Relevant CommonSense Subgraphs for "What if..." Procedural Reasoning

Anonymous ACL submission

Abstract

This work deals with the challenge of learning causal reasoning over procedural text to answer "What if..." questions when external commonsense knowledge is required. We propose a novel multi-hop graph reasoning model to 1) efficiently extract a commonsense subgraph with the most relevant information from a large knowledge graph; 2) predict the causal answer by reasoning over the representations obtained from the commonsense subgraph and the contextual interactions between the questions and context. We evaluate our model on WIQA dataset and achieve state-of-the-art performance compared to the recent models.

1 Introduction

011

014

016

017

022

037

In recent years, large-scale pre-trained language models (LMs) have made a breakthrough progress and demonstrate a high performance in many NLP tasks, including procedural text reasoning (Tandon et al., 2019; Rajagopal et al., 2020). However, since the knowledge of language that is present in the corpora is learnt only implicitly by LMs, they cannot provide explainable predictions. In some cases, the knowledge contained in a given text is sufficient to predict the answer, as it is shown in the question 1 of Figure 1. This knowledge is directly encoded and learned by LMs models (Asai and Hajishirzi, 2020; Tandon et al., 2019). However, there are many cases in which the required knowledge is not included in the procedural text itself. For example, for the question 2 in Figure 1, the information about the "nutrient" on the seeds does not exist in the procedural text. Therefore, the external commonsense knowledge is required.

There are several existing resources that contain world knowledge and commonsense. Examples are knowledge graphs (KGs) like ConceptNet (Speer et al., 2017) and ATOMIC (Sap et al., 2019). Looking back at the question 2, we can see that through providing the external knowledge triplets (nutrient,



Figure 1: WIQA task contains procedural text, and different types of questions. The bold choices are the gold answers.

relatedto, soil) and (soil, relatedto, seed) from ConceptNet, we can build an explicit reasoning chain and choose an explainable answer. 041

042

045

047

055

057

059

060

061

063

064

065

067

Two challenges exist in procedural text reasoning and using external KBs. The first challenge is effectively extracting the most relevant external information and reducing the noise from the KB. The second challenge is reasoning over the extracted knowledge. Several works enhance the QA model with commonsense knowledge (Lin et al., 2019; Lv et al., 2020). However, the noisy knowledge from KG will seriously mislead the QA model to predict the answer. Moreover, using KBs is often investigated in the tasks that perform QA directly over KB itself, such as CommonsenseQA (Talmor et al., 2019), etc. There are less sophisticated techniques proposed for using KB explicitly (i.e. not through training LMs) in reading comprehension for aiding QA over text. REM-Net (Huang et al., 2021) is the only work that uses commonsense for WIQA and uses a memory network to extract the external triplets to solve the first challenge. However, this work has no reasoning process over the extracted knowledge and uses a simple multi-head operator to predict the answer. EIGEN (Madaan et al., 2020) constructs an influence graph to find the chain of reasoning given procedural text. However, EIGEN



Figure 2: MRRG Model is composed of Candidate Triplet Extraction, KG Attention, Commonsense Subgraph Construction, Text encoder with contextual interaction, Graph Reasoning, and Answer prediction modules.

cannot deal with the challenge that the required knowledge is not in the given document.

068

071

077

085

093

099

100

101

102

103

104

105

To solve these two challenges, we propose a Multi-hop Reasoning network over Relevant CommonSense SubGraphs (MRRG) for casual reasoning over procedural Text. Our motivation is to effectively and efficiently extract the most relevant information from a large KG to help procedural reasoning. First, we extract the entities, retrieve related external triplets from KG, and learn to extract the most relevant triplets to a given text input by a novel KG attention mechanism. Then, we construct a commonsense subgraph based on the extracted triplets in a pipeline. We use the extracted subgraphs as a part of end-to-end QA model to help in filling the knowledge gaps in the procedure and performing multi-hop reasoning. The final model predicts the causal answer by reasoning over the contextual interaction representations over the question and the document and learning graph representations over the KB subgraphs. We evaluate our MRRG on the "what if" WIQA benchmark. MRRG model achieves SOTA and brings significant improvements compared to existing baselines.

The contributions of our work are: 1) We train a separate module that extracts the relevant parts of the KB to avoid the noisy and inefficient usage of the information in large KBs. 2) We design an end-to-end model that uses the extracted QA-dependent KB as a subgraph to guide the reasoning over the procedural text to answer the questions. 3) Our MRRG achieves SOTA on the WIQA benchmark.

2 Model Description

Problem Formulation: Formally, the problem is to predict an answer *a* from a set of pre-defined answers given input question *q*, a document Cwhich is composed of several sentences $C = \{s_1, \ldots, s_n\}$, and a large knowledge graph KG.

106 Overview of MRRG Model: Figure 2 shows the107 proposed architecture. (1) We extract the enti-

ties from question and context in preprocessing step and use them to retrieve the set of **candidate triples** from the ConceptNet. (2) We train the **KG Attention** module to extract the most relevant triplets and reduce the noisy concepts from candidate triplets. (3) We augment the **commonsense subgraph** based on the relevant triplets. (4) We train a model that uses two components, the commonsense subgraph as a relational graph network and a text encoder including question and document to do **procedural reasoning**. Below, we describe the details of each module. 108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

(1) Candidate Triplet Extraction from KG: Given the input q and C, we extract the contextual entities (concepts) by a open Information Extraction (OpenIE) model (Stanovsky et al., 2018). For each extracted entity t_{in} , we retrieve the relational triplets $t = (t_{in}, r, t_{out})$ from KG, where t_{out} is the concept taken from ConceptNet and r is a semantic relation type. We then apply a pre-trained Language Model, RoBERTa, to obtain the representation of each triplet $E^t = f_{LM}([t_{in}, r, t_{out}]) \in \mathbb{R}^{3 \times d}$, where f_{LM} denotes the language model operation and the triplets are given as a sequence of concepts and relation to the LM.

(2) KG Attention: The KG attention module is shown in Figure 2-A. We concatenate q and C to form Q = [[CLS]; q; [SEP]; C], where [CLS] and [SEP] are special tokens in the LMs (Liu et al., 2019). We encode Q by RoBERTa and get the contextual token representations that include $E_{[CLS]}$, E_q , and E_C .

Given triplet E^t , we generate a context-triplet pair $E_z^t = [E_{[CLS]}; E_{in}^t; E_r^t; E_{out}^t]$. Afterwards, we compute context-triplet pair attention and a softmax layer to output the Context-Triplet pairwise importance Score $CTS_t = \frac{\exp(MLP(E_z^t))}{\sum_{j=1}^m \exp(MLP(E_z^t))}$.

Then we choose the top-k relevant triplets with the top CTS scores and then use the relevant triplets to construct the subgraph. For each selected triplet, we obtain the triplet representa149 tion $E'^t = [E'_{in}^t, E_r^t, E'_{out}] \in \mathbb{R}^{3 \times d}$, where 150 $E'_{in}^t = f_{in}([CTS_t \cdot E_{in}^t; CTS_t \cdot E_r^t])$ and $E'_{out}^t = f_{out}([CTS_t \cdot E_{out}^t; CTS_t \cdot E_r^t])$. Notice that f_{in} 152 and f_{out} are MLP layers.

(3) Commonsense Subgraph Construction: We construct the subgraph G_s based on the relevant triplets from KG attention for each question and answer pair. We add more edges to the subgraph as follows: Two entities in the triplets will have an edge if a relation r in the KG exists between them. The assumption is that the augmented commonsense subgraph will contain the reasoning paths. We use E_{in}^{tt} and E_{out}^{tt} for the KG subgraph initial node representation $h^{(0)}$.

(4) **Procedural Reasoning** composes of two parts. (1) Multi-Hop Graph Reasoning: this is the Graph Reasoning part of Figure 2-B. Given the subgraph G_s , we use RGCN (Schlichtkrull et al., 2018) to learn the representations of the relational graph. RGCN learns graph representations by aggregating messages from its direct neighbors and relational semantic edges. The (l+1)-th layer node representation $h_i^{(l+1)}$ is updated based on the neighborhood node representations h_j^l from the *l*-layer multiplied by the relational matrices $W_{r1}^{(l)}, \ldots, W_{r|R|}^{(l)}$. The representation $h_i^{(l+1)}$ is computed as follows:

$$h_i^{(l+1)} = \sigma(\sum_{r \in \mathcal{R}} \sum_{j \in N_i^r} \frac{1}{|N_i^r|} W_r^{(l)} h_j^{(l)} + W_0^{(l)} h_i^{(l)}),$$

where σ denotes a non-linear activation function, N_i^r represents a set that is includes neighbor indices of node i under semantic relation r. Finally, we obtain the E_{G_s} after several hops message passing. (II) Text Contextual Interaction Encoder: We have obtained the contextual token representations $E_{[CLS]}$, E_q , and E_c in the KG attention module. Followed by Seo et al., we utilize Bi-DAF style contextual interaction module to feed E_q and E_c to Context-to-Question Attention $E_{\mathcal{C} \to q}$ = $softmax(sim(E_q^T, E_{\mathcal{C}}))E_q$ and Question-to-Context Attention $E_{q\to C}$ to obtain the contextual interaction between question and context. More details about obtaining $E_{q\to C}$ and $E_{C\to q}$ are described in Appendx A.3. Then we use LSTM to obtain the hidden state representation:

$$F_{q \to \mathcal{C}} = LSTM(E_{q \to \mathcal{C}}), F_{\mathcal{C} \to q} = LSTM(E_{\mathcal{C} \to q}).$$

(III) Answer Prediction: We concatenate $E_{[CLS]}$, $F_{q \to C}$, $F_{C \to q}$, and the compact subgraph representation E'_{G_s} obtained from attentive pooling, and use it as the final representation F =

 $[E_{[CLS]}; F_{q \to C}; F_{C \to q}; E'_{G_s}]$. We utilize a classifier MLP (F) to predict the answer. Our MRRG has two separate training modules used in a pipeline for triplet selection and procedural reasoning.

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

(I) Training KG Attention Triplet Selection: Figure 4 and the left block of Figure 2 shows the triplet selection model. The architecture is the same as KG attention except adding extra 3 MLP layers for the concatenation of $[E_{[CLS]}; E_q; E_{\mathcal{C}}; E'_1^t; \ldots; E'_k^t]$ to predict the answer. We use the cross-entropy as the loss function to train the model.

(II) Training End-to-End MRRG: After we pretrain the KG attention, we keep the learned parameters and extract the most relevant concepts and construct the multi-relational commonsense subgraph G_s . We combine subgraph representation and text interaction representation as input to train the answer prediction module by cross-entropy loss.

Models	in-para	out-of-para	no-effect	Test V1 Acc
Majority	45.46	49.47	55.0	30.66
Polarity	76.31	53.59	27.0	39.43
Adaboost (Freund and Schapire, 1995)	49.41	36.61	48.42	43.93
emphDecomp-Attn (Parikh et al., 2016)	56.31	48.56	73.42	59.48
BERT (no para) (Devlin et al., 2019)	60.32	43.74	84.18	62.41
BERT (Tandon et al., 2019)	79.68	56.13	89.38	73.80
RoBERTa (Tandon et al., 2019)	74.55	61.29	89.47	74.77
EIGEN (Madaan et al., 2020)	73.58	64.04	90.84	76.92
REM-Net (Huang et al., 2021)	75.67	67.98	87.65	77.56
Logic-Guided (Asai and Hajishirzi, 2020)	-	-	-	78.50
RoBERTa+KG Attention Triplet Selection	72.21	64.60	89.13	75.22
MRRG	79.85	69.93	91.02	80.06
Human	-	-	-	96.33

Table 1: Model Comparisons on WIQA test V1 dataset.

3 Experiments and Results

We implemented our MRRG framework using Py-Torch¹. We use a pre-trained RoBERTa (Liu et al., 2019) to encode the input. The maximum number of selected triplets from ConceptNet is 50. More details are shown in the Appendix A.1.

Datasets: WIQA is a large dataset for "what if" causal reasoning. WIQA contains three types of questions: 1) the questions can be directly answered based on the text, called in-paragraph questions. 2) the questions require external knowledge to be answered, called out-of-paragraph questions, and 3) irrelevant causes and effects, called no-effect questions. WIQA contains 29808 training samples, 6894 development samples, 3993 test samples (test V1), and 3003 test samples (test V2).

Results: Table 1 and Table 2 show the performance of MRRG on the WIQA task compared to other baselines. We show the baseline descriptions in

167

168

169

170

163

153

154

155

156

157

158

159

160

161

162

¹Our code will be available after the paper is accepted.

Question and Document Content	Extracting Triplets	85
Question: suppose the soil is rich in nutrients happens, how will it affect more seeds are produced. Content: ["A plant produces a seed", "The seed falls to the ground", "The seed is buried", "The seed germinates", "A plant grows", "The plant produces flowers", "The flowers produce more seeds."] Gold Answer: More	(nutrient, relatedto, soil) (soil, relatedto, seed)	
Question: suppose more land available happens, how will it affect less igneous rock forming. Content: ["Different kinds of rocks melt into magma", "Magma cools in the crust", "Magma goes to the surface and becomes lava", "Lava cools", "Cooled magma and lava become igneous rock."] Gold Answer: Less	(igneous rock, isa, rock) (land, relatedto, rock) (land, relatedto, surface) (surface, relatedto, igneous rock)	60 55 50 BERT REBRTA EGEN MARG

Figure 3: Left: Case study of our MRRG Framework. Right: Comparing the results over different number of hops.

Models	in-para	out-of-para	no-effect	Test v2 Acc
Random	33.33	33.33	33.33	33.33
Majority	00.00	00.00	100.0	41.80
BERT	70.57	58.54	91.08	74.26
RoBERTa	70.69	60.20	91.11	75.34
REM-Net	70.94	63.22	91.24	76.29
REM-Net (RoBERTa-large)	76.23	69.13	92.35	80.09
QUARTET (RoBERTa-large)	74.49	65.65	95.30	82.07
(Rajagopal et al., 2020)				
RGN (Zheng and Kordjamshidi, 2021)	75.91	66.15	92.12	79.95
RoBERTa+KG Attention Triplet Selection	70.02	62.30	91.23	75.86
MRRG	76.80	67.83	92.28	80.39
MRRG (RoBERTa-large)	79.12	71.10	93.53	83.46
Human	-	-	-	96.30

Table 2: Model Comparisons on WIQA test V2 dataset.

Appendix A.2. First, our KG Attention triplet selection model outperforms the RoBERTa and has 3.3% improvement on the out-of-para category. Second, our MRRG achieves SOTA result compared to all baseline models. Our MRRG improves the SOTA by 3.21% for the in-para questions and improves it by 3.95% for the out-of-para questions.

4 Analysis

209

210

211

212

213

214

215

216

Effects of Using External Knowledge: In the 217 WIQA, all the baseline models achieve significantly 218 lower accuracy in the out-of-para than in-para and 219 no-effect categories. MRRG achieves SOTA in out-of-para category because of using the highly relevant commonsense subgraphs and the combi-222 nation of reasoning over text interaction and the 223 graph modules. As is shown in table 2, the advantage of the MRRG model is reflected on out-of-para questions. MRRG improves 4.61% over REM-Net. Notice that REM-Net is the only model that utilized external knowledge on WIQA. Figure 3 shows a 228 case in which the "soil" and "nutrient" only appear 230 in the question and do not exist in the text. The baseline models fail to answer this out-of-para question due to missing external knowledge. However, our model predicts the correct answer by explicitly incorporating the (nutrient, related to, soil), (soil, relatedto, seed) that connects the critical information 235 between the question and document.

Relational Reasoning and Multi-Hops: Both inpara and out-of-para question types require multiple hops of reasoning to find the answer in the WIQA. As shown in the right side of Figure 3, our MRRG model accuracy improved 3.2% for 1 hop, 7.3% for 2 hops, and 3.9% for 3 hops compared

Ablation	Model	Dev Acc
Text only	RoBERTa-base	75.51%
Text only	+ contextual interaction	76.85%
Text only	KG Attention Triplet Selection	77.39%
	 semantic relation 	78.31%
	GNN dim=50	79.18%
Text+Graph	GNN dim=100	80.30%
	GNN dim=200	79.88%

Table 3: Ablation and hyper-para. choices on WIQA. "GNN dim" is the dimension of graph representation.

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

to EIGEN. MRRG made a sharp improvement in reasoning with multiple hops due to the relational graph reasoning and the effectiveness of the extracted commonsense subgraph. We study some cases to analyze the multi-hop reasoning and the reasoning chains. In the second case in Figure 3, the extracted relevant triplets (land, relatedto, surface), (surface, relatedto, igneous rock) construct a two-hop reasoning chain "land-surface-jgneous rock" that helps MRRG to find the correct answer. Ablation Study: Table 3 shows the ablation study of MRRG using WIQA. Firstly, we remove the commonsense subgraph and graph network. The accuracy decreases 3.4% compared to MRRG. Second, we remove the contextual interaction module and the accuracy decreases 1.3%. In an additional experiment, we use the KG attention triplet selection module to directly predict the answer without the pipeline of constructing the subgraph and using the graph reasoning module. We show the result as KG Attention Triplet Selection in Table 3. The result shows that using the selected triplets representations chosen from KG attention is helpful for the WIQA task. However, constructing the subgraph and using graph reasoning module is highly outperforming compared to Triplet Selection model.

5 Conclusion

We propose MRRG model for using external knowledge graph in reasoning over procedural text. Our model extracts a relevant subgraph for each question from the KG and uses that knowledge subgraph for answering the question. The extracted subgraph includes the reasoning path for answering the question and helps the multi-hop reasoning to predict an explainable answer. We evaluate MRRG on the WIQA benchmark and achieve SOTA performance.

References

279

281

287

290

291

292

293

296

299

305

306

307 308

309

310

311

312

313

314

315

317

318

319

321

322

323

324

325

326

327

328

- Akari Asai and Hannaneh Hajishirzi. 2020. Logicguided data augmentation and regularization for consistent question answering. In *Proceedings of the* 58th Annual Meeting of the Association for Computational Linguistics, pages 5642–5650. Association for Computational Linguistics.
- J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- Y. Freund and R. Schapire. 1995. A decision-theoretic generalization of on-line learning and an application to boosting. In *EuroCOLT*.
- Yinya Huang, Meng Fang, Xunlin Zhan, Qingxing Cao, Xiaodan Liang, and Liang Lin. 2021. Rem-net: Recursive erasure memory network for commonsense evidence refinement. In *AAAI*.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. KagNet: Knowledge-aware graph networks for commonsense reasoning. In *Proceedings* of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2829–2839, Hong Kong, China. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Shangwen Lv, Daya Guo, Jingjing Xu, Duyu Tang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, and Songlin Hu. 2020. Graph-based reasoning over heterogeneous external knowledge for commonsense question answering. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 34, pages 8449–8456.
- Aman Madaan, Dheeraj Rajagopal, Yiming Yang, Abhilasha Ravichander, Eduard Hovy, and Shrimai Prabhumoye. 2020. Eigen: Event influence generation using pre-trained language models. *arXiv preprint arXiv:2010.11764*.
- Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. 2016. A decomposable attention model for natural language inference. In *Proceedings* of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 2249–2255, Austin, Texas. Association for Computational Linguistics.
- Dheeraj Rajagopal, Niket Tandon, Peter Clark, Bhavana Dalvi, and Eduard Hovy. 2020. What-if I ask you to explain: Explaining the effects of perturbations in procedural text. In *Findings of the Association*

for Computational Linguistics: EMNLP 2020, pages 3345–3355, Online. Association for Computational Linguistics.

334

335

336

337

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

359

360

361

362

363

364

365

366

367

368

370

371

372

373

374

376

377

378

379

381

382

383

384

385

386

- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for ifthen reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *European semantic web conference*, pages 593–607. Springer.
- Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Bidirectional attention flow for machine comprehension. In *ICLR*.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31.
- Gabriel Stanovsky, Julian Michael, Luke Zettlemoyer, and Ido Dagan. 2018. Supervised open information extraction. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 885–895, New Orleans, Louisiana. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Niket Tandon, Bhavana Dalvi, Keisuke Sakaguchi, Peter Clark, and Antoine Bosselut. 2019. WIQA: A dataset for "what if..." reasoning over procedural text. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 6076– 6085, Hong Kong, China. Association for Computational Linguistics.
- Chen Zheng and Parisa Kordjamshidi. 2021. Relational gating for "what if" reasoning. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4015–4022. International Joint Conferences on Artificial Intelligence Organization. Main Track.

A Appendix

390

391

396

400

401

402

403

404

405

406

407

408

409

410

424

425

426

427

428

429

430

431

432

433

A.1 Implementation Details

We implemented our MRRG framework using Py-Torch. We use a pre-trained RoBERTa (Liu et al., 2019) to encode the contextual information in the input. The maximum number of nodes in the graph is 50. More hyper-parameter description of graph is shown in Table 3. The maximum number of words for the paragraph context is 256. For the graph construction module, we utilize a deep BiLSTM open Information Extraction model (Stanovsky et al., 2018) from AllenNLP² to extract the entities. The maximum number of hops for the graph reasoning module is 3. The learning rate is 1e - 5. The model is optimized using Adam optimizer (Kingma and Ba, 2015).

A.2 Baseline Description

EIGEN (Madaan et al., 2020) is a strong baseline that builds event influences based on the given document and leverages LMs to create the chain to predict the causal answer. However, EIGEN cannot cope with the situation in which the required knowledge is not included in the procedural text.

Logic-Guided (Asai and Hajishirzi, 2020) is a 411 baseline that combines neural networks and logic 412 rules. Specifically, the Logic-Guided model lever-413 ages LMs, symmetric logical rules, and transitive 414 logical rules to augment the training data and train 415 the model by symmetric and transitive consistency. 416 RGN (Zheng and Kordjamshidi, 2021) is the recent 417 SOTA baseline that utilizes a gating network to ef-418 fectively filter out the key entities and relationships 419 in the given document and learns the contextual 420 representations to predict the causal answer. How-421 ever, RGN cannot deal with the challenge that the 422 required knowledge is not in the document. 423

REM-Net (Huang et al., 2021) proposes a recursive erasure memory network to find out the causal evidence. Specifically, REM-Net refines the evidence by a recursive memory mechanism and then uses a generative model to predict the causal answer. REM-Net is the only work that uses commonsense for WIQA. However, this work has no reasoning process over the extracted external triplet and only uses implicate multi-head operator to encode the triplet and predict the answer.



Figure 4: The architecture of training the subgraph construction module.

A.3 Bi-DAF Implement Detail

We use the Bi-DAF (Seo et al., 2017) style contextual interaction module to feed E_q and E_c to Context-to-Question Attention $E_{C \to q}$ and Questionto-Context Attention $E_{q\to c}$ to obtain the contextual interaction between question and context. Here we introduce how to obtain $E_{q\to c}$ and $E_{C\to q}$.

Context-to-Question Attention $E_{\mathcal{C}\to q}$ aims to capture the information on which question tokens are semantically relevant to each document token. We first denote question representation to E_q and document representations to $E_{\mathcal{C}}$. The $E_{\mathcal{C}\to q}$ is computed as follows:

$$E_{\mathcal{C} \to q} = softmax(sim(E_q^T, E_{\mathcal{C}}))E_q,$$

where T is the transpose operation.

Question-to-Context Attention $E_{q\to C}$ aims to capture the information on which document tokens are semantically relevant to each question token. The process of computing the $E_{q\to C}$ is similar to $E_{C\to q}$.

445

434

435

436

437

438

439

440

²https://demo.allennlp.org/ open-information-extraction.