
Empowering Language Models with Knowledge Graph Reasoning for Question Answering

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

14 1 Introduction

15 Open-Domain Question Answering (ODQA), one of the most knowledge-intensive NLP tasks,
16 requires QA models to infer out-of-context knowledge to the given single question. Following the
17 pioneering work by Chen et al. (2017), ODQA systems often assume to access an external text corpus
18 (e.g., Wikipedia) as an external knowledge source. Due to the large scale of such textual knowledge
19 sources (e.g., 20GB for Wikipedia), it cannot be encoded in the model parameters. Therefore, most
20 works retrieve relevant passages as knowledge and thus named *Open-Book* models (Roberts et al.,
21 2020), with an analogy of referring to textbooks during an exam. Another line of *Closed-book*
22 models (Roberts et al., 2020) assume knowledge could be stored implicitly in parameters of Language
23 Models (LM, e.g. BERT (Devlin et al., 2019) and T5 (Raffel et al., 2020)). These LMs directly generate
24 answers without retrieving from an external corpus and thus benefit from faster inference speed and
25 simpler training. However, current LMs still miss a large portion of factual knowledge (Pörner et al.,
26 2020; Lewis et al., 2021a), and are not competitive with *Open-Book* models.

27 To improve the knowledge coverage of LM, one natural choice is to leverage knowledge stored in
28 Knowledge Graph (\mathcal{KG} , e.g. FreeBase (Bollacker et al., 2008) and WikiData (Vrandečić and Krötzsch,
29 2014)), which explicitly encodes world knowledge via relational triplets between entities. There are
30 several good properties of \mathcal{KG} : 1) a \mathcal{KG} triplet is a more abstract and compressed representation of
31 knowledge than text, and thus \mathcal{KG} could be stored in memory and directly enhance LM without using
32 an additional retrieval model; 2) the structural nature of \mathcal{KG} could support logical reasoning (Ren
33 et al., 2020) and infer missing knowledge through high-order paths (Lao et al., 2011; Das et al., 2018).
34 Taking the question “what cheese is used to make the desert cannoli?” as an example, even if this
35 relational fact is missing in \mathcal{KG} , we could still leverage high-order relationships, e.g., both Ricotta
36 Cheese and Cannoli are specialties in Italy, to infer the answer “Ricotta Cheese.”

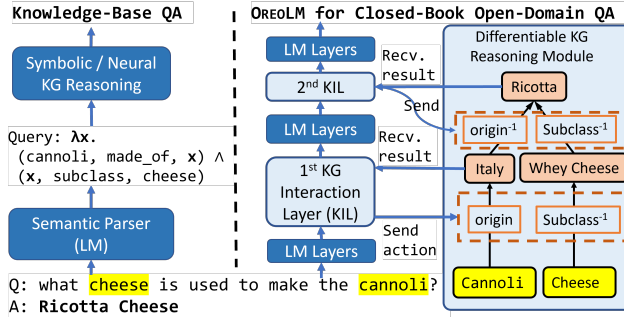


Figure 1: An Illustrative figure of OREOLM. Compared with previous KBQA systems that stack reasoner on top of LM, OREOLM enables interaction between the two.

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

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

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

61 We test OREOLM with RoBERTa and T5 as our base LMs. By evaluating on several single-hop ODQA
 62 datasets in *closed-book* setting, we show that OREOLM outperforms existing baselines with fewer
 63 model parameters. Specifically, OREOLM helps more for questions with missing relations in \mathcal{KG} , and
 64 questions that require multi-hop reasoning. We further show that OREOLM can serve as a backbone
 65 for *open-book* setting and achieves comparable performance compared with the state-of-the-art QA
 66 systems with dedicated design. In addition, OREOLM has better interpretability as it can generate
 67 reasoning paths for the answered question and summarize general rules to infer missing facts.

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

73 2 Methodology

74 **Preliminary** We denote a Knowledge Graph $\mathcal{KG} = (\mathcal{E}, \mathcal{R}, \mathcal{A} = \{A_r\}_{r \in \mathcal{R}})$, where each $e \in \mathcal{E}$
 75 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

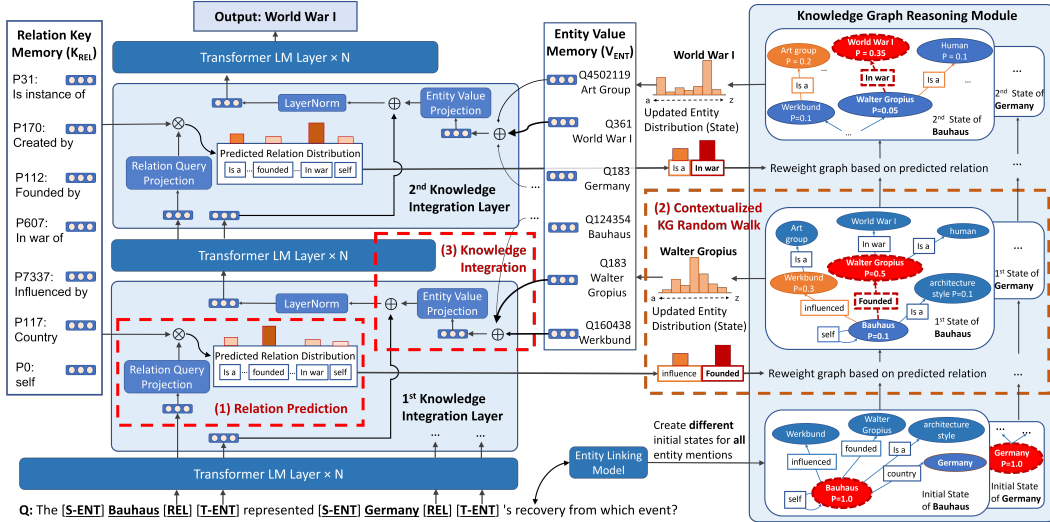


Figure 2: **Model architecture of OREOLM**. Three key procedures are highlighted in red dotted box: 1) **Relation Prediction** (Sec. 2.1.1): Knowledge Interaction Layers (KIL) predicts relation action for each entity mention. 2) **One-step State Transition** (Sec. 2.1.2): Based on the predicted relation, \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.

76 indicating whether relation r holds between a pair of entities. The task of knowledge graph
 77 reasoning aims at answering a factoid query $(s, r, ?)$, i.e., which target entity has relation r with
 78 the source entity s . If \mathcal{KG} is complete, we could simply get answers by checking the adjacency
 79 matrix, i.e., $\{\forall t : A_r[s, t] = 1\}$. For incomplete \mathcal{KG} where many relational facts are missing,
 80 path-based reasoning approaches Lao et al. (2011); Xiong et al. (2017); Das et al. (2018) have been
 81 proposed to answer the one-hop query via finding multi-hop paths. For example, to answer the
 82 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
 83 to integrate symbolic \mathcal{KG} reasoning into neural LMs and help it deal with ODQA problems.

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

89 Given a question $q = \text{“The Bauhaus represented Germany’s recovery from which event?”}$, QA model
 90 needs to extract knowledge about all n in-context entity mentions $M = \{m_i\}_{i=1}^n$, e.g., the history
 91 of “Germany” at the time when “Bauhaus” is founded, to get the answer $a = \text{“World War I”}$. Such
 92 open-domain Q&A can be abstracted as $P(a|q, M)$. Starting from each mentioned entity m_i , we
 93 desire the model to learn to walk over the graph to retrieve relevant knowledge and form a T -length
 94 reasoning path for answering this question, where T is a hyper-parameter denote the longest reasoning
 95 path required to answer the questions. We define each reasoning path starting from the entity mention
 96 m_i as a chain of entities (states) random variables $\rho_i = \{e_i^t\}_{t=0}^T$, where each mentioned entity is the
 97 initial state, i.e., $e_i^0 = m_i$. The union of all paths for this question is defined as $\varrho = \{\rho_i\}$, which
 98 contains the reasoning paths from each mentioned entity to answer the question.

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

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

100 We assume (1) reasoning paths starting from different entities are generated independently; and (2)
 101 reasoning paths can be generated autoregressively. In this way, the QA problem can be decomposed
 102 into two entangled steps: 1) \mathcal{KG} Reasoning, which autoregressively walks through the graph to get a
 103 path ρ_i starting from each entity mention m_i ; and 2) knowledge-injected LM, which benefits from the
 104 reasoning paths to obtain the out-context knowledge for answer prediction.

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

$$\underbrace{P(e_i^t|q, e_i^{<t})}_{\mathcal{KG} \text{ Reasoning (2.1)}} = \sum_r \underbrace{P_{rel}(r_i^t|q, e_i^{<t})}_{\text{relation prediction (2.1.1)}} \cdot \underbrace{P_{walk}(e_i^t|r_i^t, e_i^{<t})}_{\text{contextualized random walk (2.1.2)}}$$

109 We keep track of the entity distribution at each step t via the probability vector¹ $\pi_i^{(t)} \in \mathcal{R}^{|\mathcal{E}|}$, with
 110 $\pi_i^{(t)}[e]$ being the probability of staying at entity e , i.e., $P(e_i^t = e|q, e_i^{<t})$.

111 We highlight the three procedures in red dotted box in Figure 2. We take the first reasoning step
 112 starting from entity mention ‘‘Bauhaus’’ as an example. In the first red box within KIL, we predict
 113 which relation action should be taken for entity ‘‘Bauhaus’’, and send the prediction (e.g. ‘‘Founded’’) to
 114 \mathcal{KG} . In the second red box, \mathcal{KG} re-weights the graph and conducts contextualized random walk to
 115 update entity distribution, where ‘‘Walter’’ has the highest probability. Finally, weighted by the entity
 116 distribution, an aggregated entity embedding is sent back to KIL and added into a placeholder token
 117 as the knowledge, so the later LM layer knows to focus on the retrieved ‘‘Walter’’. We introduce these
 118 steps in the following.

119 **Input.** Initially, we first identify all N entity mentions $\{m_i\}_{i=1}^N$ in the input question q as well as the
 120 corresponding \mathcal{KG} entities². For each mention m_i we add three special tokens as the interface for
 121 Knowledge Interaction Layers (KIL) to send instruction and receive knowledge: we add a [S-ENT]
 122 token before, and [REL], [T-ENT] tokens after each entity mention m_i . KIL can be flexibly inserted
 123 into arbitrary LM intermediate layer. By default, we just insert each KIL every N Transformer-based
 124 LM layers, thus the input to the t -th KIL are contextualized embeddings of each token k as $\text{LM}_k^{(t)}$,
 125 including added special tokens.

126 2.1 LM involved \mathcal{KG} Reasoning

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

128 2.1.1 Relation Prediction.

129 For each entity mention m_i , we desire to predict which relation action should take r_i^t as instruction
 130 to transit state. We define the predicted relation probability vector $\gamma_i^{(t)} = P_{rel}(r_i^t|q, e_i^{<t}) \in \mathcal{R}^{|\mathcal{R}|}$
 131 representing the relation distribution to guide walking through the graph. Denote the corresponding
 132 [REL] token as $\text{REL}[i]$ (and similarly for other special tokens). The contextual embedding $\text{LM}_{\text{REL}[i]}^{(t)}$
 133 encode the relevant information in question q that hints next relation. We maintain a global relation
 134 key memory $K_{rel} \in \mathbb{R}^{|\mathcal{R}| \times d}$ storing each relation’s d -dimensional embedding. To calculate similarity,
 135 we first get relation query $Q_{\text{REL}[i]}^{(t)}$ by projecting relation token’s embedding into the same space of
 136 key memory via a projection head Q-Proj³ followed by a LayerNorm (abbreviated as LN), and then
 137 calculate dot-product similarity followed by softmax:

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

138 Note that the relation queries $\text{LM}_{\text{REL}[i]}^{(t)}$ are different for every mention m_i and reasoning step t
 139 depending on the context, and thus the the relation distributions $\gamma_i^{(t)}$ gives contextualized predictions

¹Throughout the paper, all vectors are row-vectors.

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

³We denote different non-linear MLP projections as X-Proj(h) = $W_2^X \sigma(W_1^X h + b_1) + b_2$.

140 based on the question q . The predicted relations are sent to the knowledge graph reasoning module as
 141 instruction to conduct state transition.

142 2.1.2 Contextualized KG Random Walk

143 Next, we introduce how we conduct state transition $P_{walk}(e_i^t|r_i^t, e_i^{<t})$. One classic transition
 144 algorithm is random walk, which is a special case of markov chain, i.e. the transition probability
 145 only depends on previous state. Consider a state at entity s , the probability walking to target t
 146 is $\frac{1}{deg(s)}$ if $A[s, t] = 1$. Based on it, we define the Markov transition matrix for random walk as
 147 $M_{rw} = D_A^{-1}A$, where the degree matrix $D_A \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$ is defined as the diagonal matrix with
 148 the degrees $deg(1), \dots, deg(|\mathcal{E}|)$ on the diagonal. With random walk Markov matrix M_{rw} we can
 149 transit the state distribution as: $\pi^{(t)} = \pi^{(t-1)}M$, The limitation of random walk is that the transition
 150 strategy is not dependent on the question q . We thus propose a Contextualized Random Walk (CRW).

151 Based on the predicted relation distribution $\gamma_i^{(t)}$, we calculate a different weighted adjacency matrix
 152 $\tilde{A}_i^{(t)} \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$ by adjusting the edge weight:

$$\tilde{A}_i^{(t)} = \sum_{r \in \mathcal{R}} w_r \cdot \gamma_{i,r}^{(t)} \cdot A_r, \quad M_{crw_i}^{(t)} = D_{\tilde{A}_i^{(t)}}^{-1} \tilde{A}_i^{(t)}, \quad \forall i \in [1, N],$$

153 where w_r is a learnable importance weight for relation r that helps solving downstream tasks, and
 154 $\gamma_{i,r}^{(t)}$ is the probability corresponding to relation r in $\gamma_i^{(t)}$. With the transition matrix $M_{crw_i}^{(t)}$, the
 155 state transition is defined as $\pi_i^{(t)} = \pi_i^{(t-1)}M_{crw_i}^{(t)}$.

156 CRW allows each reasoning path ρ_i to have its transition matrix. However, as the total number
 157 of entity nodes $|\mathcal{E}|$ could be huge (e.g., 5M for WikiData), we cannot afford to update the entire
 158 adjacency matrix for every in-batch mention. We thus adopt a scatter-gather pipeline to implement
 159 graph walking, as shown in Algorithm 1 in Appendix. The complexity is # of in-batch entities
 160 times # of edges in T -hop subgraph starting from these entities, i.e., $O(n \times \#edge)$, and thus this
 161 operation is not expensive. Another concern is why not using Graph Neural Networks (GNNs). We
 162 provide discussion in Sec. D.3 in Appendix.

163 2.2 Knowledge-Injected LM

164 After we get the updated entity distribution $\pi_i^{(t)}$, we want to inject such information back to the
 165 LM without harming its overall structure. We maintain a global entity embedding value memory
 166 $V_{ent} \in \mathbb{R}^{|\mathcal{E}| \times d}$ storing entity embeddings. We only consider the entities within the sampled local
 167 subgraph in each batch. We thus get an entity index list I as the query to sparsely retrieve a set of
 168 candidate entity embeddings and then aggregate them weighted by entity distribution and embedding
 169 table. We then use a Value Projection block to map the aggregated entity embedding into the space of
 170 LM, and then directly add the transformed embedding back to the output of T-ENT.

$$V_i^{(t)} = \text{V-Proj}^{(t)}(\pi_i^{(t)} \cdot V_{ent}[I]), \quad \widehat{\text{LM}}_{\text{T-ENT}[i]}^{(t)} = \text{LN}^{(t)}(\text{LM}_{\text{T-ENT}[i]}^{(t)} + V_i^{(t)}). \quad (1)$$

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

173 By repeating the KILL for T times, the final representation $\widehat{\text{LM}}^T$ is conditioned on the reasoning
 174 paths $\rho_i = e_i^{0:T}$, which reaches entities that are T -hop away from initial entity m_i in the question.
 175 Finally, we can predict the answer of open questions $P(a|q, \{e_i^{0:T}\}_{i=1}^n)$ by taking knowledge-injected
 176 representation $\widehat{\text{LM}}^T$ for span extraction, entity prediction or direct answer generation.

177 2.3 Pre-Train OREOLM to Reason

178 The design of OREOLM allows end-to-end training given QA datasets. However, due to the small
 179 coverage of knowledge facts for existing QA datasets, we need to pretrain OREOLM on a large-scale
 180 corpus to get good entity embeddings.

181 **Salient Span Masking.** One straightforward approach is to use Salient Span Masking (SSM)
 182 objective (Guu et al., 2020) masks out entities or noun tokens requiring specific out-of-context

183 knowledge. We mainly mask out entities for guiding OREOLM to reason. Instead of randomly
 184 masking entity mentions, we explicitly sample a set of entity IDs and mask every mentions linking
 185 to these entities. This could prevent the model copy the entity from the context to fill in the blank.
 186 We also follow (Yang et al., 2019) to mask out consecutive token spans. We then calculate the
 187 cross-entropy loss on each salient span masked (SSM) token as \mathcal{L}_{SSM} .

188 2.3.1 Weakly Supervised Training of KIL

189 Ideally, OREOLM can learn all the entity knowledge and how to access the knowledge graph by
 190 solely optimizing \mathcal{L}_{SSM} . However, without a good initialization of entity and relation embeddings,
 191 KIL makes a random prediction, and the retrieved entities by \mathcal{KG} reasoning are likely to be
 192 unrelated to the question. In this situation, KIL does not receive meaningful gradients to update
 193 the parameters, and LM learns to ignore the knowledge. To avoid this cold-start problem and provide
 194 entity and relation embedding a good initialization, We utilize the following two external signals
 195 as self-supervised guidance.

196 **Entity Linking Loss.** To initialize the large entity embedding tables in V_{ent} , we use other entities
 197 that are not masked as supervision. Similar to Févry et al. (2020), we force the output embedding
 198 of [S-ENT] token before the first KIL followed by a projection head E-Proj to be close to its
 199 corresponding entity embedding:

$$E_{S-ENT[i]} = \text{LN}(\text{E-Proj}(\text{LM}_{S-ENT[i]}^{(1)})), \quad \mathcal{L}_{ent} = \sum_i -\log \text{Softmax}(E_{S-ENT[i]} V_{ent}[I]^T) \pi_i^0[I]. \quad (2)$$

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

203 **Weakly Supervised Relation Path Loss.** Entity mentions within each Wikipedia passage are
 204 naturally grounded to WikiData \mathcal{KG} . Therefore, after we mask out several entities, we can utilize the
 205 \mathcal{KG} to get all paths from other entities to the masked entities as weakly supervised relation labels.

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

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

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

218 3 Experiments

219 The proposed KIL layers can be pugged into most Transformer-based Language Models without
 220 hurting its original structure. In this paper, we experiment with both encoder-based LM, i.e.
 221 RoBERTa-base ($d = 768, l = 12$), and encoder-decoder LM, i.e. T5-base ($d = 768, l = 12$) and
 222 T5-large ($d = 1024, l = 24$). For all LMs, add 1 KIL layer or 2 KIL layers to the encoder layers.
 223 The statistics of \mathcal{KG} are shown in Table 4 in Appendix. Altogether, it takes about 0.67B parameter
 224 for \mathcal{KG} memory, which is affordable to load as model parameter. We pre-train all LMs using the
 225 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
 226 default optimizer and learning rate in the original paper, taking approximately one week to finish
 227 pre-training of T5-large model, and 1-2 days for base model.

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 + 0.68B	28.3	30.6	32.4	20.8	24.1
+ OREOLM ($T=2$)	0.24B + 0.68B	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 + 0.68B	30.6	32.8	39.1	24.5	28.2
+ OREOLM ($T=2$)	0.76B + 0.68B	31.0	34.3	40.0	27.1	31.4
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** of Encoder-Decoder LM on (Single/Multi)-hop Dataset.

228 3.1 Evaluate for *Closed-Book QA*

229 OREOLM is designed for improving *Closed-Book QA*, so we first evaluate it in this setting. **Genera-**
 230 **tive QA** Following the hyperparameters and setting in Roberts et al. (2020), we directly fine-tune the
 231 T5-base and T5-large augmented by our OREOLM on the three single-hop ODQA datasets: Natural
 232 Question (**NQ**) (Kwiatkowski et al., 2019), WebQuestions (**WQ**) (Berant et al., 2013) and TriviaQA
 233 (**TQA**) (Joshi et al., 2017). To test OREOLM’s ability to solve complex questions, we also evaluate
 234 on two multi-hop QA datasets, i.e. **Complex WQ** (Talmor and Berant, 2018) and **HotpotQA** (Yang
 235 et al., 2018). Detailed dataset statistics and experimental setups are in Appendix C.

236 Experimental results are shown in Table 7. We use Exact Match accuracy as the metric for all the
 237 datasets. On the three single-hop ODQA datasets, OREOLM with 2 KIL blocks achieves 3.3 absolute
 238 accuracy improvement to T5-base, and 3.4 improvement to T5-large. Compared with T5 model
 239 with more model parameters (e.g., T5-3B and T5-11B), our T5-large augmented by OREOLM could
 240 outperform T5-3B on NQ and WQ datasets. In addition, OREOLM could use the generated reasoning
 241 path to interpret the model’s prediction. We show examples in Table 10 in Appendix.

242 For the two multi-hop QA datasets, the performance improvement brought by OREOLM is more
 243 significant, i.e., 7.8 to T5-base and 8.2 to T5-large. Notably, by comparing the T5-3B and T5-
 244 11B’s performance on HotpotQA (we take results from (Chen et al., 2022)), T5-large augmented by
 245 OREOLM achieves 1.2 higher than T5-11B. This shows that OREOLM is indeed very effective for
 246 improving *Closed-Book QA* performance, especially for complex questions.

247 **Entity Prediction.** Encoder-based LM (i.e. RoBERTa) in most cases cannot be directly used for
 248 *Closed-Book QA*, but more serve as reader to extract answer span. However, Verga et al. (2021) pro-
 249 pose a special evaluation setting as *Closed-Book Entity Prediction*. They add a single [MASK] token
 250 after the question, and use its output embedding to classify WikiData entity ID. This restricts that an-
 251 swers must be entities that are covered by WikiData, which they call *WikiData-Answerable* questions.
 252 We follow Verga et al. (2021) to use such reduced version of WebQuestionsSP (**WQ-SP**) (Yih et al.,
 253 2015) and TriviaQA (**TQA**) as evaluation dataset, and finetune the RoBERTa (base) model augmented
 254 by OREOLM to classify entity ID. . We mainly compare OREOLM with EaE (Férvy et al., 2020) and
 255 FILM (Verga et al., 2021), which are two \mathcal{KG} memory augmented LM. We also run experiments on
 256 KEPLER (Wang et al., 2019), a RoBERTa model pre-trained with knowledge augmented task.

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

262 3.2 Analyze \mathcal{KG} Reasoning Module

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

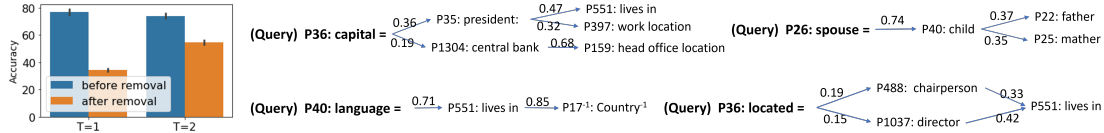


Figure 3: Testing the reasoning capacity of OREOLM to infer missing relations. On the left, the barplot shows the transfer performance on EQ before and after removing relation, OREOLM ($T=2$) is less influenced. On the right shows reasoning paths (rules) automatically generated by OREOLM.

Models	#param (B)	WQ-SP	TQA
EaE (Férvy et al., 2020)	0.11 + 0.26	62.4	24.4
FILM (Verga et al., 2021)	0.11 + 0.72	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 + 0.68	80.1	39.7
+ OREOLM ($T=2$)	0.13 + 0.68	80.9	40.3
Ablation Studies			
RoBERTa + Concat KB + \mathcal{L}_{SSM}	0.12	47.1	22.6
+ OREOLM ($T=2$) w/o PT	0.13 + 0.68	46.9	22.7
w. \mathcal{L}_{SSM}	0.13 + 0.68	51.9	26.8
w. $\mathcal{L}_{SSM} + \mathcal{L}_{ent}$	0.13 + 0.68	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 + 16	40.4	-
DPR Karpukhin et al. (2020) + BERT	0.56 + 16	41.5	56.8
+ OREOLM (DPR, $T=2$)	0.57 + 17	43.7	58.5
FiD (Base) = DPR + T5 (Base)	0.44 + 16	48.2	65.0
+ OREOLM (T5, $T=2$)	0.45 + 17	49.3	67.1
+ OREOLM (DPR & T5, $T=2$)	0.46 + 17	51.1	68.4
FiD (Large) = DPR + T5 (Large)	0.99 + 16	51.4	67.6
+ OREOLM (T5, $T=2$)	0.99 + 17	52.4	68.9
+ OREOLM (DPR & T5, $T=2$)	1.00 + 17	53.2	69.5
KG-FiD (Base) (Yu et al., 2022a)	0.44 + 16	49.6	66.7
KG-FiD (Large) (Yu et al., 2022a)	0.99 + 16	53.2	69.8
EMDR ² (Sachan et al., 2021b)	0.44 + 16	52.5	71.4

Table 3: Open-Book QA Evaluation.

270 To test whether OREOLM could recover missing relation, we mask **all** the edges corresponding to
 271 each relation separately and make the prediction again. The average results before and after removing
 272 edges are shown on the left part of Figure 3. When we remove all the edges to each relation, OREOLM
 273 with $T=1$ drops significantly, while $T=2$ could still have good accuracy. To understand why
 274 OREOLM ($T=2$) is less influenced, in the right part of Figure 3, we generate a reasoning path
 275 for each relation by averaging the predicted probability score at each reasoning step and pick the
 276 relation with the top score. For example, to predict the “Capital” of a country, the model learns
 277 to find the living place of the president, or the location of a country’s central bank. Both are very
 278 reasonable guesses. Many previous works (Xiong et al., 2017) could also learn such rules in an
 279 ad-hoc manner and require costly searching or reinforcement learning. In contrast, OREOLM could
 280 learn such reasoning capacity for all relations end-to-end during pre-training.

281 3.3 Evaluate for Open-Book QA

282 Though OREOLM is designed for *Closed-Book* QA, the learned model can serve as backbone
 283 for *Open-Book* QA. We take DPR and FiD models as baseline. For DPR retriever, we replace the
 284 question encoder to RoBERTa + OREOLM, fixing the passage embedding and only finetune on
 285 each downstream QA dataset. For FiD model, we replace the T5 + OREOLM. We also changed
 286 the retriever with our tuned DPR. Results in Table 3 show that by augmenting both retriever and
 287 generator, OREOLM improves a strong baseline like FiD, for about 3.1% for Base and 1.8% for
 288 Large, and it outperforms the very recent KG-FiD model for 1.6% in base setting, and achieve
 289 comparative performance in a large setting. Note that though our results is still lower than some
 290 recent models (e.g., EMDR²), these methods are dedicated architecture or training framework for
 291 *Open-Book* QA. We may integrate OREOLM with these models to further improve their performance.

292 4 Conclusion

293 We presented OREOLM, a novel model that incorporates symbolic \mathcal{KG} reasoning with existing LMs.
 294 We showed that OREOLM can bring significant performance gain to open-domain QA benchmarks,
 295 both for closed-book and open-book settings, as well as encoder-only and encoder-decoder models.
 296 Additionally, OREOLM produces reasoning paths that helps interpret the model prediction.

297 **References**

- 298 Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong.
299 2020. <https://openreview.net/forum?id=SJgVHkrYDH> Learning to retrieve reasoning paths over
300 wikipedia graph for question answering. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
301
- 302 Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013.
303 <https://www.aclweb.org/anthology/D13-1160/> Semantic parsing on freebase from question-answer
304 pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language
305 Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A
306 meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1533–1544. ACL.
- 307 Kurt D. Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008.
308 <https://doi.org/10.1145/1376616.1376746> Freebase: a collaboratively created graph database
309 for structuring human knowledge. In *Proceedings of the ACM SIGMOD International Conference
310 on Management of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12, 2008*, pages
311 1247–1250. ACM.
- 312 Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko.
313 2013. <http://papers.nips.cc/paper/5071-translating-embeddings-for-modeling-multi-relational-data>
314 Translating embeddings for modeling multi-relational data. In *Advances in Neural Information
315 Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013.
316 Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages
317 2787–2795.
- 318 Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021.
319 <https://openreview.net/forum?id=5k8F6UU39V> Autoregressive entity retrieval. In *9th In-
320 ternational Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7,
321 2021*. OpenReview.net.
- 322 Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [https://doi.org/10.18653/v1/P17-
323 1171](https://doi.org/10.18653/v1/P17-1171) Reading wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual
324 Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30
325 - August 4, Volume 1: Long Papers*, pages 1870–1879. Association for Computational Linguistics.
- 326 Wenhu Chen, Pat Verga, Michiel de Jong, John Wieting, and William Cohen. 2022.
327 <https://doi.org/10.48550/arXiv.2204.04581> Augmenting pre-trained language models with qa-
328 memory for open-domain question answering. *CoRR*, abs/2204.04581.
- 329 Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V. Le. 2020.
330 <https://openreview.net/forum?id=ryxjnREFwH> Neural symbolic reader: Scalable integration of
331 distributed and symbolic representations for reading comprehension. In *8th International Con-
332 ference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
333 OpenReview.net.
- 334 Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishna-
335 murthy, Alex Smola, and Andrew McCallum. 2018. <https://openreview.net/forum?id=Syg-YfWCW>
336 Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using rein-
337 forcement learning. In *6th International Conference on Learning Representations, ICLR 2018,
338 Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- 339 Rajarshi Das, Ameya Godbole, Ankita Naik, Elliot Tower, Robin Jia, Manzil Zaheer, Hannaneh Ha-
340 jishirzi, and Andrew McCallum. 2022. <http://arxiv.org/abs/2202.10610> Knowledge base question
341 answering by case-based reasoning over subgraphs. *CoRR*, abs/2202.10610.
- 342 Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan,
343 Lazaros Polymenakos, and Andrew McCallum. 2021. [https://doi.org/10.18653/v1/2021.emnlp-
344 main.755](https://doi.org/10.18653/v1/2021.emnlp-main.755) Case-based reasoning for natural language queries over knowledge bases. In *Proceedings
345 of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021,
346 Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 9594–9611.
347 Association for Computational Linguistics.

- 348 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019.
349 <https://doi.org/10.18653/v1/n19-1423> BERT: pre-training of deep bidirectional transform-
350 ers for language understanding. In *Proceedings of the 2019 Conference of the North American*
351 *Chapter of the Association for Computational Linguistics: Human Language Technologies,*
352 *NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers),*
353 pages 4171–4186. Association for Computational Linguistics.
- 354 Bhuwan Dhingra, Manzil Zaheer, Vidhisha Balachandran, Graham Neubig, Ruslan Salakhutdinov, and
355 William W. Cohen. 2020. <https://openreview.net/forum?id=SJxstlHFPH> Differentiable reasoning
356 over a virtual knowledge base. In *8th International Conference on Learning Representations,*
357 *ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.* OpenReview.net.
- 358 Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. 2019.
359 <https://doi.org/10.18653/v1/p19-1259> Cognitive graph for multi-hop reading compre-
360 hension at scale. In *Proceedings of the 57th Conference of the Association for Computational*
361 *Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers,* pages
362 2694–2703. Association for Computational Linguistics.
- 363 Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020.
364 <https://www.aclweb.org/anthology/2020.emnlp-main.99/> Scalable multi-hop relational reason-
365 ing for knowledge-aware question answering. In *Proceedings of the 2020 Conference on Empirical*
366 *Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020,* pages
367 1295–1309. Association for Computational Linguistics.
- 368 Thibault Févry, Livio Baldini Soares, Nicholas FitzGerald, Eunsol Choi, and Tom Kwiatkowski. 2020.
369 <http://arxiv.org/abs/2004.07202> Entities as experts: Sparse memory access with entity supervision.
370 *CoRR*, abs/2004.07202.
- 371 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020.
372 <http://arxiv.org/abs/2002.08909> REALM: retrieval-augmented language model pre-training. *CoRR*,
373 abs/2002.08909.
- 374 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017.
375 <https://doi.org/10.18653/v1/P17-1147> Triviaqa: A large scale distantly supervised chal-
376 lenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the*
377 *Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4,*
378 *Volume 1: Long Papers,* pages 1601–1611. Association for Computational Linguistics.
- 379 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau
380 Yih. 2020. <http://arxiv.org/abs/2004.04906> Dense passage retrieval for open-domain question
381 answering. *CoRR*, abs/2004.04906.
- 382 Pei Ke, Haozhe Ji, Yu Ran, Xin Cui, Liwei Wang, Linfeng Song, Xiaoyan Zhu, and Minlie Huang.
383 2021. <https://doi.org/10.18653/v1/2021.findings-acl.223> Jointgt: Graph-text joint representation
384 learning for text generation from knowledge graphs. In *Findings of the Association for Computa-*
385 *tional Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021,* volume ACL/IJCNLP
386 2021 of *Findings of ACL*, pages 2526–2538. Association for Computational Linguistics.
- 387 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris
388 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
389 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and
390 Slav Petrov. 2019. <https://transacl.org/ojs/index.php/tacl/article/view/1455> Natural questions: a
391 benchmark for question answering research. *Trans. Assoc. Comput. Linguistics*, 7:452–466.
- 392 Ni Lao, Tom M. Mitchell, and William W. Cohen. 2011. <https://aclanthology.org/D11-1049/> Random
393 walk inference and learning in A large scale knowledge base. In *Proceedings of the 2011 Con-*
394 *ference on Empirical Methods in Natural Language Processing, EMNLP 2011, 27-31 July 2011,*
395 *John McIntyre Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group*
396 *of the ACL*, pages 529–539. ACL.

- 397 Patrick S. H. Lewis, Pontus Stenetorp, and Sebastian Riedel. 2021a.
398 <https://www.aclweb.org/anthology/2021.eacl-main.86/> Question and answer test-train overlap in
399 open-domain question answering datasets. In *Proceedings of the 16th Conference of the European*
400 *Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online,*
401 *April 19 - 23, 2021*, pages 1000–1008. Association for Computational Linguistics.
- 402 Patrick S. H. Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra
403 Piktus, Pontus Stenetorp, and Sebastian Riedel. 2021b. <http://arxiv.org/abs/2102.07033> PAQ: 65
404 million probably-asked questions and what you can do with them. *CoRR*, abs/2102.07033.
- 405 Chen Liang, Jonathan Berant, Quoc V. Le, Kenneth D. Forbus, and Ni Lao. 2017.
406 <https://doi.org/10.18653/v1/P17-1003> Neural symbolic machines: Learning semantic parsers
407 on freebase with weak supervision. In *Proceedings of the 55th Annual Meeting of the Association*
408 *for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long*
409 *Papers*, pages 23–33. Association for Computational Linguistics.
- 410 Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. [https://doi.org/10.18653/v1/D19-](https://doi.org/10.18653/v1/D19-1282)
411 [1282](https://doi.org/10.18653/v1/D19-1282) Kagnet: Knowledge-aware graph networks for commonsense reasoning. In *Proceedings*
412 *of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*
413 *International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong*
414 *Kong, China, November 3-7, 2019*, pages 2829–2839. Association for Computational Linguistics.
- 415 Bill Yuchen Lin, Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Xiang Ren, and William W. Cohen.
416 2020. <http://arxiv.org/abs/2010.14439> Differentiable open-ended commonsense reasoning. *CoRR*,
417 abs/2010.14439.
- 418 Ye Liu, Yao Wan, Lifang He, Hao Peng, and Philip S. Yu. 2021.
419 <https://ojs.aaai.org/index.php/AAAI/article/view/16796> KG-BART: knowledge graph-augmented
420 BART for generative commonsense reasoning. In *Thirty-Fifth AAAI Conference on Artificial*
421 *Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial*
422 *Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial*
423 *Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 6418–6425. AAAI Press.
- 424 Sewon Min, Danqi Chen, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2019.
425 <http://arxiv.org/abs/1911.03868> Knowledge guided text retrieval and reading for open
426 domain question answering. *CoRR*, abs/1911.03868.
- 427 Nina Pörner, Ulli Waltinger, and Hinrich Schütze. 2020. [https://doi.org/10.18653/v1/2020.findings-](https://doi.org/10.18653/v1/2020.findings-emnlp.71)
428 [emnlp.71](https://doi.org/10.18653/v1/2020.findings-emnlp.71) E-BERT: efficient-yet-effective entity embeddings for BERT. In *Findings of the Association*
429 *for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume
430 EMNLP 2020 of *Findings of ACL*, pages 803–818. Association for Computational Linguistics.
- 431 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
432 Zhou, Wei Li, and Peter J. Liu. 2020. <http://jmlr.org/papers/v21/20-074.html> Exploring the limits
433 of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- 434 Siva Reddy, Oscar Täckström, Michael Collins, Tom Kwiatkowski, Dipanjan Das, Mark Steedman,
435 and Mirella Lapata. 2016. <https://tacl2013.cs.columbia.edu/ojs/index.php/tacl/article/view/807>
436 Transforming dependency structures to logical forms for semantic parsing. *Trans. Assoc. Comput.*
437 *Linguistics*, 4:127–140.
- 438 Hongyu Ren, Weihua Hu, and Jure Leskovec. 2020. <https://openreview.net/forum?id=BJgr4kSFDS>
439 Query2box: Reasoning over knowledge graphs in vector space using box embeddings. In *8th*
440 *International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April*
441 *26-30, 2020*. OpenReview.net.
- 442 Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. [https://doi.org/10.18653/v1/2020.emnlp-](https://doi.org/10.18653/v1/2020.emnlp-main.437)
443 [main.437](https://doi.org/10.18653/v1/2020.emnlp-main.437) How much knowledge can you pack into the parameters of a language model? In
444 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing,*
445 *EMNLP 2020, Online, November 16-20, 2020*, pages 5418–5426. Association for Computational
446 Linguistics.

- 447 Devendra Singh Sachan, Mostofa Patwary, Mohammad Shoeybi, Neel Kant, Wei Ping, William L.
448 Hamilton, and Bryan Catanzaro. 2021a. End-to-end training of neural retrievers for open-domain
449 question answering. In *Proceedings of the 59th Annual Meeting of the Association for Computa-*
450 *tional Linguistics and the 11th International Joint Conference on Natural Language Processing,*
451 *ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 6648–6662.
452 Association for Computational Linguistics.
- 453 Devendra Singh Sachan, Siva Reddy, William L. Hamilton, Chris Dyer, and Dani Yogatama.
454 2021b. [https://proceedings.neurips.cc/paper/2021/hash/da3fde159d754a2555eaa198d2d105b2-](https://proceedings.neurips.cc/paper/2021/hash/da3fde159d754a2555eaa198d2d105b2-Abstract.html)
455 [Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/da3fde159d754a2555eaa198d2d105b2-Abstract.html) End-to-end training of multi-document reader and retriever for open-domain question
456 answering. In *Advances in Neural Information Processing Systems 34: Annual Conference on*
457 *Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages
458 25968–25981.
- 459 Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee, and Danqi Chen. 2021.
460 <https://doi.org/10.18653/v1/2021.emnlp-main.496> Simple entity-centric questions chal-
461 lenge dense retrievers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural*
462 *Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11*
463 *November, 2021*, pages 6138–6148. Association for Computational Linguistics.
- 464 Haitian Sun, Tania Bedrax-Weiss, and William W. Cohen. 2019. [https://doi.org/10.18653/v1/D19-](https://doi.org/10.18653/v1/D19-1242)
465 [1242](https://doi.org/10.18653/v1/D19-1242) Pullnet: Open domain question answering with iterative retrieval on knowledge bases and text.
466 In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
467 *and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP*
468 *2019, Hong Kong, China, November 3-7, 2019*, pages 2380–2390. Association for Computational
469 Linguistics.
- 470 Alon Talmor and Jonathan Berant. 2018. <https://doi.org/10.18653/v1/n18-1059> The web as a
471 knowledge-base for answering complex questions. In *Proceedings of the 2018 Conference of the*
472 *North American Chapter of the Association for Computational Linguistics: Human Language*
473 *Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long*
474 *Papers)*, pages 641–651. Association for Computational Linguistics.
- 475 Pat Verga, Haitian Sun, Livio Baldini Soares, and William W. Cohen. 2021.
476 <https://doi.org/10.18653/v1/2021.naacl-main.288> Adaptable and interpretable neural memory over
477 symbolic knowledge. In *Proceedings of the 2021 Conference of the North American Chapter of*
478 *the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*
479 *2021, Online, June 6-11, 2021*, pages 3678–3691. Association for Computational Linguistics.
- 480 Denny Vrandečić and Markus Krötzsch. 2014. <https://doi.org/10.1145/2629489> Wikidata: a free
481 collaborative knowledgebase. *Commun. ACM*, 57(10):78–85.
- 482 Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2019.
483 <http://arxiv.org/abs/1911.06136> KEPLER: A unified model for knowledge embedding and pre-
484 trained language representation. *CoRR*, abs/1911.06136.
- 485 Wenhan Xiong, Thien Hoang, and William Yang Wang. 2017. <https://doi.org/10.18653/v1/d17-1060>
486 Deeppath: A reinforcement learning method for knowledge graph reasoning. In *Proceedings*
487 *of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017,*
488 *Copenhagen, Denmark, September 9-11, 2017*, pages 564–573. Association for Computational
489 Linguistics.
- 490 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V.
491 Le. 2019. [https://proceedings.neurips.cc/paper/2019/hash/dc6a7e655d7e5840e66733e9ee67cc69-](https://proceedings.neurips.cc/paper/2019/hash/dc6a7e655d7e5840e66733e9ee67cc69-Abstract.html)
492 [Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/dc6a7e655d7e5840e66733e9ee67cc69-Abstract.html) Xlnet: Generalized autoregressive pretraining for language understanding. In
493 *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information*
494 *Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages
495 5754–5764.
- 496 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov,
497 and Christopher D. Manning. 2018. <https://doi.org/10.18653/v1/d18-1259> Hotpotqa: A dataset

- 498 for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on*
499 *Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4,*
500 *2018*, pages 2369–2380. Association for Computational Linguistics.
- 501 Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021.
502 <https://doi.org/10.18653/v1/2021.naacl-main.45> QA-GNN: reasoning with language models and
503 knowledge graphs for question answering. In *Proceedings of the 2021 Conference of the North*
504 *American Chapter of the Association for Computational Linguistics: Human Language Technolo-*
505 *gies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 535–546. Association for Computational
506 Linguistics.
- 507 Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. 2015.
508 <https://doi.org/10.3115/v1/p15-1128> Semantic parsing via staged query graph generation:
509 Question answering with knowledge base. In *Proceedings of the 53rd Annual Meeting of the*
510 *Association for Computational Linguistics and the 7th International Joint Conference on Natural*
511 *Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July*
512 *26-31, 2015, Beijing, China, Volume 1: Long Papers*, pages 1321–1331. The Association for
513 Computer Linguistics.
- 514 Donghan Yu, Chenguang Zhu, Yuwei Fang, Wenhao Yu, Shuohang Wang, Yichong Xu, Xiang
515 Ren, Yiming Yang, and Michael Zeng. 2022a. <https://aclanthology.org/2022.acl-long.340> Kg-
516 fid: Infusing knowledge graph in fusion-in-decoder for open-domain question answering. In
517 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume*
518 *1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 4961–4974. Association for
519 Computational Linguistics.
- 520 Donghan Yu, Chenguang Zhu, Yiming Yang, and Michael Zeng. 2022b.
521 <http://arxiv.org/abs/2010.00796> JAKET: joint pre-training of knowledge graph and language
522 understanding. Conference on Artificial Intelligence, AAAI.
- 523 Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D.
524 Manning, and Jure Leskovec. 2022. <http://arxiv.org/abs/2201.08860> Greaselm: Graph reasoning
525 enhanced language models for question answering. *CoRR*, abs/2201.08860.
- 526 Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019.
527 <https://doi.org/10.18653/v1/p19-1139> ERNIE: enhanced language representation with informative
528 entities. In *Proceedings of the 57th Conference of the Association for Computational Linguistics,*
529 *ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 1441–1451.
530 Association for Computational Linguistics.
- 531 Victor Zhong, Caiming Xiong, and Richard Socher. 2017. <http://arxiv.org/abs/1709.00103> Seq2sql:
532 Generating structured queries from natural language using reinforcement learning. *CoRR*,
533 abs/1709.00103.

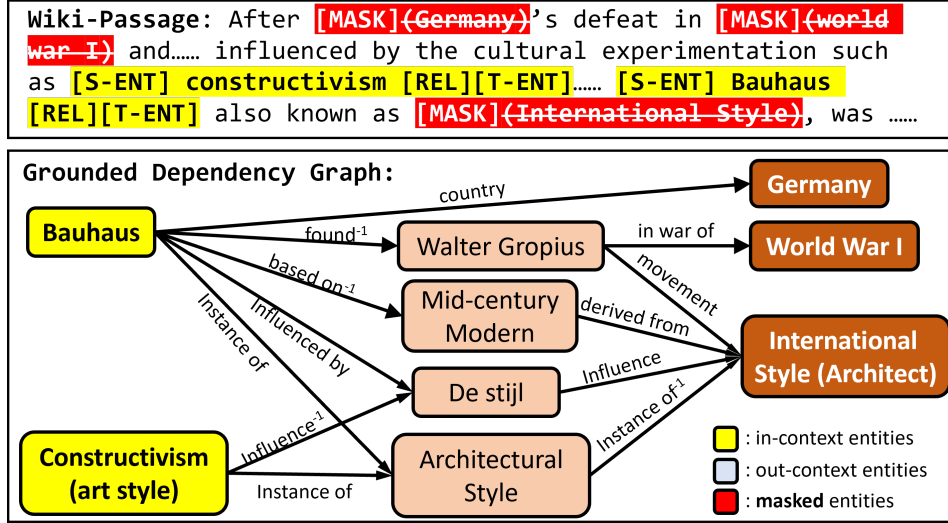


Figure 4: Pre-training sample w/ golden reasoning path. More real examples are shown in Table 8 in Appendix.

Name	Number	dimension	#param (M)
Number of Entity	4,947,397	128	633
Number of Relation	2,008	768	1.5
Number of Edges	45,217,947	-	47

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

534 A Related Work

535 To encode knowledge (significantly smaller than the web corpus) as *memory* into LM parameter, a
 536 line of works try compressed knowledge including QA pairs (Chen et al., 2022; Lewis et al., 2021b),
 537 entity embedding (Férvy et al., 2020) and reasoning cases (Das et al., 2021, 2022). There’s also
 538 several works utilizing Knowledge Graph(KG). FILM (Verga et al., 2021) turns KG triplets into
 539 memory. Given a question, LM retrieves most relevant triplet as answer. GreaseLM (Zhang et al.,
 540 2022) propose to interact LM with KG via a interaction node. JAKET (Yu et al., 2022b) encode text
 541 and KG independently and fuse information at late stage. We introduce and discuss with other related
 542 works in Sec. D in Appendix.

543 B Implementation Details

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

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

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

Dataset	Train	Dev	Test
Natural Questions	58880	8757, 3610	
Trivia QA	60413	8837	11313
Web Questions	2474	361	2032
HotpotQA			
Complex WebQ	27623	3518	3531
WebQ-SP (Wiki-answerable)	1388	153	841
FreebaseQA (Wiki-answerable)	12535	2464	2440

Table 5: Dataset Train/Valid/Test splits.

Models	#param (B)	WQ-SP	TQA
RoBERTa (Base)	0.12	47.5	40.3
+ OREOLM ($T=1$)	0.12 + 0.68	89.7	61.4
+ OREOLM ($T=2$)	0.13 + 0.68	92.4	66.8

Table 6: **Closed-Book Entity Prediction** validation performance of Encoder RoBERTa on WikiData-Answerable Dataset.

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 = ll_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

```

556 **Hyperparameters** In this work, we don't have too much hyperparameters to be tuned, as most
557 parameters as well as optimizing setting of LM is fixed. Our random walk part is non-parametric.
558 The only tunable hyperparameter is hidden dimension size. We simply choose one setting, which is
559 128 for entity embedding, and 768 for relation embedding. The former is because entity is super large
560 (over 5M), so we use a relatively smaller dimension size. Detailed statistics about wikidata memory
561 is in Table 4.

562 C Dataset Details

563 Below shows details for each dataset, and the detailed dataset split is shown in Figure 5

564 **Natural Questions** Kwiatkowski et al. (2019) contains questions from Google search queries, and
565 the answers are text spans in Wikipedia. We report short answer Exact Match (EM) performance.
566 The open version of this dataset is obtained by discarding answers with more than 5 tokens.

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

569 **TriviaQA** Joshi et al. (2017) contains trivia questions and answers are text spans from the Web. We
570 report Exact Match (EM) performance. We use its unfiltered version for evaluation.

571 **HotpotQA** Yang et al. (2018) is a multi-hop QA dataset. There are two evaluation settings. In the
572 *distractor setting*, 10 candidate paragraphs are provided for each question, of which there are two

Models	#param	NQ	WQ	TQA	ComplexWQ	HotpotQA
T5 (Large)	0.74B	-	-	-	-	-
+ OREOLM ($T=2$)	0.76B + 0.68B	33.6	38.9	42.7	29.6	35.5

Table 7: Closed-Book Generative QA validation performance of T5.

573 golden paragraphs. In the *full-wiki setting*, a model is required to extract paragraphs from the entire
574 Wikipedia. We report Exact Match (EM) on full-wiki setting.

575 **Complex WebQuestions** Talmor and Berant (2018) is a dataset that composite simple one-hot
576 questions in WebQuestionsSP by extending entities or adding constraints, so that each question
577 requires complex reasoning to solve.

578 **WebQuestionsSP** Yih et al. (2015) is annotated dataset from WebQuestions, such taht each quetsion
579 is answerable using Freebase via a SQL query.

580 D Other Related Works

581 D.1 Introduce other related works

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

591 **Knowledge-Base Question Answering** Traditional parsing-based methods parse the question into
592 some intermediate query (e.g., SQL language, query graphs), which can execute on a knowledge base
593 to get answer (Berant et al., 2013; Yih et al., 2015; Reddy et al., 2016; Zhong et al., 2017; Liang
594 et al., 2017). However, existing knowledge bases suffer from low coverage of entities and relations
595 required for open-ended questions. As an alternative, several works try to incorporate the structured
596 knowledge into neural QA models for differentiable reasoning. Lin et al. (2019) and Feng et al.
597 (2020) parse the question into a sub-graph of knowledge base, and apply graph neural networks as
598 reasoner to extract answers. Chen et al. (2020) integrates general symbolic operations as basic units,
599 and parse questions into compositional programs to answer general questions.

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

610 D.2 Discussion with Previous Works

611 **Compare with FILM** Though FILM has the advantage of end-to-end training and easily modifica-
612 tion of knowledge memory, it simply stacks \mathcal{KG} module on top of LM without interaction, and can
613 only handle one-hop relational query that is answerable by \mathcal{KG} . Our approach, OREOLM, follows the

614 same *memory* idea by encoding \mathcal{KG} into LM parameter, and we desire LM and \mathcal{KG} reasoning module
615 could interact and collaboratively improve each other.

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

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

631 D.3 Discussion of Graph Walking-based Reasoning vs Graph Neural Networks

632 Recently, Graph Neural Networks (GNNs) have shown superior performance for structured represen-
633 tation learning. There’s also a lot of works trying to use GNNs for Question Answering (Yasunaga
634 et al., 2021; Zhang et al., 2022). The one that has very similar motivation with us is GreaseLM.
635 Therefore, a natural question is, whether could we use GNN instead of the non-parametric random
636 walk module, for ODQA?

637 To answer this question, let’s consider a simplest setup of GNN. We could identify initial entities,
638 connected them via a k-hop subgraph, and encode graph with text (Zhang et al., 2022) or indepen-
639 dently (Yu et al., 2022b). When we want to retrieve knowledge from graph to LM, normally we just
640 take the contextualized node embedding as input for knowledge fusion.

641 In this setup, say the answer is K -hop away from an initial entity, the ground-truth reasoning path
642 is $e_0, r_1, e_1, r_2, \dots, e_{k-1}, r_k, e_k = a$. Using our method, we first predict r_1 , transit to e_1 , and step
643 by step conduct reasoning via walking. However, if we use GNN’s final embedding, it requires to
644 pass information from neighbor to itself. Therefore, suppose we have a K -layer GNN, the first step
645 should be identify r_k , and pass information from answer $e_k = a$ to e_{k-1} . This is conter-intuitive
646 as we normally cannot assume to know the answer, nor knowing the last step to reach the answer.
647 In situations where all candidate answer is given, like CommonSenseQA, where GreaseLM mainly
648 works on, this problem is less harmful as it’s guaranteed to contain the answer in a restricted small
649 graph. However, in open-domain setup, we need to try best to narrow down the search space by
650 following the forward reasoning instead of the backward manner. Therefore, in this work we adopt
651 walking-based reasoning.

652 E Illustration of Pre-Trained Data and Reasoning Paths

653 The pre-training samples and reasoning paths (generated by T5-large on NQ dataset) is shown from
654 Table 8-11.

655 F Ablation Studies

656 We conduct several ablation studies to evaluate which model design indeed contributes to the model.
657 As shown in the bottom blocks in Table 2, we first remove the \mathcal{KG} reasoning component and provide
658 RoBERTa base model via concatenated KB triplets and train such a model using \mathcal{L}_{SSM} over the
659 same WikiDataset. Such a model’s results are close to the KEPLER results but much lower than
660 other models with explicit knowledge memory. We further investigate the role of pre-training tasks.
661 Without pre-training, the OREOLM only performs slightly better than RoBERTa baseline, due to the
662 cold-start problem of entity and relation embedding. We further show that removing \mathcal{L}_{ent} and \mathcal{L}_{ent}

663 could significantly influence final performance. The current combination is the best choice to train
664 OREOLM to reason.

665 G Limitations

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

670 Though we didn't spend much time exploring before the paper submission, we indeed try using
671 $T = 3$, but currently it didn't get better results. We hypothesize the following reasons: 1) A large
672 portion of our current model's improvement relies on the weakly supervised relation pre-training. To
673 do it, we construct a K-hop ($K=2$ now) subgraph, and sample dependency graph based on it. The
674 larger K we choose, the more noise is included into the generated relation label, in an exponential
675 increasing speed. Thus, it's harder to get accurate reasoning path ground-truth for high-order T .
676 Another potential reason is that within Transformer model, the representation space in lower and
677 upper layer might be very different, say, encode more syntax and surface knowledge at lower layers,
678 while more semantic knowledge at upper layers. Currently we adopt a MLP projection head, wishing
679 to map integrated knowledge into the same space, but it might have many flaws and need further
680 improvement.

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

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

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

691 **Unreasonable Assumption of Path In-dependency** When we derive equation 1, we have the
692 assumption that reasoning paths starting from different entities should be independent. This is not
693 always correct, especially for questions that require logical reasoning, say, have conjunction or
694 disjunction operation over each entity state. And thus our current methods might not work for those
695 complex QA with logical dependencies.

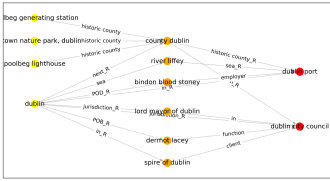
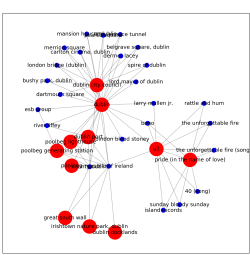
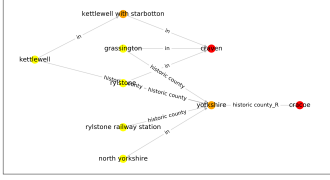
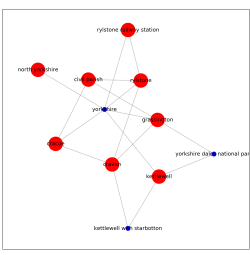
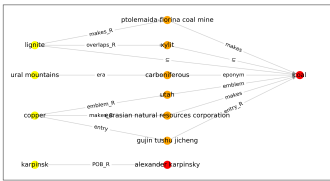
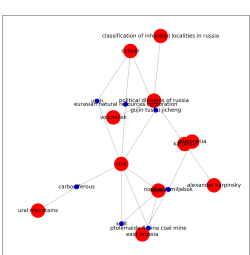
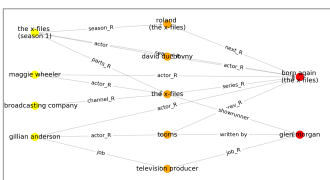
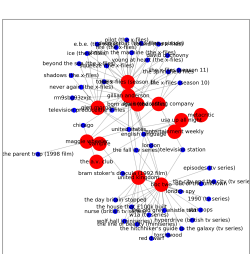
Title	Masked Text	Ground Truth	Dependency Graph	2-Hop Graph
Poolbeg	the lighthouse was [mask] [s-ent] [mask] [rel] [t-ent] completed in 1795. overview. the [s-ent] poolbeg[rel] [t-ent] "peninsula" is home to a number of landmarks including the [s-ent] [mask][rel] [t-ent] , the [s-ent] pool[mask] lighthouse[rel] [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][ington][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', ' egal', ' metacritic', ' as "wretched", ' fact that', ' background noise for a', ' heavy-handed attempts at', ' glen morgan', ' doing an episode on']		

Table 8: Example of Pre-training data points (Part 1).

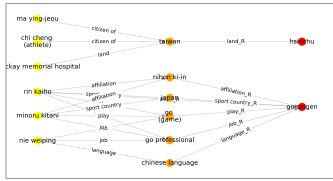
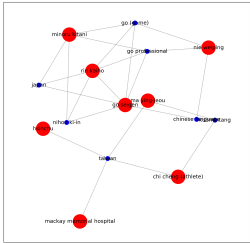
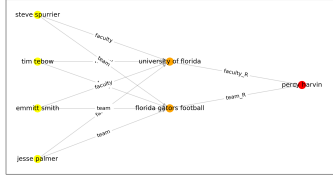
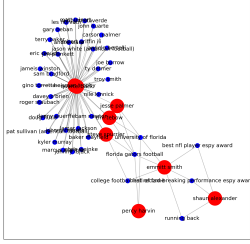
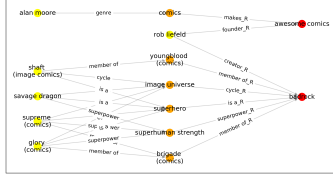
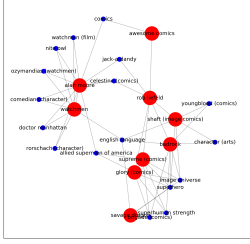
Title	Masked Text	Ground Truth	Dependency Graph	K-Hop Graph
Shen Chunshan	his memoirs, he suffered his second stroke[mask][mask], even after his second stroke, he continued writing; his series of biographies of five go masters [s-ent] [mask][mask][mask][rel] [t-ent] , [s-ent] minoru kit[mask][rel] [t-ent]	[. 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] [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	[' 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.', ' ators', ' touchdown']		
Judgment Day (Awesome 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 toms/ [s-ent] knightsabre[rel] [t-ent].....	[' 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 9: 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 10: 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 11: Example of QA prediction with reasoning path on NQ (part 2).