

# Empowering Language Models with Knowledge Graph Reasoning for Open-Domain Question Answering

Anonymous EMNLP submission

## Abstract

Answering open-domain questions requires world knowledge about in-context entities. As pre-trained Language Models (LMs) lack the power to store required knowledge, external knowledge sources, such as knowledge graphs, are often used to augment LMs. In this work, we propose knOwledge REasOning empowered Language Model (OREOLM), which consists of a novel Knowledge Interaction Layer that can be flexibly plugged into existing Transformer-based LMs to collaboratively interact with a differentiable Knowledge Graph Reasoning module. In this way, LM guides KG to walk towards the desired answer, while the retrieved knowledge improves LM. By adopting OREOLM to RoBERTa and T5, we show significant performance gain, achieving state-of-art results in the *Closed-Book* setting. The performance enhancement is mainly from the KG reasoning’s capacity to infer missing relational facts. In addition, OREOLM provides reasoning paths as rationales to interpret model’s decision.

## 1 Introduction

Open-Domain Question Answering (ODQA), one of the most knowledge-intensive NLP tasks, requires QA models to infer out-of-context knowledge to the given single question. Following the pioneering work by Chen et al. (2017), ODQA systems often assume to access an external text corpus (e.g., Wikipedia) as an external knowledge source. Due to the large scale of such textual knowledge sources (e.g., 20GB for Wikipedia), it cannot be encoded in the model parameters. Therefore, most works retrieve relevant passages as knowledge and thus named *Open-Book* models (Roberts et al., 2020), with an analogy to referring textbooks during an exam. Another line of *Closed-book* models (Roberts et al., 2020) assume knowledge could be stored implicitly in parameters of Language Models (LM, e.g. BERT (Devlin et al., 2019) and

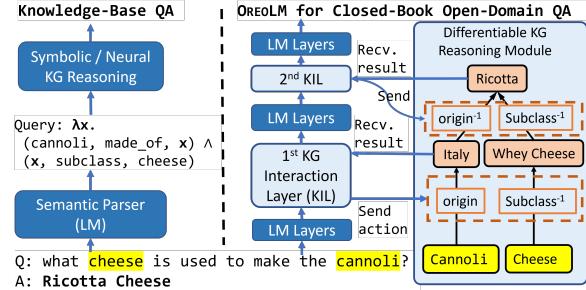


Figure 1: An Illustrative figure of OREOLM. Compared with previous KBQA systems that stack reasoner on top of LM, OREOLM enables interaction between the two.

T5 (Raffel et al., 2020)). These LMs directly generate answers without retrieving from an external corpus and thus benefit from faster inference speed and simpler training. However, current LMs still miss a large portion of factual knowledge (Pörner et al., 2020; Lewis et al., 2021a), and are not competitive with *Open-Book* models.

To improve the knowledge coverage of LM, one natural choice is to leverage knowledge stored in Knowledge Graph ( $\mathcal{KG}$ , e.g. FreeBase (Bollacker et al., 2008) and WikiData (Vrandecic and Krötzsch, 2014)), which explicitly encodes world knowledge via relational triplets between entities. There are several good properties of  $\mathcal{KG}$ : 1) a  $\mathcal{KG}$  triplet is a more abstract and compressed representation of knowledge than text, and thus  $\mathcal{KG}$  could be stored in memory and directly enhance LM without using an additional retrieval model; 2) the structural nature of  $\mathcal{KG}$  could support logical reasoning (Ren et al., 2020) and infer missing knowledge through high-order paths (Lao et al., 2011; Das et al., 2018). Taking the question “what cheese is used to make the desert cannoli?” as an example, even if this relational fact is missing in  $\mathcal{KG}$ , we could still leverage high-order relationships, e.g., both Ricotta Cheese and Cannoli are specialties in Italy, to infer the answer “Ricotta Cheese.”

In light of the good properties of  $\mathcal{KG}$ , there are

several efforts on building Knowledge Base Question Answering (KBQA) systems. As is illustrated in Figure 1(a), most KBQA models use LM as a parser to map textual questions into a structured form (e.g., SQL query or subgraph), and then based on  $\mathcal{KG}$ , the queries could be executed by symbolic reasoning (Berant et al., 2013) or neural reasoning (e.g. Graph Neural Networks) (Sun et al., 2019) to get the answer. Another recent line of research (Verga et al., 2021; Yu et al., 2022b) tries to encode the knowledge graph as the *memory* into LM parameters. However, for most methods discussed above, LM is not interacting with  $\mathcal{KG}$  to correctly understand the question, and the answer is usually restricted to a node or edge in  $\mathcal{KG}$ .

In this paper, we propose knOwledge REasOning empowered Language Model (OREOLM), a model architecture that can be applied to Transformer-based LMs to improve *Closed-Book* ODQA. As is illustrated in Figure 1(b), the key component is the Knowledge Interaction Layers (KIL) inserted amid LM layers, which is like cream filling within two waffles, leading to our model’s name OREO. KIL interacts with a  $\mathcal{KG}$  reasoning module, in which we maintain different reasoning paths for each entity in the question. We formulate the retrieval and reasoning process as a contextualized *random walk* over the  $\mathcal{KG}$ , starting from the in-context entities. Each KIL is responsible for one reasoning step. It first predicts a relation distribution for every in-context entity, and then the  $\mathcal{KG}$  reasoning module traverses the graph following the predicted relation distribution. The reasoning result in each step is summarized as a weighted averaged embedding over the retrieved entities from the traversal.

By stacking  $T$  layers of KIL, OREOLM can retrieve entities that are  $T$ -hop away from in-context entities and help LM to answer open questions that require out-of-context knowledge or multi-hop reasoning. The whole procedure is fully differentiable, and thus OREOLM learns and infers in an end-to-end manner. We further introduce how to pre-train OREOLM over unlabelled Wikipedia corpus. In addition to the entity salient span masking objective, we introduce two self-supervised objectives to guide OREOLM to learn better entity and relation representations and how to reason over them.

We test OREOLM with RoBERTa and T5 as our base LMs. By evaluating on several single-hop ODQA datasets in *closed-book* setting, we show

that OREOLM outperforms existing baselines with fewer model parameters. Specifically, OREOLM helps more for questions with missing relations in  $\mathcal{KG}$ , and questions that require multi-hop reasoning. We further show that OREOLM can serve as a backbone for *open-book* setting and achieves comparable performance compared with the state-of-the-art QA systems with dedicated design. In addition, OREOLM has better interpretability as it can generate reasoning paths for the answered question and summarize general relational rules to infer missing relations.

This key contributions are as follows: (1) We propose OREOLM to integrate symbolic knowledge graph reasoning with neural LMs. Different from prior works, OREOLM can be seamlessly plugged into existing LMs. (2) We pretrain OREOLM with RoBERTa and T5 to on the Wikipedia corpus. OREOLM can bring significant performance gain on ODQA. (3) OREOLM offers interpretable reasoning paths for answering the question and high-order reasoning rules as rationales.

## 2 Methodology

**Preliminary** We denote a Knowledge Graph  $\mathcal{KG} = (\mathcal{E}, \mathcal{R}, \mathcal{A} = \{A_r\}_{r \in \mathcal{R}})$ , where each  $e \in \mathcal{E}$  and  $r \in \mathcal{R}$  is entity node and relation label.  $A_r \in \{0, 1\}^{|\mathcal{E}| \times |\mathcal{E}|}$  is a sparse adjacency matrix indicating whether relation  $r$  holds between a pair of entities. The task of knowledge graph reasoning aims at answering a factoid query  $(s, r, ?)$ , i.e., which target entity has relation  $r$  with the source entity  $s$ . If  $\mathcal{KG}$  is complete, we could simply get answers by checking the adjacency matrix, i.e.,  $\{\forall t : A_r[s, t] = 1\}$ . For incomplete  $\mathcal{KG}$  where many relational facts are missing, path-based reasoning approaches (Lao et al., 2011; Xiong et al., 2017; Das et al., 2018) have been proposed to answer the one-hop query via finding multi-hop paths. For example, to answer the query  $(s, \text{Mother}, ?)$ , a path  $s \xrightarrow{\text{Father}} j \xrightarrow{\text{Wife}} t$  could reach the target answer  $t$ . In this paper we try to integrate symbolic  $\mathcal{KG}$  reasoning into neural LMs and help it deal with ODQA problems.

**Overview of OREOLM** We illustrate the overall architecture of OREOLM in Figure 2. All the light blue blocks are our added components to support  $\mathcal{KG}$  reasoning, while the dark blue Transformer layers are knowledge-injected LM. The key component of OREOLM for conducting  $\mathcal{KG}$  reasoning is the Knowledge Interaction Layers (KIL), which are added amid LM layers to enable deeper

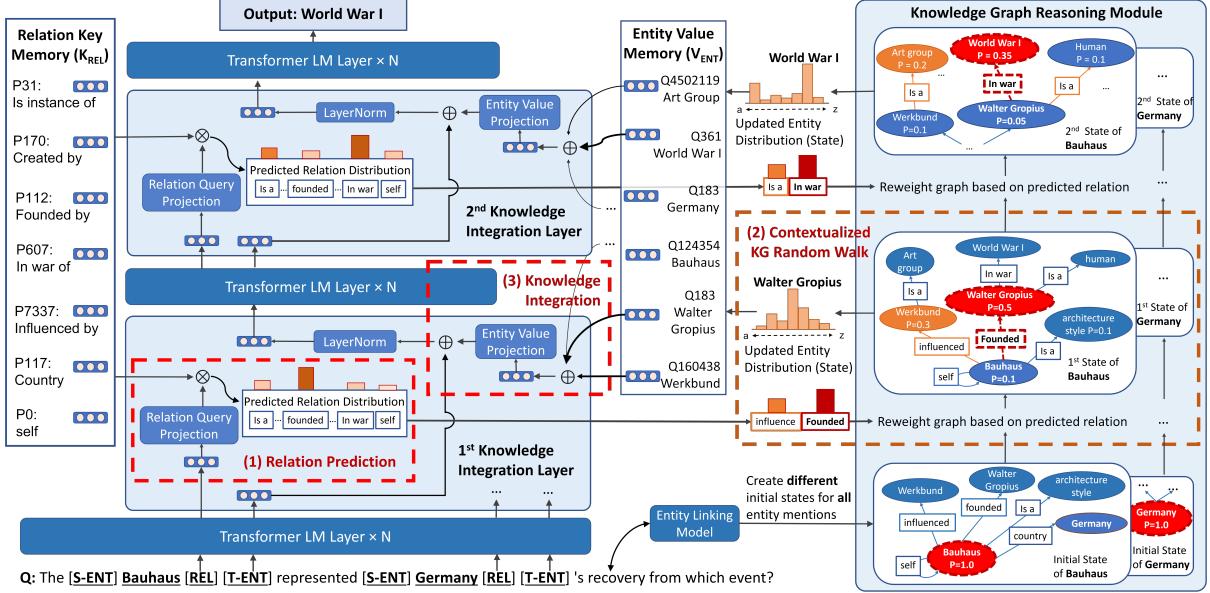


Figure 2: **Model architecture of OREOLM.** Three key procedures are highlighted in red dotted box: 1) **Relation Prediction** (Sec. 2.1.1): Knowledge Interaction Layers (KIL) predicts relation action for each entity mention. 2) **One-step State Transition** (Sec. 2.1.2): Based on the predicted relation,  $\mathcal{KG}$  re-weights each graph and conduct contextualized random walk to update entity distribution state. 3) **Knowledge Integration** (Sec. 2.2): An weighted aggregated entity embedding is added into a placeholder token as retrieved knowledge.

interaction with the  $\mathcal{KG}$ .

Given a question  $q = \text{"The Bauhaus represented Germany's recovery from which event?"}$ , QA model needs to extract knowledge about all  $n$  in-context entity mentions  $M = \{m_i\}_{i=1}^n$ , e.g., the history of “Germany” at the time when “Bauhaus” is founded, to get the answer  $a = \text{"World War I"}$ . Such open-domain Q&A can be abstracted as  $P(a|q, M)$ .

Starting from each mentioned entity  $m_i$ , we desire the model to learn to walk over the graph to retrieve relevant knowledge for answering this question. We define each reasoning path starting from the entity mention  $m_i$  as a chain of entities (states) random variables  $\rho_i = \{e_i^t\}_{t=0}^T$ , where each mentioned entity is the initial state, i.e.,  $e_i^0 = m_i$ . The union of all paths for this question is defined as  $\varrho = \{\rho_i\}$ , which contains the reasoning paths from each mentioned entity to answer the question.

OREOLM factorizes  $P(a|q, M)$  by incorporating possible paths  $\varrho$  as a latent variable, yielding:

$$\begin{aligned} P(a|q, M) &= \sum_{\varrho} P(\varrho|q, \{m_i\}_{i=1}^n) \cdot P(a|q, M, \varrho) \\ &= \sum_{\varrho} \left( \prod_{i=1}^n P(\rho_i|q, m_i) \right) \cdot P(a|q, \{m_i, \rho_i\}_{i=1}^n) \\ &= \sum_{\varrho} \left( \prod_{i=1}^n \prod_{t=1}^T \underbrace{P(e_i^t|q, e_i^{<t})}_{\mathcal{KG} \text{ Reasoning (2.1)}} \right) \underbrace{P(a|q, \{e_i^{0:T}\}_{i=1}^n)}_{\text{knowledge-injected LM (2.2)}} \end{aligned}$$

We assume (1) reasoning paths starting from dif-

ferent entities are generated independently; and (2) reasoning paths can be generated autoregressively.

In this way, the QA problem can be decomposed into two entangled steps: 1)  $\mathcal{KG}$  Reasoning, which autoregressively walks through the graph to get a path  $\rho_i$  starting from each entity mention  $m_i$ ; and 2) knowledge-injected LM, which benefits from the reasoning paths to obtain the out-context knowledge for answer prediction.

The relational path  $\rho_i$  in  $\mathcal{KG}$  Reasoning requires the selection of next entity  $e_i^t$  at each step  $t$ . We further decompose it into two steps: 1.a) relation prediction, in which LM is involved to predict the next-hop relation based on the current state and context; and 1.b) the non-parametric state transition, which is to predict the next-hop entity based on the  $\mathcal{KG}$  and the predicted relation. Formally:

$$P(e_i^t|q, e_i^{<t}) = \sum_r \underbrace{P_{rel}(r_i^t|q, e_i^{<t})}_{\mathcal{KG} \text{ Reasoning (2.1)}} \cdot \underbrace{P_{walk}(e_i^t|r_i^t, e_i^{<t})}_{\text{contextualized random walk (2.1.2)}}$$

We keep track of the entity distribution at each step  $t$  via the probability vector<sup>1</sup>  $\pi_i^{(t)} \in \mathcal{R}^{|\mathcal{E}|}$ , with  $\pi_i^{(t)}[e]$  being the probability of staying at entity  $e$ , i.e.,  $P(e_i^t = e|q, e_i^{<t})$ .

We highlight the three procedures in red dotted box in Figure 2. Still take the first reasoning step as an example. In the first red box within

<sup>1</sup>Throughout the paper, all vectors are row-vectors

221 KIL, we predict which relation action should be  
 222 taken for entity “Bauhaus”, and send the prediction  
 223 (e.g. “Founded”) to  $\mathcal{KG}$ . In the second red  
 224 box,  $\mathcal{KG}$  re-weights the graph and conduct context-  
 225 ualized random walk to update entity distribution,  
 226 where “Walter” has the highest probability. Finally,  
 227 weighted by the entity distribution, an aggregated  
 228 entity embedding is sent back to KIL and added  
 229 into a placeholder token as the knowledge, so the  
 230 the later LM layer knows to focus on the retrieved  
 231 “Walter”. We introduce these steps in the following.

232 **Input.** Initially, we first identify all  $N$  entity men-  
 233 tions  $\{m_i\}_{i=1}^N$  in the input question  $q$  as well as  
 234 the corresponding  $\mathcal{KG}$  entities<sup>2</sup>. For each men-  
 235 tion  $m_i$  we add three special tokens as the inter-  
 236 face for Knowledge Interaction Layers (KIL) to  
 237 send instruction and receive knowledge: we add  
 238 a [S-ENT] token before, and [REL], [T-ENT] to-  
 239 kens after each entity mention  $m_i$ . KIL can be  
 240 flexibly inserted into arbitrary LM intermediate  
 241 layer. By default, we just insert each KIL every  $N$   
 242 Transformer-based LM layers, thus the input to the  
 243  $t$ -th KIL are contextualized embeddings of each  
 244 token  $k$  as  $\text{LM}_k^{(t)}$ , including added special tokens.

## 245 2.1 LM involved $\mathcal{KG}$ Reasoning

246 We first introduce the reasoning process  
 247  $\mathbf{P}(e_i^t|q, e_i^{<t}) = \sum_r \mathbf{P}(r_i^t|q, e_i^{<t}) \cdot \mathbf{P}(e_i^t|r_i^t, e_i^{<t})$ .

### 248 2.1.1 Relation Prediction.

249 For each entity mention  $m_i$ , we desire to predict  
 250 which relation action should take  $r_i^t$  as instruction  
 251 to transit state. We define the predicted relation  
 252 probability vector  $\gamma_i^{(t)} = \mathbf{P}_{\text{rel}}(r_i^t|q, e_i^{<t}) \in \mathcal{R}^{|\mathcal{R}|}$   
 253 representing the relation distribution to guide walk-  
 254 ing through the graph. Denote the corresponding  
 255 [REL] token as  $\text{REL}[i]$  (and similarly for other spe-  
 256 cial tokens). The contextual embedding  $\text{LM}_{\text{REL}[i]}^{(t)}$   
 257 encode the relevant information in question  $q$  that  
 258 hints next relation. We maintain a global relation  
 259 key memory  $K_{\text{rel}} \in \mathbb{R}^{|\mathcal{R}| \times d}$  storing each relation’s  
 260  $d$ -dimentional embedding. To calculate similarity,  
 261 we first get relation query  $Q_{\text{REL}[i]}^{(t)}$  by projecting re-  
 262 lation token’s embedding into the same space of  
 263 key memory via a projection head Q-Proj<sup>3</sup> followed  
 264 by a LayerNorm (abbreviated as LN), and then cal-

<sup>2</sup>For Wikipedia pretraining, we use the ground-truth entity label as one-hot initialization for  $\pi_i^0$ . For downstream tasks we use GENRE (Cao et al., 2021) to get top 5 entity links.

<sup>3</sup>We denote a non-linear MLP projection as  $X\text{-Proj}(h) = W_2^X \sigma(W_1^X h + b_1) + b_2$ , where X have different instantiation.

265 culate dot-product similarity followed by softmax:

$$266 Q_{\text{REL}[i]}^{(t)} = \text{LN}^{(t)}(\text{Q-Proj}^{(t)}(\text{LM}_{\text{REL}[i]}^{(t)})) \\ 267 \gamma_i^{(t)} = \mathbf{P}_{\text{rel}}(r_i^t|q, e_i^{<t}) = \text{Softmax}(Q_{\text{REL}[i]}^{(t)} K_{\text{rel}}^T)$$

268 Note that the relation queries  $\text{LM}_{\text{REL}[i]}^{(t)}$  are dif-  
 269 ferent for every mention  $m_i$  and reasoning step  $t$   
 270 depending on the context, and thus the the relation  
 271 distributions  $\gamma_i^{(t)}$  gives contextualized predictions  
 272 based on the question  $q$ . The predicted relations  
 273 are sent to the knowledge graph reasoning module  
 274 as instruction to conduct state transition.

### 275 2.1.2 Contextualized KG Random Walk

276 Next, we introduce how we conduct state trans-  
 277 ition  $\mathbf{P}_{\text{walk}}(e_i^t|r_i^t, e_i^{<t})$ . One classic transi-  
 278 on algorithm is random walk, which is a special  
 279 case of markov chain, i.e. the transition proba-  
 280 bility only depends on previous state. Consider  
 281 a state at entity  $s$ , the probability walking to tar-  
 282 get  $t$  is  $\frac{1}{\deg(s)}$  if  $A[s, t] = 1$ . Based on it, we  
 283 define the Markov transition matrix for random  
 284 walk as  $M_{rw} = \mathbf{D}_A^{-1} A$ , where the degree matrix  
 285  $\mathbf{D}_A \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$  is defined as the diagonal matrix  
 286 with the degrees  $\deg(1), \dots, \deg(|\mathcal{E}|)$  on the diag-  
 287 onal. With random walk Markov matrix  $M_{rw}$  we can  
 288 transit the state distribution as:  $\pi^{(t)} = \pi^{(t-1)} M$ ,  
 289 The limitation of random walk is that the transition  
 290 strategy is not dependent on the question  $q$ . We thus  
 291 propose a Contextualized Random Walk (CRW).

292 Based on the predicted relation distribution  $\gamma_i^{(t)}$ ,  
 293 we calculate a different weighted adjacency matrix  
 294  $\tilde{A}_i^{(t)} \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$  by adjusting the edge weight:

$$295 \tilde{A}_i^{(t)} = \sum_{r \in \mathcal{R}} w_r \gamma_i^{(t)} A_r \\ 296 M_{crw_i}^{(t)} = \mathbf{D}_{\tilde{A}_i^{(t)}}^{-1} \tilde{A}_i^{(t)}, \quad \forall i \in [1, N],$$

297 where  $w_r$  is a learnable importance weight for re-  
 298 lation  $r$  that helps solving downstream tasks, and  
 299  $\gamma_{i,r}^{(t)}$  is the probability corresponding to relation  $r$  in  
 300  $\gamma_i^{(t)}$ . With the transition matrix  $M_{crw_i}^{(t)}$ , the state  
 301 transition is defined as  $\pi_i^{(t)} = \pi_i^{(t-1)} M_{crw_i}^{(t)}$ .

302 CRW is expressive that allows each reasoning  
 303 path  $\rho_i$  has its own transition matrix. However, as  
 304 the total number of entity nodes  $|\mathcal{E}|$  could be very  
 305 large (e.g., 5M for WikiData), we cannot afford to  
 306 update the full adjacency matrix for every mention  
 307 in every batch. For efficient implementation, we  
 308 adopt a scatter-gather pipeline to implement graph  
 309 walking as is shown in Algorithm 1 in Appendix.

310  
311  
312  
313  
The complexity is number of in-batch entities times  
number of edges in  $T$ -hop subgraph starting from  
these entities, i.e.  $\mathcal{O}(n \times \#edge)$ , and thus this  
operation is not expensive.

## 314 2.2 Knowledge-Injected LM

315 After we get the updated entity distribution  $\pi_i^{(t)}$ ,  
316 we want to inject such information back to the LM  
317 without harming its overall structure. We maintain  
318 a global entity embedding value memory  $V_{ent} \in$   
319  $\mathbb{R}^{|\mathcal{E}| \times d}$  storing entity embeddings. In each batch,  
320 we only consider the entities within the sampled  
321 local subgraph. We thus get an entity index list  $I$  as  
322 query to sparsely retrieve a set of candidate entity  
323 embeddings, and then aggregate them weighted by  
324 entity distribution and embedding table. We then  
325 use a Value Projection block to map the aggregated  
326 entity embedding into the space of LM, and then  
327 directly add the transformed embedding back to  
328 the output of T-ENT.

$$329 V_i^{(t)} = \text{V-Proj}^{(t)}(\pi_i^{(t)} \cdot V_{ent}[I]) \quad (1)$$

$$330 \widehat{LM}_{T-ENT[i]}^{(t)} = \text{LN}^{(t)}(LM_{T-ENT[i]}^{(t)} + V_i^{(t)}) \quad (2)$$

331 Then, we just take all  $\widehat{LM}_{T-ENT}^{(t)}$  as input to next  
332 Transformer-based LM layer to learn the interaction  
333 between the retrieved knowledge with in-context  
334 words via self-attention.

335 By repeating the KIL for  $T$  times, the final rep-  
336 resentation  $\widehat{LM}^T$  is conditioned on the reasoning  
337 paths  $\rho_i = e_i^{0:T}$ , which reaches entities that are  $T$ -  
338 hop away from initial entity  $m_i$  in the question. Fi-  
339 nally, we can predict the answer of open questions  
340  $P(a|q, \{e_i^{0:T}\}_{i=1}^n)$  by taking knowledge-injected  
341 representation  $\widehat{LM}^T$  for span extraction, entity pre-  
342 diction or direct answer generation.

## 343 2.3 Pre-Train OREOLM to Reason

344 The design of OREOLM allows it to be trained  
345 end-to-end given question answering datasets.  
346 However, due to small coverage of knowledge  
347 facts for existing QA datasets, we need to pretrain  
348 OREOLM on a large-scale corpus to get good entity  
349 embeddings.

350 **Salient Span Masking** One straightforward  
351 approach is to use Salient Span Masking (SSM)  
352 objective (Guu et al., 2020), which mask out  
353 entities or noun tokens that require a certain  
354 out-of-context knowledge. We mainly mask out  
355 entities for guiding OREOLM to reason. Instead  
356 of randomly mask entity mentions, we explicitly

357 sample a set of entity IDs and mask out every  
358 mentions linking to these entities. This could avoid  
359 model simply copy the entity from the context  
360 to fill in the blank. We also follow (Yang et al.,  
361 2019) to mask out consecutive token spans. We  
362 then calculate cross-entropy loss on each salient  
363 span masked (SSM) token as  $\mathcal{L}_{SSM}$ .

### 364 2.3.1 Weakly Supervised Training of KIL

365 Ideally, OREOLM can learn all the entity  
366 knowledge and how to access knowledge graph by  
367 solely optimizing  $\mathcal{L}_{SSM}$ . However, without a good  
368 initialization of entity and relation embeddings,  
369 KIL makes random prediction and the retrieved  
370 entities by  $\mathcal{KG}$  reasoning is likely to be irrelevant  
371 to the question. In this situation, KIL does  
372 not receive meaningful gradients to update the  
373 parameters and LM learns to ignore the knowledge.  
374 To avoid this cold-start problem and provide entity  
375 and relation embedding a good initialization,  
376 We utilize the following two external signals as  
377 self-supervised guidance.

378 **Entity Linking Loss** To initialize the large entity  
379 embedding tables in  $V_{ent}$ , we use other entities that  
380 are not masked as supervision. Similar to Févry  
381 et al. (2020), we force the output embedding of  
382 [S-ENT] token before the first KIL followed by a  
383 projection head E-Proj to be close to its correspond-  
384 ing entity embedding:

$$385 E_{S-ENT[i]} = \text{LN}(E\text{-Proj}(LM_{S-ENT[i]}^{(1)}))$$

$$386 P_{ent}^{(0)}(e|m_i, q) = \text{Softmax}(E_{S-ENT[i]} V_{ent}[I]^T)$$

$$387 \mathcal{L}_{ent} = \sum_{m_i} -\log P_{ent}^{(0)}(e|m_i, q) \cdot \pi_i^0[I]$$

388 Similar to Section 2.2, we only consider entities  
389 within the batch, denoted by index  $I$ . This  
390 contrastive loss guides each entity’s embedding  
391  $V_{ent}[e]$  closer to all its previously mentioned  
392 contextualized embedding, and thus memorizes  
393 those context as a good initialization for later  
394 knowledge integration.

395 **Weakly Supervised Relation Path Loss** As the  
396 entity mentions within each Wikipedia passage are  
397 naturally grounded to WikiData  $\mathcal{KG}$ , after we mask  
398 out several entities, we can utilize the  $\mathcal{KG}$  to get all  
399 reasoning paths from other in-context entities to the  
400 masked entities as weakly supervised relation label.

401 Formally, we define a **Grounded Dependency**  
402 **Graph**  $\mathcal{DG}$ , which contains all reasoning paths  
403 within  $T$ -step from other in-context entities to  
404 masked entities, and then define  $R_{\mathcal{DG}}(m_i, t)$  as

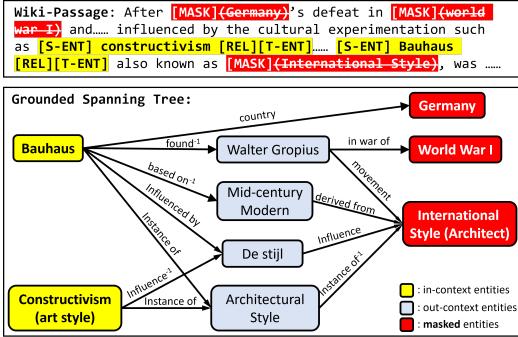


Figure 3: Pre-training sample w/ golden reasoning path. More real examples are shown in Table 5 in Appendix.

the set of all relations over every edges for entity mention  $m_i$  at  $t$ -th hop. Based on it, we define the weakly supervised relation label  $q_i^{(t)} \in \mathbb{R}^{|\mathcal{R}|}$  as the probabilistic vector which uniformly distributed on each relation in set. Note that we call uniformly-weighted  $q_i^{(t)}$  as weakly supervised because 1) some paths lead to multiple entities rather than only the target masked entity; 2) the correct relation is dependent on the context. Therefore,  $q_i^{(t)}$  only provides all potential candidates for reachability, and more fine-grained signals for reasoning should be learned from unsupervised  $\mathcal{L}_{SSM}$ . We adopt a list-wise ranking loss to guide the model to assign a higher score on these relations than others.

$$\mathcal{L}_{rel} = \sum_{m_i} \sum_{t=1}^T -\log P_{rel}^{(t)}(r|m_i, q) \cdot q_i^{(t)}.$$

Overall,  $\mathcal{L}_{ent}$  and  $\mathcal{L}_{rel}$  provide OREOLM with good initialization of the large  $\mathcal{KG}$  memory. Afterwards, via optimizing  $\mathcal{L}_{SSM}$ , the reasoning paths that provide informative knowledge receive positive gradient, guiding OREOLM to reason.

### 3 Experiments

The proposed KIL layers can be pugged into most Transformer-based Language Models without hurting its original structure. In this paper, we experiment with both encoder-based LM, i.e. RoBERTa-base ( $d = 768, l = 12$ ), and encoder-decoder LM, i.e. T5-base ( $d = 768, l = 12$ ) and T5-large ( $d = 1024, l = 24$ ). For all LMs, add 1 KIL layer or 2 KIL layers to the encoder layers. The statistics of  $\mathcal{KG}$  are shown in Table 4 in Appendix. Altogether, it takes about 0.67B parameter for  $\mathcal{KG}$  memory, which is affordable to load as model parameter. We pre-train all LMs using the combination of  $\mathcal{L}_{SSM}$ ,  $\mathcal{L}_{ent}$  and  $\mathcal{L}_{rel}$  for 200k steps on 8 V100 GPUs, with a batch size of

128 and default optimizer and learning rate in the original paper, taking approximately one week to finish pre-training of T5-large model, and 1-2 days for base model.

#### 3.1 Evaluate for *Closed-Book* QA

OREOLM is designed for improving *Closed-Book* QA, so we first evaluate it in this setting.

**Generative QA** Following the hyperparameters and setting in (Roberts et al., 2020), we directly fine-tune the T5-base and T5-large augmented by our OREOLM on the three single-hop ODQA datasets: Natural Question (NQ) (Kwiatkowski et al., 2019), WebQuestions (WQ) (Berant et al., 2013) and TriviaQA (TQA) (Joshi et al., 2017). To test OREOLM’s ability to solve complex questions, Apart from the , we also evaluate on two multi-hop QA datasets, i.e. **Complex WQ** (Talmor and Berant, 2018) and **HotpotQA** (Yang et al., 2018). Detailed dataset statistics and experimental setups are in Appendix B.

Experimental results are shown in Table 1. For all the datasets we use Exact Match accuracy as evaluation metric. On the three single-hop ODQA dataset, OREOLM with 2 KIL blocks achieves 3.3 absolute accuracy improvement to T5 (base), and 3.4 improvement to T5 (large). Compare with T5 model with more model parameter (e.g., T5-3B and T5-11B), our T5 (large) augmented by OREOLM could outperform T5-3B on NQ and WQ datasets that has much larger model parameter. In addition, OREOLM could use the generated reasoning path as rationale to interpret model’s prediction. We show examples on NQ in Table 7 in Appendix.

For the two multi-hop QA datasets, the performance improvement brought by OREOLM is more significant, i.e., 7.8 to T5 (base) and 8.2 to T5 (large). Notably, by comparing the T5-3B and T5-11B’s performance on HotpotQA (we take result from (Chen et al., 2022)), T5 (large) augmented by OREOLM achieves 1.2 higher than T5-11B. This shows that OREOLM is indeed very effective for improving *Closed-Book* QA performance, especially for complex questions.

**Entity Prediction** Encoder-based LM (i.e. RoBERTa) in most cases cannot be directly used for *Closed-Book* QA, but more serve as reader to extract answer span. However, Verga et al. (2021) propose a special evaluation setting as *Closed-Book Entity Prediction*. They add a single [MASK] token after the question, and

Models	#param	NQ	WQ	TQA	ComplexWQ	HotpotQA
T5 (Base)	0.22B	25.9	27.9	29.1	11.6	22.8
+ OREOLM ( $T=1$ )	0.23B + <u>0.68B</u>	28.3	30.6	32.4	20.8	24.1
+ OREOLM ( $T=2$ )	0.24B + <u>0.68B</u>	28.9	31.2	33.7	23.7	26.3
T5 (Large)	0.74B	28.5	30.6	35.9	16.7	25.3
+ OREOLM ( $T=1$ )	0.75B + <u>0.68B</u>	30.6	32.8	39.1	24.5	28.2
+ OREOLM ( $T=2$ )	0.76B + <u>0.68B</u>	<b>31.0</b>	<b>34.3</b>	<b>40.0</b>	<b>27.1</b>	<b>31.4</b>
T5-3B (Roberts et al., 2020)	3B	30.4	33.6	43.4	-	27.8
T5-11B (Roberts et al., 2020)	11B	32.6	37.2	50.1	-	30.2

Table 1: **Closed-Book Generative QA** performance of Encoder-Decoder LM on Single- and Multi-hop Dataset.

use its output embedding to classify WikiData entity ID. This restricts that answers must be entities that are covered by WikiData, which they call *WikiData-Answerable* questions. We follow Verga et al. (2021) to use such reduced version of WebQuestionsSP (**WQ-SP**) (Yih et al., 2015) and TriviaQA (**TQA**) as evaluation dataset, and finetune the RoBERTa (base) model augmented by OREOLM to classify entity ID. . We mainly compare OREOLM with EaE (Févry et al., 2020) and FILM (Verga et al., 2021), which are two  $\mathcal{KG}$  memory augmented LM. We also run experiments on KEPLER (Wang et al., 2019), a RoBERTa model pre-trained with knowledge augmented task.

Experimental results are shown in Table 2. Similar to the observation reported by Verga et al. (2021), adding  $\mathcal{KG}$  memory for this entity prediction task could significantly improve over vanilla LM, as most of the factual knowledge required to predict such entities are stored in  $\mathcal{KG}$ . By comparing with FILM (Verga et al., 2021), which is the state-of-the-art model in this setup, OREOLM with reasoning step ( $T = 2$ ) outperform FILM by 2.9, with smaller memory consumption.

### 3.2 Analyze $\mathcal{KG}$ Reasoning Module

In our previous studies, we find that using a higher reasoning step, i.e.  $T = 2$ , generally gets better performance than  $T = 1$ . We hypothesize that the  $\mathcal{KG}$  we use has many missing one-hop facts, and high-order reasoning helps to recover them and empower the model to answer related questions. To test whether OREOLM indeed has the ability to infer missing facts, we use **EntityQuestions (EQ)** (Sciavolino et al., 2021), which is a synthetic dataset by mapping each WikiData triplet to natural questions. We take RoBERTa-base model augmented by OREOLM trained on NQ as entity predictor and directly test its transfer performance on EQ dataset without further fine-tuning.

To test whether OREOLM could recover missing relation, we mask **all** the edges corresponding to each relation separately and make prediction again. The average results before and after removing edges are shown on the left part of Figure 4. When we remove all the edges to each relation, OREOLM with  $T = 1$  drops significantly, while  $T = 2$  could still have promising accuracy. To understand why OREOLM ( $T = 2$ ) is less influenced, in the right part of Figure 4, we generate reasoning path for each relation, by averaging the predicted probability score at each reasoning step, and pick the relation with top score. For example, to predict the “Capital” of a country, the model learns to find the living place of the president, or location of a country’s central bank. Both are very reasonable guess. Many previous works (Xiong et al., 2017) could also learn such rules in an ad-hoc manner and require costly searching or reinforcement learning, while OREOLM could learn such reasoning capacity for all relations end-to-end during pre-training.

**Ablation Studies** We conduct several ablation studies to evaluate which model design indeed contribute to the model. As is shown in the bottom blocks in Table 2, we first remove the  $\mathcal{KG}$  reasoning component and provide RoBERTa base model via concatenated KB triplets. For a fair comparison, we also train such model using  $\mathcal{L}_{SSM}$  over the same WikiDataset. The results of such model just close the KEPLER results, but much lower than other models with explicit knowledge memory. We further investigate the role of pre-training tasks. Without pre-training, the OREOLM only performs slightly better than RoBERTa baseline, due to the cold-start problem of entity and relation embedding. We further show that removing  $\mathcal{L}_{ent}$  and  $\mathcal{L}_{rel}$  could significantly influence final performance. The current combination is the best choice to train OREOLM to reason.

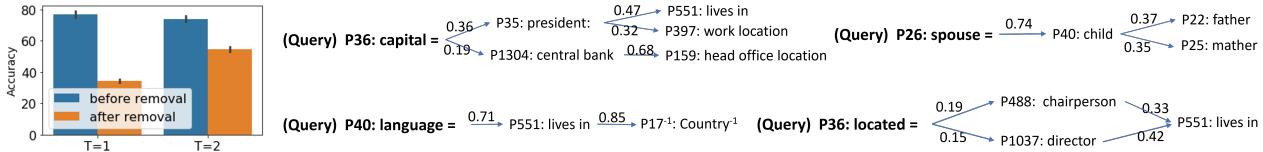


Figure 4: **Testing the reasoning capacity of OREOLM to infer missing relations.** On the **left**, the barplot shows the transfer performance on EQ before and after removing relation edges, OREOLM ( $T = 2$ ) is less influenced. On the **right** shows reasoning paths (rules) automatically generated by OREOLM for each missing relation.

Models	#param (B)	WQ-SP	TQA
EaE (Févy et al., 2020)	0.11 + <u>0.26</u>	62.4	24.4
FILM (Verga et al., 2021)	0.11 + <u>0.72</u>	78.1	37.3
KEPLER (Wang et al., 2019)	0.12	48.3	24.1
RoBERTa (Base)	0.12	43.5	21.3
+ OREOLM ( $T=1$ )	0.12 + <u>0.68</u>	80.1	39.7
+ OREOLM ( $T=2$ )	0.13 + <u>0.68</u>	<b>80.9</b>	<b>40.3</b>
Ablation Studies			
RoBERTa + Concat KB + $\mathcal{L}_{SSM}$	0.12	47.1	22.6
+ OREOLM ( $T=2$ ) w/o PT	0.13 + <u>0.68</u>	46.9	22.7
w. $\mathcal{L}_{SSM}$	0.13 + <u>0.68</u>	51.9	26.8
w. $\mathcal{L}_{SSM} + \mathcal{L}_{ent}$	0.13 + <u>0.68</u>	68.4	35.7

Table 2: **Closed-Book Entity Prediction** performance of Encoder LM on WikiData-Answerable Dataset.

Models	#param (B)	NQ	TQA
Graph-Retriever (Min et al., 2019)	0.11	34.7	55.8
REALM (Guu et al., 2020)	0.33 + <u>16</u>	40.4	-
DPR (Karpukhin et al., 2020) + BERT	0.56 + <u>16</u>	41.5	56.8
+ OREOLM (DPR, $T=2$ )	0.57 + <u>17</u>	43.7	58.5
FiD (Base) = DPR + T5 (Base)	0.44 + <u>16</u>	48.2	65.0
+ OREOLM (T5, $T=2$ )	0.45 + <u>17</u>	49.3	67.1
+ OREOLM (DPR & T5, $T=2$ )	0.46 + <u>17</u>	51.1	68.4
FiD (Large) = DPR + T5 (Large)	0.99 + <u>16</u>	51.4	67.6
+ OREOLM (T5, $T=2$ )	0.99 + <u>17</u>	52.4	68.9
+ OREOLM (DPR & T5, $T=2$ )	1.00 + <u>17</u>	<b>53.2</b>	69.5
KG-FiD (Base) (Yu et al., 2022a)	0.44 + <u>16</u>	49.6	66.7
KG-FiD (Large) (Yu et al., 2022a)	0.99 + <u>16</u>	<b>53.2</b>	69.8
EMDR <sup>2</sup> (Sachan et al., 2021b)	0.44 + <u>16</u>	52.5	<b>71.4</b>

Table 3: **Open-Book QA Evaluation.**

### 3.3 Evaluate for *Open-Book QA*

Though OREOLM is designed for *Closed-Book QA*, the learned model can serve as backbone for *Open-Book QA*. We take DPR and FiD models and change the underlying LM encoder to OREOLM. For DPR retriever, we replace the question encoder to RoBERTa trained with OREOLM, fixing the passage embedding and only finetune on each downstream QA dataset. For FiD model, we replace the T5 to our pre-trained version. We also change the retriever with our fine-tuned DPR. Results in Table 3 show that by augmenting both retriever and generator, OREOLM improves a strong baseline like FiD, for about 3.1% for Base and 1.8% for Large, and it outperforms very recent KG-FiD model for 1.6% in base setting, and achieve comparative performance in large setting. Note that though our results is still lower than some recent models (e.g. EMDR<sup>2</sup>), these methods are dedicated architecture or training framework for *Open-Book QA*. We may integrate OREOLM with these models to further improve their performance.

## 4 Related Work

Due to the limitation of current LMs to memorize all knowledge, a line of works try to encode knowledge (significantly smaller than the

web corpus) as *memory* into LM parameter. Potential choices for such compressed knowledge include QA pairs (Chen et al., 2022; Lewis et al., 2021b), entity embedding (Févy et al., 2020) and reasoning cases (Das et al., 2021, 2022). There’s also several works utilizing Knowledge Graph(KG). FILM (Verga et al., 2021) turns KG triplets into memory. Given a question, LM retrieves most relevant triplet as answer. GreaseLM (Zhang et al., 2022) propose to interact LM with KG via a interaction node. JAKET (Yu et al., 2022b) encode text and KG independently and fuse information at late stage. We introduce and discuss with other related works in Sec. C in Appendix.

## 5 Conclusion

We presented OREOLM, a novel model that incorporates symbolic  $\mathcal{KG}$  reasoning with existing LMs. We showed that OREOLM can bring significant performance gain to open-domain QA benchmarks, both for closed-book and open-book settings, as well as encoder-only and encoder-decoder models. Additionally, OREOLM produces reasoning paths that helps interpret the model prediction. In future, we’d like to apply OREOLM to a broader range of in knowledge-intensive NLP tasks.

## 6 Limitations

**Limited Reasoning Steps** In our experiments, we show that using reasoning step  $T = 2$  has better performance to  $T = 1$  on one-hop and multi-hop (mostly two) QA datasets. Thus, it's a natural question about whether we could extending reasoning steps more? As previous KG reasoning mostly could support very long path (with LSTM design)

Though we didn't spend much time exploring before the paper submission, we indeed try using  $T = 3$ , but currently it didn't get better results. We hypothesize the following reasons: 1) A large portion of our current model's improvement relies on the weakly supervised relation pre-training. To do it, we construct a K-hop ( $K=2$  now) subgraph, and sample dependency graph based on it. The larger  $K$  we choose, the more noise is included into the generated relation label, in an exponential increasing speed. Thus, it's harder to get accurate reasoning path ground-truth for high-order  $T$ . Another potential reason is that within Transformer model, the representation space in lower and upper layer might be very different, say, encode more syntax and surface knowledge at lower layers, while more semantic knowledge at upper layers. Currently we adopt a MLP projection head, wishing to map integrated knowledge into the same space, but it might have many flaws and need further improvement.

**Large Entity Embedding Table requires Pre-Training and GPU resources** Our current design has a huge entity embedding table, which should be learned through additional supervision and could not directly fine-tune to downstream tasks. This is restricts our approach's usage.

**Require Entity Linking** Current model design requires an additional step of entity linking for incoming questions, and then add special tokens as interface. A truly end-to-end model should identify which elements to start conducting reasoning by its own without relying on external models.

**Only support relational path-based reasoning** Though there are lots of potential reasoning tasks, such as logical reasoning, commonsense reasoning, physical reasoning, temporal reasoning, etc. Our current model design mainly focus on path-based relational reasoning, and it should not work for other reasoning tasks at current stage.

**Unreasonable Assumption of Path Independence** When we derive equation 1,

we have the assumption that reasoning paths starting from different entities should be independent. This is not always correct, especially for questions that require logical reasoning, say, have conjunction or disjunction operation over each entity state. And thus our current methods might not work for those complex QA with logical dependencies.

## References

- Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong. 2020. [Learning to retrieve reasoning paths over wikipedia graph for question answering](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1533–1544. ACM.
- Kurt D. Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. [Freebase: a collaboratively created graph database for structuring human knowledge](#). In *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12, 2008*, pages 1247–1250. ACM.
- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. [Translating embeddings for modeling multi-relational data](#). In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 2787–2795.
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. [Autoregressive entity retrieval](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1870–1879. Association for Computational Linguistics.
- Wenhu Chen, Pat Verga, Michiel de Jong, John Wieting, and William Cohen. 2022. [Augmenting pre-trained language models with qa-memory for open-domain question answering](#). *CoRR*, abs/2204.04581.

722	Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V. Le. 2020. Neural symbolic reader: Scalable integration of distributed and symbolic representations for reading comprehension. In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	780
723		781
724		782
725		783
726		784
727		785
728		
729	Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. 2018. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.	786
730		787
731		788
732		789
733		
734		
735		
736		
737		
738	Rajarshi Das, Ameya Godbole, Ankita Naik, Elliot Tower, Robin Jia, Manzil Zaheer, Hannaneh Hajishirzi, and Andrew McCallum. 2022. Knowledge base question answering by case-based reasoning over subgraphs. <i>CoRR</i> , abs/2202.10610.	790
739		791
740		792
741		793
742		
743	Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan, Lazaros Polymenakos, and Andrew McCallum. 2021. Case-based reasoning for natural language queries over knowledge bases. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 9594–9611. Association for Computational Linguistics.	794
744		795
745		796
746		797
747		798
748		799
749		800
750		801
751		
752		
753	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 4171–4186. Association for Computational Linguistics.	802
754		803
755		804
756		805
757		
758		
759		
760		
761		
762		
763	Bhuwan Dhingra, Manzil Zaheer, Vidhisha Balachandran, Graham Neubig, Ruslan Salakhutdinov, and William W. Cohen. 2020. Differentiable reasoning over a virtual knowledge base. In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	815
764		816
765		817
766		818
767		819
768		820
769		821
770	Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. 2019. Cognitive graph for multi-hop reading comprehension at scale. In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 2694–2703. Association for Computational Linguistics.	822
771		823
772		
773		
774		
775		
776		
777		
778	Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. Scalable multi-hop relational reasoning for knowledge-aware question answering. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 1295–1309. Association for Computational Linguistics.	824
779		825
780		826
781		827
782		828
783		829
784		830
785		831
786		
787		
788		
789		
790	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: retrieval-augmented language model pre-training. <i>CoRR</i> , abs/2002.08909.	890
791		891
792		892
793		893
794	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers</i> , pages 1601–1611. Association for Computational Linguistics.	894
795		895
796		896
797		897
798		898
799		899
800		900
801		
802	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. <i>CoRR</i> , abs/2004.04906.	901
803		902
804		903
805		904
806	Pei Ke, Haozhe Ji, Yu Ran, Xin Cui, Liwei Wang, Linfeng Song, Xiaoyan Zhu, and Minlie Huang. 2021. Jointgt: Graph-text joint representation learning for text generation from knowledge graphs. In <i>Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021</i> , volume ACL/IJCNLP 2021 of <i>Findings of ACL</i> , pages 2526–2538. Association for Computational Linguistics.	905
807		906
808		907
809		908
810		909
811		910
812		911
813		912
814		913
815	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. <i>Trans. Assoc. Comput. Linguistics</i> , 7:452–466.	914
816		915
817		916
818		917
819		918
820		919
821		920
822		921
823		922
824	Ni Lao, Tom M. Mitchell, and William W. Cohen. 2011. Random walk inference and learning in A large scale knowledge base. In <i>Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, EMNLP 2011, 27-31 July 2011, John McIntyre Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group of the ACL</i> , pages 529–539. ACL.	923
825		924
826		925
827		926
828		927
829		928
830		929
831		930
832	Patrick S. H. Lewis, Pontus Stenetorp, and Sebastian Riedel. 2021a. Question and answer test-train overlap in open-domain question answering datasets. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19</i>	931
833		932
834		933
835		934
836		935
837		936

838	- 23, 2021, pages 1000–1008. Association for Computational Linguistics.	894
839		895
840		896
841		
842		
843		
844		
845		
846		
847		
848		
849		
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		
864		
865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878		
879		
880		
881		
882		
883		
884		
885		
886		
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926		
927		
928		
929		
930		
931		
932		
933		
934		
935		
936		
937		
938		
939		
940		
941		
942		
943		
944		
945		
946		
947		
948		
949		
950		
951		

- 952                  June 1-6, 2018, Volume 1 (Long Papers), pages 641–  
 953                  651. Association for Computational Linguistics.
- 954                  Pat Verga, Haitian Sun, Livio Baldini Soares, and  
 955                  William W. Cohen. 2021. [Adaptable and interpretable neural memory over symbolic knowledge](#).  
 956                  In *Proceedings of the 2021 Conference of the North*  
 957                  *American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages  
 958                  3678–3691. Association for Computational Linguistics.
- 963                  Denny Vrandecic and Markus Krötzsch. 2014. [Wikidata: a free collaborative knowledgebase](#). *Commun. ACM*, 57(10):78–85.
- 966                  Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhiyuan  
 967                  Liu, Juanzi Li, and Jian Tang. 2019. [KEPLER: A unified model for knowledge embedding and pre-trained](#)  
 968                  [language representation](#). *CoRR*, abs/1911.06136.
- 970                  Wenhan Xiong, Thien Hoang, and William Yang Wang.  
 971                  2017. [Deeppath: A reinforcement learning method](#)  
 972                  [for knowledge graph reasoning](#). In *Proceedings of*  
 973                  *the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 564–  
 975                  573. Association for Computational Linguistics.
- 977                  Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell,  
 978                  Ruslan Salakhutdinov, and Quoc V. Le. 2019.  
 979                  [Xlnet: Generalized autoregressive pretraining for](#)  
 980                  [language understanding](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 5754–5764.
- 985                  Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,  
 986                  William W. Cohen, Ruslan Salakhutdinov, and  
 987                  Christopher D. Manning. 2018. [Hotpotqa: A dataset](#)  
 988                  [for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*,  
 989                  pages 2369–2380. Association for Computational  
 990                  Linguistics.
- 994                  Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jian-  
 995                  feng Gao. 2015. [Semantic parsing via staged query](#)  
 996                  [graph generation: Question answering with knowl-](#)  
 997                  [edge base](#). In *Proceedings of the 53rd Annual Meet-*  
 998                  *ing of the Association for Computational Linguistics and the 7th International Joint Conference on Natu-*  
 999                  *ral Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31,*  
 1000                 *2015, Beijing, China, Volume 1: Long Papers*, pages  
 1001                 1321–1331. The Association for Computer Linguis-  
 1002                 *tics*.
- 1005                  Donghan Yu, Chenguang Zhu, Yuwei Fang, Wenhao  
 1006                  Yu, Shuohang Wang, Yichong Xu, Xiang Ren, Yim-  
 1007                  ing Yang, and Michael Zeng. 2022a. [Kg-fid: Infusing](#)  
 1008                  [knowledge graph in fusion-in-decoder for open-](#)  
 1009                  [domain question answering](#). In *Proceedings of the*
- 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 4961–4974. Association for Computational Linguistics. 1010
- Donghan Yu, Chenguang Zhu, Yiming Yang, and Michael Zeng. 2022b. [JAKET: joint pre-training of knowledge graph and language understanding](#). Conference on Artificial Intelligence, AAAI,. 1011
- Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D. Manning, and Jure Leskovec. 2022. [Greaselm: Graph reasoning enhanced language models for question answering](#). *CoRR*, abs/2201.08860. 1012
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. [ERNIE: enhanced language representation with informative entities](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 1441–1451. Association for Computational Linguistics. 1013
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. [Seq2sql: Generating structured queries from natural language using reinforcement learning](#). *CoRR*, abs/1709.00103. 1014
- 1015
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034

<b>Name</b>	<b>Number</b>	<b>dimension</b>	<b>#param (M)</b>	The answers are text spans in Wikipedia. We report short answer Exact Match (EM) performance.
# Entity	4,947,397	128	633	1063
# Relation	2,008	768	1.5	1064
# Edges	45,217,947	-	47	1065 1066 1067

Table 4: Statistics and parameter of  $\mathcal{KG}$  Memory.

## A Implementation Details

**Entity Linking during pre-training** We use the 2021 Jan. English dump of Wikidata and Wikipedia. For each wikipedia page, we link all entity mentions with hyperlinks to WikiData entity entry, augment all other mentions with same aliases, tokenize via each LM’s tokenizer and split into chunks with maximum token length allowed. We then construct induced k-hop subgraphs connecting entities within each chunk for quickly get grounded computational graph.

For entities, Wikipedia provides hyperlinks with ground-truth entity ID, but it doesn’t cover all the entity mentions, mostly hyperlinks only appear when this entity appears for the first time. Therefore, we first collect all entities appeared in hyperlinks as well as their aliases stored in WikiData, and then search any mentions that have any of these alias and link it to the corresponding entity.

**Implementation of Contextualized Random Walk** We first gather the entity and relation probability to each edge, and then scatter the probability to target nodes. This allows us to simultaneously conduct message passing with modified adjacency weight  $\tilde{A}_i^t$  for all entity mention  $m_i$  in parallel.

---

### Algorithm 1: Pytorch Pseudocode of CRW

---

```

def ContextualizedRandomWalk(
    i_init, KG,      # initial entity index and Graph
    w_deg, w_rel,   # inv(degree) and relation weights
    p_ent, p_rel    # entity and predicted relation dis-
                    # tribution tensor @ t-th step.
): -> FloatTensor
    # Get <src, rel, tgt> edge list of k-hop subgraph
    i_src, i_rel, i_tgt = k_hop_subgraph(i_init, KG)
    # Gather entity and relation probability to edge
    p_src = (p_ent * w_deg)[:, i_src] # N x n_edge
    p_rel = (p_rel * w_rel)[:, i_rel] # N x n_edge
    p_edge = l1_normalize(p_src * p_rel, dim=1)
    # Scatter edge probability to target node
    p_ent = scatter_add(src=p_edge, idx=i_tgt, dim=1)
    return p_ent  #(t+1)-th step's entity distribution

```

---

## B Dataset Details

**Natural Questions** (Kwiatkowski et al., 2019) contains questions from Google search queries, and

**WebQuestions (WQ)** (Berant et al., 2013) contains questions from Google Suggest API, and the answers are entities in Freebase.

**TriviaQA** (Joshi et al., 2017) contains trivia questions and answers are text spans from the Web. We report Exact Match (EM) performance.

**HotpotQA** (Yang et al., 2018) is a multi-hop QA dataset. There are two evaluation settings. In the *distractor setting*, 10 candidate paragraphs are provided for each question, of which there are two golden paragraphs. In the *full-wiki setting*, a model is required to extract paragraphs from the entire Wikipedia. We report Exact Match (EM) on full-wiki setting.

## C Other Related Works

### C.1 Introduce other related works

**Open-Domain Question Answering** aims at answering factoid questions by referring to a large-scale corpus. Most works adopt a two-stage pipeline proposed in (Chen et al., 2017) that combines a retriever with a neural reader. There also exists several QA works using  $\mathcal{KG}$  to help ODQA. For example, (Asai et al., 2020) and (Min et al., 2019) expand the entity graph following wikipedia hyperlinks or triplets in knowledge base. (Ding et al., 2019) extract entities from current context via entity-linking and turn them into a cognitive graph, and a graph neural network is applied on top of it to extract answer. (Dhingra et al., 2020) and (Lin et al., 2020) construct an entity-mention bipartite graph and then model the QA reasoning as graph traversal by filtering only the contexts that are relevant to the question.

**Knowledge-Base Question Answering** Traditional parsing-based methods parse the question into some intermediate query (e.g., SQL language, query graphs), which can execute on a knowledge base to get answer (Berant et al., 2013; Yih et al., 2015; Reddy et al., 2016; Zhong et al., 2017; Liang et al., 2017). However, existing knowledge bases suffer from low coverage of entities and relations required for open-ended questions. As an alternative, several works try to incorporate the structured knowledge into neural QA models for differentiable reasoning. (Lin et al., 2019) and (Feng et al., 2020) parse the question into a sub-graph

of knowledge base, and apply graph neural networks as reasoner to extract answers. (Chen et al., 2020) integrates general symbolic operations as basic units, and parse questions into compositional programs to answer general questions.

**Knowledge-augmented Language Models** explicitly incorporate external knowledge (e.g. knowledge graph) into LM. Overall, these approaches can be grouped into two categories: The first one is to explicitly inject knowledge representation into language model pre-training, where the representations are pre-computed from external sources (Zhang et al., 2019; Liu et al., 2021). For example, ERNIE (Zhang et al., 2019) encodes the pre-trained TransE (Bordes et al., 2013) embeddings as input. The second one is to implicitly model knowledge information into language model by performing knowledge-related tasks, such as entity category prediction (Yu et al., 2022b) and graph-text alignment (Ke et al., 2021). For example, JAKET jointly pre-trained both the KG representation and language representation by adding two self-supervised learning objectives (i.e., entity category prediction, relation type prediction) on knowledge graphs (Yu et al., 2022b).

## C.2 Discussion with Previous Works

**Compare with FILM** Though FILM has the advantage of end-to-end training and easily modification of knowledge memory, it simply stacks  $\mathcal{KG}$  module on top of LM without interaction, and can only handle one-hop relational query that is answerable by  $\mathcal{KG}$ . Our approach, OREOLM, follows the same *memory* idea by encoding  $\mathcal{KG}$  into LM parameter, and we desire LM and  $\mathcal{KG}$  reasoning module could interact and collaboratively improve each other.

Notably, OREOLM with  $T = 1$  shares a similar design with FILM. The major differences are: 1) they store every triplet as a key-value pair, while we explicitly keep the  $\mathcal{KG}$  adjacency matrix and conduct a random walk, which has smaller search space and is more controllable. 2) They add the memory on top of LM, and thus the knowledge could not help language understanding, and FILM could mainly help wikipedia-answerable questions. Instead, we insert the KIL layer amid LM layers to encourage interaction, and thus the model could also benefit encoder-decoder model (as shown above).

## Compare with Preious Path-Based Reasoning

**and Retrieval Pre-Training** Note that as our definition of entity state  $\pi_i$  and relation action  $\gamma_i$  are both continuous probabilistic vector, the whole  $\mathcal{KG}$  Reasoning is fully differentiable and thus could be integrated into LM seamlessly and trained end-to-end. This is different from previous path traversal works such as DeepPath (Xiong et al., 2017) and MINERVA (Das et al., 2018), which defines state and action as discrete and could only be trained via reinforcement learning rewards. The reasoner training is also different from passage retrieval pre-training (Guu et al., 2020; Sachan et al., 2021a), as the passage are naturally consisted of discrete tokens, and thus the reader is still required to re-encode the question with each passage, and different objectives are required to train retriever and reader separately.

## D Illustration of Pre-Trained Data and Reasoning Paths

The pre-training samples and reasoning paths (generated by T5-large on NQ dataset) is shown from Table 5-8.

1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177

1178  
1179  
1180  
1181  
1182

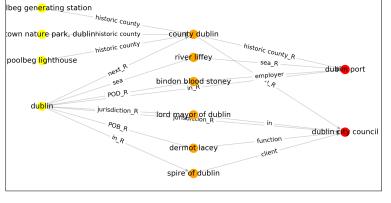
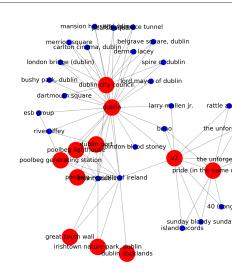
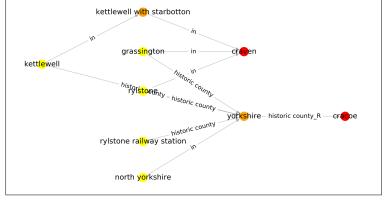
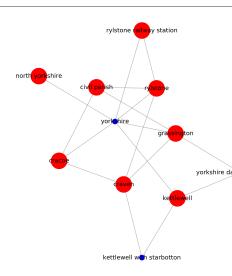
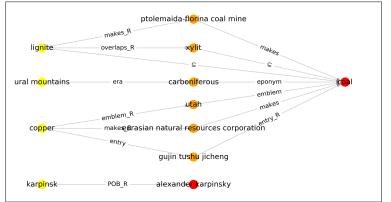
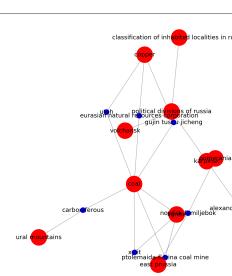
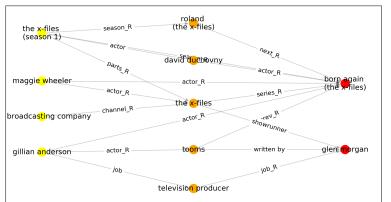
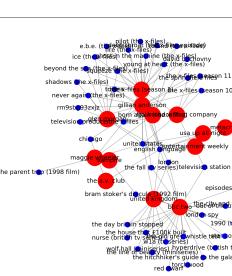
Title	Masked Text	Ground Truth	Dependency Graph	2-Hop Graph
Poolbeg	the lighthouse was [mask] [s-ent] [mask] [rel] [t-ent] completed in 1795. overview. the [s-ent] poolbeg[rel] [t-ent] "peninsula" is home to a number of landmarks including the [s-ent] [mask][rel] [t-ent] , the [s-ent] pool[mask] lighthouse[rel] [t-ent] , the [s-ent] irishtown nature park[rel] [t-ent] , the southern part of [s-ent] [mask][rel] [t-ent] ...	[ ' connected to land by the', ' great south wall', ' beg', ' dublin port', "s main power station", ' structures in', '48', ' a process to list the', ' after the station', ' including 3,', ' dublin city council', ' quarter' on the'		
Rylstone	it is situated very near to [s-ent] [mask][rel] [t-ent] and about 6 miles south west[mask] [s-ent] [mask]jlington[rel] [t-ent]. the population of the [s-ent] civil parish[rel] [t-ent] as of the 2011 census was 160. [s-ent] rylstone railway station[rel] [t-ent] opened in 1902, closed to passengers in 1930, and closed completely in 1969....	[ ' craven', ' cracoe', ' of', ' grass', ' the inspiration for', ' tour de france', ' stone', ' by will'...]		
Karpinsk	ologist [s-ent] [mask] [rel] [t-ent] . history.[mask]the settlement of bogoslovsk () was founded in either 1759 or in 1769. it remained one of the largest [s-ent] copper[rel] [t-ent] production centers in the [s-ent] urals[rel] [t-ent] [mask] [s-ent] [mask][rel] [t-ent] deposits started to be mined in 1911.....	[ ' alexander karpinsky', ' until 1917', ' coal', ' erman civilians, who', ' and', ' years of', ' forest laborers. moreover', ' in', ' the', ' framework of the', ' districts', ' karpinsk', 'insk'...]		
3 (The X-Files)	[s-ent] [mask][mask][rel] [t-ent] ". [s-ent] gillian anderson[rel] [t-ent] is absent[mask][mask] episode as she was on leave to give birth to her daughter piper at the time. this episode was the first[mask] not appear. reception. ratings. "3" premiered on the [s-ent] fox network[rel] [t-ent] on, and was first broadcast in the [s-ent] united kingdom[rel] [t-ent]....	[ 'ny had', ' episode', 'born again', ' from the', ' in which scully did', ' it was', 'egall', ' metacritic', ' as "wretched", ' fact that', ' background noise for a', ' heavy-handed attempts at', ' glen morgan', ' doing an episode on']		

Table 5: Example of Pre-training data points (Part 1).

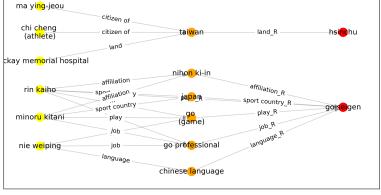
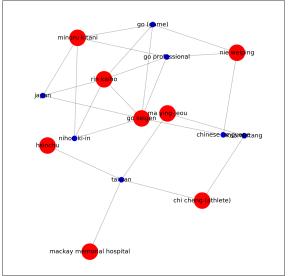
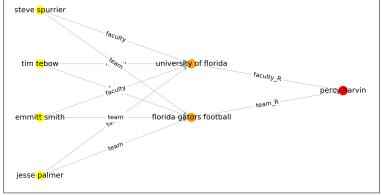
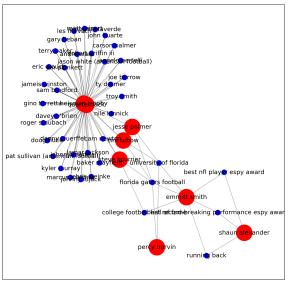
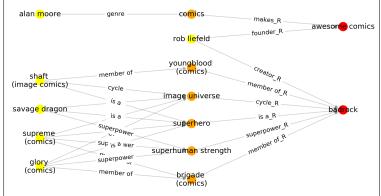
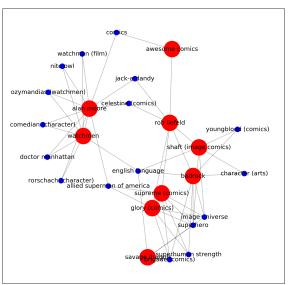
Title	Masked Text	Ground Truth	Dependency Graph	K-Hop Graph
Shen Chun-shan	his memoirs, he suffered his second stroke[mask][mask], even after his second stroke, he continued writing; his series of biographies of five go masters [s-ent] [mask][mask][mask][rel] [t-ent] , [s-ent] minoru kit[mask][rel] [t-ent] .....	he however', ' go seigen', 'ani', ' 2007, he', ' was hospital', 'hsinchu', 'after surgery', 'scale', ' continuing to improve.', ' his coma. in'....]		
2007 Florida Gators football team	[s-ent] tim[mask][mask][rel] [t-ent] completed 22 of 27 passes for 281 yards passing and also ran for [mask] yards on 6 carries. [s-ent] [mask] [rel] [t-ent] carried the ball 11 times for 113 yards[mask] two touchdowns and also caught 9 passes for 110[mask] receiving, becoming the first player in school history .....	I' tebow', '35', ' percy harvin', 'and', ' yards', '30-9', ' renewed their budding', 'gamecocks', 'gator', 'quarterback', 'set a career-high', 'of these five rushing', 'percy harvin', 'sinus infection.', 'actors', 'touchdown'		
Judgment Day (Awe-some Comics)	[s-ent] alan moore[rel] [t-ent] used "judgment day" to reject the violent, deconstructive clichés of 1990s comics inadvertently caused by his own work on " [s-ent] watchmen[rel] [t-ent] ", "" and " [s-ent] saga of the[mask][mask][rel] [t-ent] " and uphold the values of classic superhero comics. the series deals with a metacommentary of the notion of retcons to super-hero histories as [s-ent] alan moore[rel] [t-ent] [mask] for the characters of [s-ent] [mask][mask][rel] [t-ent] , to replace the shared universe they left when [s-ent] rob liefeld[rel] [t-ent] left image several years earlier. plot. in[mask], mick tombs/ [s-ent] knightsabre[rel] [t-ent]....	I' swamp thing', 'himself creates a new backstory', 'awesome comics', '1997', 'riptide', 'knightsabre appears to be', 'and sw', 'badrock', 'supreme', 'by', 'analyzing', 'cybernetic young', 'it, and it has', 'ue out', 'administrator for youngblood'		

Table 6: Example of Pre-training data points (Part 2).

Question	Answer	Reasoning Paths as Rationale
southern soul was considered the sound of what independent record label	[’Motown’]	soul music $\xrightarrow{\text{genre-R}}$ ? $\xrightarrow{\text{label}}$ ? independent record label $\xrightarrow{\text{belong}}$ ? $\xrightarrow{\text{is a-R}}$ ?
who is the bad guy in lord of the rings	[’Sauron’]	the lord of the rings (film series) $\xrightarrow{\text{theme}}$ ? $\xrightarrow{\text{characters}}$ ?
where was the mona lisa kept during ww2	[’the Ingres Museum’, "Château d’Amboise", ’Château de Chambord’, ’the Loc - Dieu Abbey’]	mona lisa $\xrightarrow{\text{creator}}$ ? $\xrightarrow{\text{tomb}}$ ? world war 2 $\xrightarrow{\text{take place}}$ ? $\xrightarrow{\text{located-R}}$ ?
who have won the world cup the most times	[’Brazil’]	fifa world cup $\xrightarrow{\text{parts}}$ ? $\xrightarrow{\text{land}}$ ?
who wrote the song the beat goes on	[’Sonny Bono’]	song $\xrightarrow{\text{album type-R}}$ ? $\xrightarrow{\text{author}}$ ?
who plays mrs. potato head in toy story	[’Estelle Harris’]	toy story $\xrightarrow{\text{series}}$ ? $\xrightarrow{\text{VO}}$ ?
who plays caroline on the bold and beautiful	[’Linsey Godfrey’]	the bold and the beautiful $\xrightarrow{\text{in work-R}}$ ? $\xrightarrow{\text{actor}}$ ?
where are the fruits of the spirit found in the bible	[’Epistle to the Galatians’]	bible $\xrightarrow{\text{parts}}$ ? $\xrightarrow{\text{parts}}$ ?
who is the only kaurava who survived the kurukshetra war	[’Yuyutsu’]	kaurava $\xrightarrow{\text{in work}}$ ? $\xrightarrow{\text{in work-R}}$ ? Kurukshetra War $\xrightarrow{\text{location}}$ $\xrightarrow{\text{live in-R}}$
what is the deepest depth in the oceans	[’Mariana Trench’]	ocean $\xrightarrow{\text{in}}$ ? $\xrightarrow{\text{lowest point}}$ ?
where did the french national anthem come from	[’Strasbourg’]	national anthem $\xrightarrow{\text{is a-R}}$ ? $\xrightarrow{\text{released in}}$ ?

Table 7: Example of QA prediction with reasoning path on NQ (part 1).

Question	Answer	Generated Reasoning Paths as Rationale
who sings the song where have all the flowers gone	[‘Pete Seeger’]	song $\xrightarrow{\text{album type-R}}$ ? $\xrightarrow{\text{actor}}$ ?
who discovered some islands in the bahamas in 1492	[‘Christopher Columbus’]	the bahamas $\xrightarrow{\text{entry}}$ ? $\xrightarrow{\text{entry-R}}$ ?
which type of wave requires a medium for transmission	[‘mechanical waves’, ‘heat energy’, ‘Sound’]	wave $\xrightarrow{\text{belong-R}}$ ? $\xrightarrow{\text{belong-R}}$ ?
land conversion through burning of biomass releases which gas	[‘traces of methane’, ‘carbon monoxide’, ‘hydrogen’]	gas $\xrightarrow{\text{belong-R}}$ ? $\xrightarrow{\text{as-R}}$ ?
the sum of the kinetic and potential energies of all particles in the system is called the	[‘internal energy’]	kinetic energy $\xrightarrow{\text{belong}}$ ? $\xrightarrow{\text{belong-R}}$ ? potential energy $\xrightarrow{\text{belong}}$ ? $\xrightarrow{\text{belong-R}}$ ?
who did seattle beat in the super bowl	[‘Denver Broncos’]	super bowl $\xrightarrow{\text{organizer}}$ ? $\xrightarrow{\text{league-R}}$ ?
what is the name of the girl romeo loved before juliet	[‘Rosaline’]	romeo $\xrightarrow{\text{in work}}$ ? $\xrightarrow{\text{in work-R}}$ ?
who will get relegated from the premier league 2016/17	[‘Hull City’, ‘Sunderland’, ‘Middlesbrough’]	premier league $\xrightarrow{\text{league-R}}$ ? $\xrightarrow{\text{POB}}$ ?
actress in the girl with the dragon tattoo swedish	[‘Noomi Rapace’]	sweden $\xrightarrow{\text{speaking}}$ ? $\xrightarrow{\text{mother tongue-R}}$ ?

Table 8: Example of QA prediction with reasoning path on NQ (part 2).