context to language model(pawan)
- Explainability (Saransh, Siddhant,Pratyush)
- Ability to adapt to new knowledge(Siddhant)
- Greener alternative to T5(Vipul)
- Modular approach(Pratyush,Vipul)

Reviews (Cons)

- Lot of hyper-parameters (Atishya)
- Answer to be continuous span of keywords (Atishya, Siddhant, saransh)
- Doesn't allow multi-hop reasoning (Soumya, Rajas, Jigyasa, Siddhant, saransh)
- Conflicting information during retrieval due to time (Rajas)
- Oversell their paper (Keshav)
- Pre-training before pre-training (Lovish)
- Not actually explainable (Pawan)
- Started with issues with Bert and used BERT in the end (Pawan, Makkunda)
- Document embedding is fixed but input embedding is allowed to train during fine-tuning resulting these embeddings might go into different spaces. (Vipul)

Reviews(Extension)

- Using attention, copy mechanism to copy certain entities from retrieved documents - not vocab dependent and no need of answer span to be continuous (Atishya,siddhant)
 - How will you define $P(y/z,x)$ (Lovish)
- Retrieve Subgraph of big KB to augment sentence generation(Atishya)
 - Combining text is better then graphs, graph2text? (Soumya)
- Concat top k retrieved document to allow for multi-hop answering.(Soumya)
 - Extend current SOTA for multi-hop answering with current paper(Keshav)
 - May exceed Bert's capacity(Rajas)
 - Extract top N sentences/paragraph instead of documents(Saransh)
 - Multiple - retrieve-and-rank framework, in second retrieve only select from top documents selected in 1st step (Pratyush)
 - Retrieval Multiple times for multi-hop answering. Append the answer of 1st hop to retrieve relevant documents for 2nd hop(Makkunda)

Reviews(Extensions)

- Separate pre-training and fine-tuning to make system actually modular(Soumya)
 - Use openIE triplets, construct a graph and then use of GNNs to predict missing nodes for pre-training. Similarly GNN can operate on retrieved graph for fine-tuning(keshav)
 - Use of GNNs is moving away from the focus which is knowledge learning by adapting pre-training in language models. How do we incorporate multi-hop answering in pre-training(Saransh)
 - We should focus on building end - to -end pipelines for Graphs similar to current task.(Vipul)
- Instead of using Bert like architecture for retrieve/rank, how to extract knowledge from its pre-trained parameters (Pratyush)

Reviews(Extensions)

- Add time component to documents/questions to counter conflicting answers after updating knowledge source(Rajas)
- Multiple hstart/hend over multiple documents for multi-hop answering(Jigyasa)

Thanks !!!

