A Graph Fusion Approach for Cross-Lingual Machine Reading Comprehension

Anonymous ACL submission

Abstract

Although great progress has been made for Machine Reading Comprehension (MRC) in English, scaling out to a large number of languages remains a huge challenge due to the lack of large amounts of annotated training data in non-English languages. To address this challenge, some recent efforts of cross-lingual MRC employ machine translation to transfer knowledge from English to other languages, through either explicit alignment or implicit attention. For effective knowledge transition, it is beneficial to leverage both semantic and syntactic information. However, the existing 013 methods fail to explicitly incorporate syntax information in model learning. Consequently, the models are not robust to errors in alignment and noises in attention. In this work, we 017 propose a novel approach, named GraFusion-MRC, which jointly models the cross-lingual alignment information and the mono-lingual syntax information using a graph. We develop a series of algorithms including graph construction, learning, and pre-training. The experiments on two benchmark datasets for cross-lingual MRC show that our approach outperforms all strong baselines, which verifies the effectiveness of syntax information for 027 cross-lingual MRC. The code will be made open-sourced on Github.

1 Introduction

041

Machine Reading Comprehension (MRC) (Campos et al., 2016; Rajpurkar et al., 2016; Joshi et al., 2017; Rajpurkar et al., 2018), which aims to improve the ability of machines to read and understand human texts, is a challenging task in Natural Language Understanding (NLU) (Fader et al., 2014; Rajpurkar et al., 2016; Lewis et al., 2020a; Shou et al., 2020; Wang et al., 2020). Various large-scale human-annotated corpora, such as SQuAD (Rajpurkar et al., 2016), have greatly advanced the progress in the MRC task (Seo et al., 2017; Wang et al., 2017; Hu et al., 2017; Yu et al., 2018; Devlin et al., 2019). However, those large-scale human-annotated datasets are mostly in resource-rich languages, such as English. For most languages in the world, there is, however, scarce annotated data for MRC, which limits the corresponding MRC performance.

043

044

045

046

047

049

051

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

081

To tackle the challenge of data scarcity in lowresource languages, recent attempts in cross-lingual NLU adopt machine translation to transfer the knowledge learned from the high quality annotated data in resource-rich languages (i.e., the source languages) to low-resource languages (i.e., the target languages) (Upadhyay et al., 2018; Schuster et al., 2019). For example, several methods (Zhu et al., 2019; Hu et al., 2020; Liang et al., 2020) translate training data in English to target languages, and use the translated data to train the cross-lingual MRC models. Some other methods (Cui et al., 2019; Fang et al., 2020) translate test cases in a target language to English, and use the representation of the translated cases in English to enhance the representations of the original test cases.

For effective knowledge transfer across languages, both semantic and syntactic information is highly valuable and thus should be well represented. However, all previous translation-based approaches carry over knowledge across languages only through unstructured texts, where semantic and syntactic information is implicitly represented and complicatedly entangled. To represent the correlation among words in different languages, previous works either build translation alignments or learn attention matrices. However, it is very challenging to learn the connection between words across languages solely relying on texts only. Admitted by previous studies, misalignments often happen and badly hurt model performance (Xu et al., 2020; Li et al., 2020; Pei et al., 2020). Moreover, deep learning models may pay attention to less relevant words in long text (Zhang et al., 2020).

Can we use syntax information explicitly to en-



Figure 1: A passage (P), a question (Q), and an answer (A) in English with their translations in German. The words in red are the correct answers, and the links in gray and blue represent the semantic alignments and syntax information between words, respectively. The red dashed link indicates a misalignment of answers.

hance knowledge transfer across languages and im-084 prove cross-lingual MRC? In this paper, we tackle this challenge. Figure 1 shows a motivating exam-086 ple. Suppose the source training example is in English: the question is "Where are egg tubes found inside of an insect?", and the answer "ovaries" is in the sentence "The ovaries are made up of a number of egg tubes ..." After the English example is translated into German, the corresponding answer "Eierstöcke" in "Die Eierstöcke bestehen aus einer Anzahl von Eierröhrchen" is not correctly identified due to misalignment by an off-the-shelf alignment tool GIZA++ (Och and Ney, 2003). Checking many cases manually, we find misalignments 097 commonly happen in complex sentence structures (e.g., involving passive voice where word orders are different from usual) and usages of rare words (e.g., "ovaries" and "Eierstöcke" belong to the do-101 main of biology). In such cases, syntax informa-102 tion can help the model to figure out the correct 103 alignment. In the example in Figure 1, although "ovaries" and "Eierstöcke" are not correctly aligned, 105 their parents "made/bestehen", and siblings "num-106 ber/Anzahl" are correctly aligned. Therefore, if we can leverage the syntax structure to propagate the 108 alignment information, we can learn better repre-109 sentation for the target language. 110

Carrying the above insights, in this paper, we 111 jointly model the cross-lingual alignment informa-112 tion and the mono-lingual syntax information us-113 ing a graph. We make the following contributions. 114 First, we propose using syntax information to en-115 hance knowledge transfer across languages. Sec-116 ond, we develop a novel graph fusion approach to 117 model the syntax structure as well as the alignment 118 across the source and target inputs. We design 119 a series of algorithms including graph construc-120 tion, learning, and pre-training. Last, we evaluate 121

our approach on two public cross-lingual MRC benchmarks. The experimental results show that our model effectively transfers knowledge from source language to target language through attention guided by syntax information, and hence outperforms all the strong baselines. The results verify the effectiveness of syntax information for crosslingual MRC. The code will be made open-sourced on Github. 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

149

150

151

152

153

154

155

156

157

158

2 Related Work

Given a question and a passage, the MRC task (Campos et al., 2016; Rajpurkar et al., 2016; Shen et al., 2017; Joshi et al., 2017; Shen et al., 2017; Rajpurkar et al., 2018) builds a model to find the span of the correct answer for the question from the given passage. Thanks to the emergence of large-scale pre-trained language models, such as BERT (Devlin et al., 2019), and the availability of the large-scale human-annotated corpora in English, such as SQuAD (Rajpurkar et al., 2016), great progress has been made in English MRC task (Seo et al., 2017; Shen et al., 2017; Yu et al., 2018). Interested readers can consult some recent surveys (Zhu et al., 2021; Zeng et al., 2020).

Among the various MRC approaches, Zhang et al. (2020) argue that the syntax information can prevent a model from attending to some dispensable words. The authors incorporate the syntactic dependencies between words into a pre-trained model and show significant gains in the English MRC task. This paper also considers syntax information in the MRC task but in a totally different problem setting. Zhang et al. (2020) assume a single-language setting and use the syntax information to explore the correlation between questions and passages, while we target at the cross-lingual setting where the syntax information is used to

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

210

211

212

213

guide the correlation between the inputs in sourceand target languages.

161

162

163

164

167

168

170

171

172

173

174

175

176

177

178

179

180

181

182

183

186

187

190

193

194

195

196

197

198

199

201

202

206

207

Limited by the availability of large-scale annotated data, for most languages in the world, the MRC task relies on cross-lingual MRC models, which transfer knowledge from a resource-rich language to some low-resource languages. As a baseline, some multi-lingual pre-trained models, such as mBERT (Devlin et al., 2019), XLM (Lample and Conneau, 2019), and XLM-R (Conneau et al., 2020), are fine-tuned by training data in English and then directly applied to other languages. Several previous studies (Upadhyay et al., 2018; Schuster et al., 2019; Li et al., 2021) show that the baseline usually works well for classification tasks, such as intent detection. However, for sequence labeling tasks, such as slot tagging, the results in target languages are much poorer than that in English.

Since MRC is also a sequence labeling problem, i.e., labeling the answer span in a passage, the second approach to cross-lingual MRC employs machine translators to translate training data from English to some target languages, and then adds the translated training sets into the fine-tuning stage (Cui et al., 2019; Zhu et al., 2019; Hu et al., 2020; Liang et al., 2020; Yuan et al., 2020; Liu et al., 2020). Although this approach improves the MRC results substantially, one weakness remained is the alignment quality between the training example in English and the translated example. Previous studies (Xu et al., 2020; Li et al., 2020; Pei et al., 2020) indicate that misalignments often happen and can badly degrade the model performance.

The third approach to cross-lingual MRC translates test cases in a target language into English, and combines the representation of the original test cases and the representation of the translated case to English through the attention mechanism. For example, Cui et al. (2019) apply two BERT models (Devlin et al., 2019) to encode the inputs in a target language and English, respectively, and then learn an attention matrix between the representations of the two inputs. Fang et al. (2020) concatenate the target example with the translated example to English and feed it through a stack of transformers (Vaswani et al., 2017) to jointly learn the cross-lingual attention. This approach avoids hard alignments across the inputs in different languages, and shows promising results. However, all the previous methods learn the attention matrix only from plain text without considering internal

syntax structures. To the best of our knowledge, we are the first to leverage both the semantic and syntactic information to guide the attention between the inputs in different languages.

Another approach to cross-lingual MRC is to synthesize training data. Instead of translating training data in English to other languages, this approach starts from passages in the target languages, and then generates questions and answers (Puri et al., 2020; Riabi et al., 2020; Shakeri et al., 2021). Both translation and generation can be considered as data augmentation methods, and are orthogonal to each other.

3 Methodology

In this section, we first describe the overall network architecture of our proposed *GraFusionMRC* approach in Section 3.1. We then elaborate the details of the major components of our approach, including the construction of the Syntax-Enhanced and Alignment-Aware Graph (SA-Graph) in Section 3.2 and the learning algorithm of the graph in Section 3.3. Limited by space, the pre-training methods for the SA-Graph with different masking strategies and the implementation details of our model are described in Appendix B and C, respectively.

3.1 Overview of Network Architecture

Figure 2 shows the overview of our approach. The backbone is a stack of bidirectional Transformers (Vaswani et al., 2017) with N + 2 layers. The first layer encodes the inputs, and the last layer learns the final representations for decoding.

Our major technical contribution is in the middle N layers, where a graph neural network is constructed and trained to model both syntax and alignment information. Such information jointly contributes to the knowledge transfer across languages and results in better representation for the target language through enhanced attention matrices.

Given an instance S in the source language, we first apply a machine translator to translate it to an instance T in the target language (or given an instance T in the target language, we can translate it to S in English). We then unify the length of Sand T to be l by padding or truncating operations, and input $S \in \mathbb{R}^{l \times d}$ and $T \in \mathbb{R}^{l \times d}$ in parallel to our model, where d is the dimensionality of the token embedding vectors. The input is then



Figure 2: The overview of our model *GraFusionMRC*, where the red and green nodes represent the words in the source and target language, respectively.

encoded by a Transformer as follows:

260

261

262

263

264

265

269

273

274

275

276

277

278

279

281

282

$$\boldsymbol{A}_{0}^{s} = Transformer(\boldsymbol{S}), \qquad (1)$$

$$\mathbf{A}_0^t = Transformer(\mathbf{T}). \tag{2}$$

We then take the concatenation of A_0^s and A_0^t , and apply *N Syntax-enhanced and Alignment-aware Fusion Transformer* layers (or SA-Transformer for short) to produce the representation by

$$[\boldsymbol{A}_{n}^{s}; \boldsymbol{A}_{n}^{t}] = \textit{Transformer}_{sa}([\boldsymbol{A}_{n-1}^{s}; \boldsymbol{A}_{n-1}^{t}]), \ (3)$$

where the subscripts $n \in [1, N]$ indicate that the variables are at the *n*-th SA-Transformer layer, and $[\mathbf{A}_n^s; \mathbf{A}_n^t] \in \mathbb{R}^{2l \times d}$ is the concatenation of the representations of the parallel sentences in the source and the target languages.

Each SA-Transformer layer applies a multi-head self-attention operation (Vaswani et al., 2017) followed by a feed-forward layer. Specifically, the multi-head self-attention operation first obtains a triplet consisting of the query Q_i , the key K_i and the value $V_i \in \mathbb{R}^{2l \times d_h}$ for each *head_i* by applying linear transformations W_i^q , W_i^k , and $W_i^v \in \mathbb{R}^{d \times d_h}$ on the input matrix $[A_{n-1}^s; A_{n-1}^t]$, respectively, where d_h is the dimensionality of each head, and matrices W_i^q , W_i^k , and W_i^v are parameters to be learned. Then, each *head_i* conducts the following attention operation:

$$head_i = softmax(\frac{\boldsymbol{Q}_i \boldsymbol{K}_i^{\top}}{\sqrt{d_k}} + \boldsymbol{G})\boldsymbol{V}_i, \quad (4)$$

where *i* denotes the *i*-th head of the multi-head operation, and $G \in \mathbb{R}^{2l \times 2l}$ is the attention matrix to be described in Section 3.2. After obtaining the representation $[\mathbf{A}_n^s; \mathbf{A}_n^t]$ via (3), we further add another Transformer layer to separately project \mathbf{A}_n^s and \mathbf{A}_n^t back to the individual language spaces and obtain \mathbf{A}^s and $\mathbf{A}^t \in \mathbb{R}^{l \times d}$, respectively, since our final goal is to predict the labels in the individual languages.

290

291

294

297

298

299

300

302

303

304

306

308

309

310

311

312

313

314

We then use A^s and A^t to predict the answer span in the source and target languages, respectively. Let us take A^t as an example to elaborate. Following Liu et al. (2020), we feed A^t to two separate linear layers, each followed by a softmax operation to produce the final span prediction p_{str}^t and $p_{end}^t \in \mathbb{R}^l$, i.e., the predictions of the start and the end positions, respectively. For example, p_{str}^t is calculated by $p_{str}^t = softmax(A^t \cdot u_{str} + b_{str})$, where $u_{str} \in \mathbb{R}^d$ and $b_{str} \in \mathbb{R}^l$ are two trainable parameters. We then calculate the standard cross entropy loss for the predicted start and end positions in the target language by

$$\mathcal{L}^{t} = -\frac{1}{\|\mathcal{D}\|} \sum_{i=1}^{\|\mathcal{D}\|} (\boldsymbol{y}_{str,i}^{t} \cdot \log(\boldsymbol{p}_{str,i}^{t}) + \boldsymbol{y}_{end,i}^{t} \cdot \log(\boldsymbol{p}_{end,i}^{t})),$$
(5)

where $\|\mathcal{D}\|$ is the total number of training examples, $\boldsymbol{y}_{str,i}^t$ and $\boldsymbol{y}_{end,i}^t \in \mathbb{R}^l$ are the ground-truth labels for the start and end positions of the *i*-th training example.

3.2 Syntax-Enhanced and Alignment-Aware Graph (SA-Graph)

To incorporate the syntax and alignment informa-
tion into Transformer, we learn an attention matrix315G, where an element $G_{i,j}$ in G is the attention317score indicating the attention that word i pays to318the word j. To learn the matrix G, we first con-
struct the syntax-enhanced and alignment-aware320

graph (or SA-Graph for short), where each node corresponds to a word, and the edges represent the syntax and alignment information. Given a pair of parallel sentences as input, we build a graph to represent the relations among the words in the sentences. Each word in the parallel sentences corresponds to a node in the graph, and the edges between the nodes are based on the relations between the words. As introduced in Section 1, we consider two types of relations of words, cross-lingual word alignment and mono-lingual syntactic dependency. We build edges for those two relations.

321

326

327

332

333

334

339

341

343

344

345

347

351

361

362

364

In machine translation, the corresponding words in source and target languages can be aligned with each other. Taking German sentence "*Wir sollten die Umwelt schützen*" and its parallel sentence "*We should protect the environment*" in English as an example, we can apply some off-the-shelf alignment tools, such as GIZA++¹ (Och and Ney, 2003), to compute the word alignment. The aligned words often share similar semantic meaning, for example, "*Wir*" and "*We*", "*sollten*" and "*should*", "*die*" and "*the*", "*Umwelt*" and "*environment*", as well as "*schützen*" and "*protect*". We then add *wordalignment edges* between the nodes corresponding to those words.

In addition to the edges between words across languages, we also consider the syntactic structures of sentences and build edges between words within the same language. Specifically, we first split a given passage into sentence-level and then apply the Stanza toolkit² (Qi et al., 2020) to extract the dependency between words for each sentence. Two words are connected by a *word-dependency edge* if there exists a dependency between them. We also add a special word-dependency edge between the same words in a passage.

Based on the graph, the representation f_i of a word *i* is derived by

$$\boldsymbol{f}_{i,n} = \mathcal{F}(\boldsymbol{h}_{i,n}, \mathcal{N}(i)), \tag{6}$$

where $h_{i,n}$ is the representation of word *i* from the *n*-th layer, $\mathcal{N}(i)$ denotes the neighbors of word *i* in the SA-Graph, and $\mathcal{F}(\cdot, \cdot)$ is the aggregation function of word *i* and its neighbors that will be described in Equation (8), Section 3.3. Once $f_{i,n}$ is computed, the attention matrix G is obtained by

$$\boldsymbol{G}_{i,j}^{n} = (\boldsymbol{W}_{att}^{n} \cdot \boldsymbol{f}_{i,n} + \boldsymbol{b}_{att}^{n}) \cdot (\boldsymbol{W}_{att}^{n} \cdot \boldsymbol{f}_{j,n} + \boldsymbol{b}_{att}^{n}), \quad (7)$$

where $W_{att}^n \in \mathbb{R}^{d \times d}$ and $b_{att}^n \in \mathbb{R}^d$ are trainable parameters. For convenient representation, we use f_i instead of $f_{i,n}$ in the following. Next, we present the learning process of the representation f_i .

3.3 Graph Learning

After we construct the SA-Graph, we perform a learning algorithm over the graph. For each node *i*, we want to learn a better representation $f_i = \mathcal{F}(h_i, \mathcal{N}(i))$ than its original representation h_i by aggregating the information from its neighbors $\mathcal{N}(i)$. As described in Section 3.2, there are two types of edges in the graph. Correspondingly, the node representation f_i consists of two parts:

$$\boldsymbol{f}_i = \frac{1}{2} (\boldsymbol{f}_i^a + \boldsymbol{f}_i^d), \qquad (8)$$

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

385

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

where f_i^a is the representation of word *i* aggregated from the alignment information, i.e., $f_i^a = \mathcal{F}_a(h_i, \mathcal{N}_a(i))$, where $\mathcal{N}_a(i)$ is the set of neighbors of word *i* that are connected by word-alignment edges. Similarly, f_i^d aggregates the dependency information, i.e., $f_i^d = \mathcal{F}_d(h_i, \mathcal{N}_d(i))$, where $\mathcal{N}_d(i)$ is the set of dependency neighbors. In Equation (8), in addition to the average function, other combination operators, such as weighted sum or maxpooling, may also be considered. Here we choose the simple but effective average method based on our empirical study. Experiments with other combination operators are presented and discussed in Appendix D.1.

To learn aggregation by alignment $\mathcal{F}_a(\cdot, \cdot)$, for a word *i*, the representation f_i^a aggregates the information from its neighbors $\mathcal{N}_a(i)$ connected by the word-alignment edges. As indicated in the previous studies (Xu et al., 2020; Li et al., 2020; Pei et al., 2020), word alignment is a challenging task and misalignments may exist in results produced by existing methods. To mitigate the alignment errors, we develop a gate mechanism to guard against irrelevant alignment.

$$g_{i} = \sigma(V_{1} \cdot h_{i} + W_{1} \cdot h_{j}),$$

$$f_{i}^{a} = (1 - g_{i}) \odot (V_{2} \cdot h_{i}) + g_{i} \odot (W_{2} \cdot \bar{h}_{j}),$$
(9)

where $\bar{h}_j = \operatorname{avg}\{h_j | h_j \in \mathcal{N}_a(i)\}\$ is the average of the representations of the nodes in the neighbor set $\mathcal{N}_a(i)$, σ is the sigmoid function, \odot denotes element-wise multiplication, $g_i \in \mathbb{R}^d$ serves as the role of gating, and matrices V_1, W_1, V_2 and $W_2 \in \mathbb{R}^{d \times d}$ are model parameters. g_i is the gate to control whether the aligned information should

¹https://github.com/moses-smt/giza-pp

²https://github.com/stanfordnlp/stanza

contribute to the representation of word *i*. If the nodes connected by the alignment edge bear very different semantic meanings, the weights in the gate are close to zero, which switch off the information flow.

In addition to cross-lingual word alignment information, the mono-lingual syntax information discloses the inherent dependency among words and thus also benefits the representation of words. The representation f_i^d aggregates the syntax information for node *i* using a graph attention network (Velickovic et al., 2018) as follows.

$$\boldsymbol{f}_{i}^{d} = \sigma(\sum (\alpha_{iu} \boldsymbol{W}_{3} \boldsymbol{h}_{u}, \forall u \in \mathcal{N}_{d}(i))), \quad (10)$$

where $W_3 \in \mathbb{R}^{d \times d}$ is a model parameter, σ is the sigmoid function, and $\alpha_{iu} \in \mathbb{R}$ is the attention coefficient that indicates the importance of word *i* to its neighbor *u*, calculated as follows:

$$\boldsymbol{\alpha}_{iu} = \frac{\exp(LR(\boldsymbol{W}_{4}[\boldsymbol{h}_{i};\boldsymbol{h}_{u}]))}{\sum_{k\in\mathcal{N}_{d}(i)}\exp(LR(\boldsymbol{W}_{4}[\boldsymbol{h}_{i};\boldsymbol{h}_{k}]))}, \quad (11)$$

where LR is the Leaky ReLU activate function, and $W_4 \in \mathbb{R}^{2d}$ is a model parameter.

To enhance the representation power of the SA-Graph, we use translated parallel data to pre-train the graph. The basic idea is to randomly mask some nodes in the graph and use the representations of its semantic and syntactic neighbors to recover it. Limited by space, the details are left in Appendix B.

Experiments

We evaluate the proposed *GraFusionMRC* approach on two benchmark datasets. In this section, we first describe the experiment setup. We then report and analyze the experimental results. We also illustrate how SA-Graph affects the attention weights through a case study. The further analysis of our model is presented in Appendix D.

4.1 Experimental Setup

4.1.1 Datasets and Evaluation Metrics

MLQA (Lewis et al., 2020b) and TyDiQA-GoldP dataset (Clark et al., 2020) are two recent public benchmark datasets for cross-lingual machine reading comprehension. The details of these two datasets are given in the Appendix A.

Although MLQA and TydiQA provide sufficient test data, their training data is quite limited. Following Fang et al. (2020), we use SQuAD v1.1 (Rajpurkar et al., 2016) English training data as additional data during the fine-tuning stage of our model. Moreover, the English training data in SQuAD v1.1 is further translated into the target languages in the MLQA and TyDiQA-GoldP test data via the Google Machine Translation system³. Besides, to pre-train our SA-Graph model, we further collect additional parallel sentences following Huang et al. (2019). There are one million pairs of parallel sentences in English and each target language.

We adopt the standard evaluation metrics from the SQuAD dataset (Rajpurkar et al., 2016), including F1 and Exact Match (EM) scores. The F1 score is used to measure the overlap of tokens between the predicted and ground-truth answer spans, while the EM score only counts the cases where the predicted answer spans exactly match the ground-truth answer spans. We run the official evaluation script provided by Lewis et al. (2020b) and Clark et al. (2020) for the MLQA and TyDiQA-GoldP datasets, respectively, to report the results.

4.1.2 Baselines

We compare *GraFusionMRC* with the following two groups of approaches. The first group consists of models fine-tuned using training data in English only. The second group is the methods that employ parallel translation data.

Fine-tuning with English training data only In this group of baselines, we pick the existing cross-lingual models, including **mBERT** (Devlin et al., 2019), **XLM** (Lample and Conneau, 2019), and **XLM-R** (Conneau et al., 2020). We also include the encoder **MMTE** (Siddhant et al., 2020), which is a large-scale cross-lingual neural machine translation model for up to 102 languages to and from English. These models are fine-tuned using English training data only.

Models using translation In this group, we first select XLM-R as the representative for cross-lingual models, since it performs the best among all the models in the first group in our experiments for the cross-lingual MRC task. We then fine-tune the XLM-R model with the combined translated training data of all languages jointly, which is denoted as **XLM-R (translate-train)**. We also include baselines, the **FILTER** (Fang et al., 2020), which leverages the intrinsic cross-lingual correlation between different languages, and the **XLM-ALIGN**_{base} (Chi et al., 2021), which introduces

³https://console.cloud.google.com/storage/browser/xtreme _translations

Table 1: MLQA results (F1 / EM) for each language.

Model	en	ar	de	es	hi	vi	zh	Avg.
mBERT	80.2 / 67.0	52.3 / 34.6	59.0/43.8	67.4 / 49.2	50.2/35.3	61.2 / 40.7	59.6 / 38.6	61.4 / 44.2
XLM	68.6 / 55.2	42.5 / 25.2	50.8 / 37.2	54.7 / 37.9	34.4 / 21.1	48.3 / 30.2	40.5 / 21.9	48.5 / 32.7
MMTE	78.5 / -	56.1 / -	58.4 / -	64.9 / -	46.2 / -	59.4 / -	58.3 / -	60.3 / 41.4
XLM-R _{base}	78.5 / 65.3	56.1 / 36.8	61.7 / 47.1	66.0 / 48.7	60.1 / 42.4	63.6 / 43.5	60.1 / 35.5	63.7 / 45.6
XLM-R _{base} (translate-train)	77.8 / 64.4	58.0/38.1	63.4 / 49.1	68.7 / 51.9	62.8 / 46.1	65.3 / 45.9	61.8 / 36.9	65.4 / 47.5
FILTER _{base}	77.2 / 63.9	60.2 / 41.2	66.9 / 52.7	70.5 / 53.2	64.5 / 47.2	66.8 / 47.7	63.4 / 42.1	67.1 / 49.7
XLM-ALIGN _{base}	81.5 / 68.3	60.7 / 41.2	64.5 / 49.8	70.3 / 52.2	65.2 / 47.5	69.8 / 48.9	64.4 / 40.4	68.1 / 49.8
GraFusionMRC _{base} +a	77.7/64.0	61.2 / 42.2	67.6/53.4	72.9/55.7	66.0/48.4	67.1/48.6	64.6/43.3	68.2 / 50.8
GraFusionMRC _{base} +ad	77.5 / 63.8	61.9 / 42.8	68.9 / 54.9	73.4 / 56.4	66.6 / 49.2	68.4 / 49.1	65.2 / 44.0	68.8 / 51.5
$GraFusionMRC_{base}$ +adp	77.4 / 63.8	62.8 / 43.4	69.3 / 55.3	74.0 / 56.8	67.1 / 49.5	68.8 / 49.5	65.6 / 44.3	69.3 / 51.8
XLM-R _{large} (translate-train)	83.5 / 70.6	66.6 / 47.1	70.1 / 54.9	74.1 / 56.6	70.6 / 53.1	74.0 / 52.9	62.1 / 37.0	71.6 / 53.2
FILTER _{large}	84.0 / 70.8	72.1 / 51.1	74.8 /60.0	78.1 / 60.1	76.0 / 57.6	78.1 /57.5	70.5 / 47.0	76.2 / 57.7
GraFusionMRC _{large} +a	84.2 / 71.5	73.0/52.0	75.4/60.3	78.8760.9	77.9/58.4	79.0/57.8	71.4748.6	77.1/58.5
GraFusionMRC _{large} +ad	83.9 / 71.0	73.7 / 52.5	75.9/61.2	79.6/61.2	78.6 / 58.7	79.9 / 59.8	72.4 / 48.8	77.7 / 59.0
GraFusionMRC _{large} +adp	83.5 / 70.7	74.2 / 52.7	76.2 / 61.7	80.1 / 62.0	79.2 / 59.0	80.4 / 60.1	73.0 / 49.3	78.1 / 59.4

Table 2: TyDiQA-GoldP results (F1 / EM) for each language. As the Stanza toolkit doesn't support languages *Bengali* and *Swahili*, we don't report results on these two languages in the syntactic fusion setting. Please note that we correct the text segment module of FILTER when handling ko language and bring its performance back from 33.1 to 68.9 of the F1 score.

Model	en	ar	bn	fi	id	ko	ru	SW	te	Avg.
mBERT	75.3 / 63.6	62.2 / 42.8	49.3 / 32.7	59.7 / 45.3	64.8 / 45.8	58.8 / 50.0	60.0 / 38.8	57.5 / 37.9	49.6 / 38.4	59.7 / 43.9
XLM	66.9 / 53.9	59.4 / 41.2	27.2 / 15.0	58.2/41.4	62.5 / 45.8	14.2 / 5.1	49.2 / 30.7	39.4 / 21.6	15.5 / 6.9	43.6 / 29.1
MMTE	62.9 / 49.8	63.1 / 39.2	55.8 / 41.9	53.9 / 42.1	60.9 / 47.6	49.9 / 42.6	58.9 / 37.9	63.1 / 47.2	54.2 / 45.8	58.1/43.8
XLM-R _{base}	71.9/57.1	54.3 / 32.1	57.8 / 44.6	63.9 / 50.4	68.5 / 51.3	61.2 / 38.7	60.4 / 33.8	65.2 / 55.6	64.9 / 47.6	63.1 / 45.7
XLM-R _{base} (translate-train)	71.6 / 56.4	57.8 / 34.5	60.8 / 48.6	67.1 / 53.0	71.9 / 53.7	63.3 / 40.4	62.2 / 34.7	62.8 / 53.7	67.3 / 50.9	65.0/47.3
FILTER _{base}	68.4 / 55.5	58.3 / 34.6	61.3 / 46.7	67.7 / 54.0	72.2 / 54.5	65.5/41.1	63.3 / 35.4	72.3 / 63.1	67.1 / 49.6	66.2 / 48.3
XLM-ALIGN _{base}	69.4 / 56.2	68.7 / 49.4	56.0 / 38.9	64.2 / 47.2	73.9 / 57.9	53.0 / 40.4	62.3 / 38.0	60.1 / 42.8	51.0/31.9	62.1 / 44.8
GraFusionMRC _{base} +a	71.2/55.7	60.1 / 37.1	62.4 / 48.8	69.0/54.6	73.5757.2	67.2/43.5	64.6/38.0	73.9/64.1	68.1/53.4	67.8/50.3
GraFusionMRC _{base} +ad	70.8 / 55.4	61.4 / 38.6	/	69.7 / 55.1	74.9 / 59.0	68.1 / 44.7	64.8 / 36.7	/	70.4 / 53.2	68.6/49.0
$GraFusionMRC_{base}$ +adp	70.6 / 55.1	62.6 / 39.5	/	70.4 / 55.7	75.7 / 59.3	69.0 / 46.1	65.5 / 37.2	/	71.0 / 53.6	69.3 / 49.5
XLM-R _{large} (translate-train)	75.1 / 62.0	66.9 / 39.8	63.8 / 47.5	70.1 / 52.8	77.1/61.7	67.8 / 43.4	66.5 / 41.8	65.7 / 47.8	69.6/43.4	69.2 / 48.9
FILTER _{large}	72.4 / 59.1	72.8 / 50.8	70.5 / 56.6	73.3 / 57.2	76.8 / 59.8	68.9 / 45.7	68.9 / 46.6	77.4 / 65.7	69.9 / 50.4	72.3 / 54.7
GraFusionMRC _{large} +a	74.1/61.3	73.4751.6	71.7 / 57.5	74.1/58.0	77.8/62.4	69.5/46.2	69.8 / 46.7	78.0/65.7	70.3 / 53.0	73.2/55.8
GraFusionMRClarge+ad	73.9/61.2	74.2 / 53.2	/	74.9 / 59.5	79.2 / 64.2	70.5 / 47.5	70.4 / 47.6	/	72.0 / 54.9	73.6/55.4
GraFusionMRClarge+adp	73.5 / 60.8	75.1 / 53.8	/	76.2 / 61.0	79.8 / 64.2	71.3 / 48.3	71.3 / 48.2	/	72.5 / 55.7	74.2 / 56.0

denoising word alignment pre-training task. Please note that the **XLM-ALIGN** only provides the base model.

4.2 Experimental Results

510

511

512

513

514

515

516

518

519

520

521

522

523

524

We conduct experiments with three variants of our *GraFusionMRC* approach: (1) *GraFusion-MRC*+a: only the word-alignment edges are used in graph learning and all the word-dependency edges are ignored; (2) *GraFusionMRC*+ad: both the word-alignment and word-dependency edges are included in the graph learning stage to obtain the node representation; and (3) *GraFusionMRC*+adp: the pre-training stage is added before the graph learning stage to enhance the representation power of the SA-Graph.

The results for all the methods on the MLQA and the TyDiQA-GoldP datasets are presented in Table 1 and Table 2, respectively. In the first group of baselines, the XLM- R_{base} model consistently outperforms all other baselines in most of the target languages, demonstrating itself as a strong baseline for the cross-lingual MRC task. Based on this observation, we use XLM-R as the representative for cross-lingual models, and further fine-tune this model with translated training data in target languages.

As shown in the first row ("**XLM-R** (**translatetrain**)") of the second and third group of baselines, adding translated data in target languages substantially improves the model performance, which suggests that the translated data strengthen knowledge transfer effectively. We also observe the FILTER method performs better than the strong baseline XLM-R (translate-train) on both datasets. It indicates that the attention between the source sentence



Figure 3: Visualization of the FILTER model (left) and the proposed model (right). The triple of (P,Q,A) in German is ("Jahren stieg die objektorientierte -Programmierung an ...", "Welche Art von Programmierung hat in den 1990er Jahren den Umgang mit Datenbanken verändert?", "objektorientierte"). The corresponding translation in English is ("... along with a rise in object-oriented programming ...", "In the 1990s, what type of programming changed the handling of databases?", "object-oriented"). In order to distinguish the two languages more clearly, the German sentence is presented in green color and the corresponding English sentence in orange color. We also highlight the correct answer spans in both German and English.

and the translated target sentence leads to better representation of words, and further contributes to the cross-lingual MRC task.

All three variants of our *GraFusionMRC* approach outperform the XLM-R and FILTER models on both datasets. In particular, the *GraFusionMRC*+adp method achieves an average improvement of 2 points over the FILTER model in both MLQA and TyDiQA-GoldP. The major difference of *GraFusionMRC* from FILTER is that we enhance the learning of the attention matrix between the inputs in the source and target languages through explicit syntax and alignment information. Moreover, our gate mechanism and graph attention network increase the model robustness against the errors in alignment and syntactic parsing.

When we compare the three variants of the *Gra-FusionMRC* approach, the general trend is that using both the syntax edges and alignment edges is better than using alignment edges alone. This justifies the effectiveness of injecting syntax information into representation learning. Moreover, the pre-training using large-scale parallel data also boosts the model performance with a clear gain.

An interesting observation is that our base $GraFusionMRC_{base}$ model even outperforms the large XLM-R_{large} model in the TyDiQA-GoldP dataset among languages *fi*, *ko*, *sw*, and *te*. The use of alignment and syntactic information successfully bridges the model performance gap caused by the number of parameters, and once again confirms

the effectiveness of utilizing SA-Graph.

4.3 Visualization

To showcase the effectiveness of our SA-Graph, we compare the attention distributions from the last fusion layer of the FILTER model with that of our proposed *GraFusionMRC* in Figure 3. Please note that the original answer in German "objektorientierte" is misaligned to the word "were" in English by the GIZA++ toolkit. With the help of syntactic information, our model is able to learn a higher attention weight between "objektorientierte" and the correct parallel word "object-oriented". The visualization of the example illustrates the benefit of the SA-Graph, which improves knowledge transfer through the enhanced attention matrix.

5 Conclusion

In this paper, we develop a novel *GraFusionMRC* approach that leverages both cross-lingual alignment information and mono-lingual syntactic information for cross-lingual MRC. To the best of our knowledge, we are the first to explicitly inject both information to enhance the representation learning in the cross-lingual MRC task. We develop a systematic approach including the construction of the Syntax-Enhanced and Alignment-Aware Graph, the learning algorithms, as well as the pre-training strategies. The experimental results show that our approach outperforms all strong baselines on two public cross-lingual MRC benchmarks.

References

604

606

611

612

613

614

615

618

622

623

633

634

635

642

654

655

- Daniel Fernando Campos, T. Nguyen, M. Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, L. Deng, and Bhaskar Mitra. 2016. Ms marco: A human generated machine reading comprehension dataset. *ArXiv*, abs/1611.09268.
 - Zewen Chi, Li Dong, Bo Zheng, Shaohan Huang, Xian-Ling Mao, Heyan Huang, and Furu Wei. 2021. Improving pretrained cross-lingual language models via self-labeled word alignment. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing* (Volume 1: Long Papers), pages 3418–3430, Online. Association for Computational Linguistics.
- J. Clark, Eunsol Choi, M. Collins, Dan Garrette, T. Kwiatkowski, V. Nikolaev, and J. Palomaki. 2020. Tydiqa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8440– 8451, Online. Association for Computational Linguistics.
 - Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2019. Cross-lingual machine reading comprehension. 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 1586–1595, Hong Kong, China. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. 2014. Open question answering over curated and extracted knowledge bases. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1156–1165.
- Yuwei Fang, Shuohang Wang, Zhe Gan, S. Sun, and Jing jing Liu. 2020. Filter: An enhanced fusion method for cross-lingual language understanding. *ArXiv*, abs/2009.05166.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization. In *International Conference on Machine Learning*. 660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

703

704

706

708

709

712

- Minghao Hu, Yuxing Peng, Zhen Huang, Xipeng Qiu, Furu Wei, and Ming Zhou. 2017. Reinforced mnemonic reader for machine reading comprehension. *arXiv preprint arXiv:1705.02798*.
- H. Huang, Yaobo Liang, N. Duan, Ming Gong, Linjun Shou, D. Jiang, and M. Zhou. 2019. Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks. In *EMNLP/IJCNLP*.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Guillaume Lample and Alexis Conneau. 2019. Crosslingual language model pretraining. In *NeurIPS*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. pages 7871–7880.
- Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020b. MLQA: Evaluating cross-lingual extractive question answering. pages 7315–7330.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2021. MTOP: A comprehensive multilingual task-oriented semantic parsing benchmark.
- Xin Li, Lidong Bing, Wenxuan Zhang, Zheng Li, and Wai Lam. 2020. Unsupervised cross-lingual adaptation for sequence tagging and beyond.
- Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, et al. 2020. Xglue: A new benchmark dataset for cross-lingual pretraining, understanding and generation. *arXiv preprint arXiv:2004.01401*.
- Junhao Liu, Linjun Shou, Jian Pei, Ming Gong, Min Yang, and Daxin Jiang. 2020. Cross-lingual machine reading comprehension with language branch knowledge distillation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2710–2721, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

714

715

716

717

718

719

721

723

724

726

730

731

733

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

753 754

755

756

757

758

759

760

761

762

763

765

768

- Shichao Pei, Lu Yu, Guoxian Yu, and Xiangliang Zhang. 2020. Rea: Robust cross-lingual entity alignment between knowledge graphs. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2175– 2184.
- Raul Puri, Ryan Spring, Mohammad Shoeybi, Mostofa Patwary, and Bryan Catanzaro. 2020. Training question answering models from synthetic data. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5811–5826, Online. Association for Computational Linguistics.
 - Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for squad. *ArXiv*, abs/1806.03822.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Arij Riabi, Thomas Scialom, Rachel Keraron, Benoît Sagot, Djamé Seddah, and Jacopo Staiano. 2020. Synthetic data augmentation for zero-shot crosslingual question answering.
- Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. Cross-lingual transfer learning for multilingual task oriented dialog. 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 3795–3805, Minneapolis, Minnesota. Association for Computational Linguistics.
- Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hananneh Hajishirzi. 2017. Bidirectional attention flow for machine comprehension. In *International Conference on Learning Representations*.
- Siamak Shakeri, Noah Constant, Mihir Kale, and Linting Xue. 2021. Towards zero-shot multilingual synthetic question and answer generation for crosslingual reading comprehension. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 35–45, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Yelong Shen, Po-Sen Huang, Jianfeng Gao, and Weizhu Chen. 2017. Reasonet: Learning to stop reading in machine comprehension. In *Proceedings* of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 1047–1055. 770

774

776

779

780

781

782

783

784

785

786

787

788

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

- Linjun Shou, Shining Bo, Feixiang Cheng, Ming Gong, Jian Pei, and Daxin Jiang. 2020. Mining implicit relevance feedback from user behavior for web question answering. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2931–2941.
- Aditya Siddhant, Melvin Johnson, Henry Tsai, N. Arivazhagan, Jason Riesa, Ankur Bapna, Orhan Firat, and Karthik Raman. 2020. Evaluating the cross-lingual effectiveness of massively multilingual neural machine translation. *ArXiv*, abs/1909.00437.
- S. Upadhyay, M. Faruqui, G. Tur, H. Dilek, and L. Heck. 2018. (almost) zero-shot cross-lingual spoken language understanding. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 6034–6038.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, L. Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.
- Petar Velickovic, Guillem Cucurull, A. Casanova, A. Romero, P. Liò, and Yoshua Bengio. 2018. Graph attention networks. *ArXiv*, abs/1710.10903.
- Qifan Wang, Li Yang, Bhargav Kanagal, Sumit Sanghai, D Sivakumar, Bin Shu, Zac Yu, and Jon Elsas. 2020. Learning to extract attribute value from product via question answering: A multi-task approach. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 47–55.
- Zhiguo Wang, Haitao Mi, Wael Hamza, and Radu Florian. 2017. Multi-perspective context matching for machine comprehension. *arXiv preprint arXiv:1612.04211*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Weijia Xu, Batool Haider, and Saab Mansour. 2020. End-to-end slot alignment and recognition for crosslingual NLU. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5052–5063, Online. Association for Computational Linguistics.
- Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V Le. 2018. Qanet: Combining local convolution with global self-attention for reading comprehension. In

826

tions.

tics.

International Conference on Learning Representa-

Fei Yuan, Linjun Shou, Xuanyu Bai, Ming Gong, Yaobo Liang, Nan Duan, Yan Fu, and Daxin Jiang. 2020. Enhancing answer boundary detection for multilingual machine reading comprehension. In

Proceedings of the 58th Annual Meeting of the As-

sociation for Computational Linguistics, pages 925-

934, Online. Association for Computational Linguis-

Changchang Zeng, Shaobo Li, Qin Li, Jie Hu, and Jian-

Zhuosheng Zhang, Yu-Wei Wu, Junru Zhou, Sufeng

Fengbin Zhu, Wengiang Lei, C. Wang, J. Zheng, Sou-

Junnan Zhu, Qian Wang, Yining Wang, Yu Zhou, Jiajun Zhang, Shaonan Wang, and Chengqing Zong.

2019. NCLS: Neural cross-lingual summarization.

In Proceedings of the 2019 Conference on Empirical

Methods in Natural Language Processing and the

9th International Joint Conference on Natural Lan-

guage Processing (EMNLP-IJCNLP), pages 3054-

3064, Hong Kong, China. Association for Computa-

We evaluate our model on two public cross-

lingual machine reading comprehension datasets: MLQA (Lewis et al., 2020b) and TyDiQA-GoldP

• MLQA (Lewis et al., 2020b), is a cross-

lingual machine reading comprehension

benchmark that covers 7 languages, includ-

ing English, Arabic, German, Spanish, Hindi,

Vietnamese and Simplified Chinese. The num-

ber of question-answering instances in the test

set for those languages is 11590, 5335, 4517,

5254, 4918, 5495, and 5137, respectively.

More Details for the Datasets

janya Poria, and T. Chua. 2021. Retrieving and read-

ing: A comprehensive survey on open-domain ques-

net: Syntax-guided machine reading comprehension.

Duan, Hai Zhao, and Rui Wang. 2020.

tion answering. ArXiv, abs/2101.00774.

datasets. Applied Sciences, 10(21).

ArXiv, abs/1908.05147.

tional Linguistics.

dataset (Clark et al., 2020).

Α

jun Hu. 2020. A survey on machine reading compre-

hension-tasks, evaluation metrics and benchmark

- 836 837 838 839
- 842

- 850 853
- 854
- 858

• TyDiQA-GoldP (Clark et al., 2020), is an-870 other cross-lingual machine reading compre-871 hension benchmark covering 9 typologically diverse languages, including English, Arabic, Bengali, Finnish, Indonesian, Korean, Russian, Swahili, and Telugu. The number of 875 question-answering instances in the develop-876 ment set for those languages is 440, 921, 113, 782, 565, 276, 812, 499, and 669, respectively.



Figure 4: Illustration of the graph masking strategies, where dotted node stands for the masked node, and nodes connected by dotted lines are used to predict the masked nodes. (left) Mono-lingual Masking Strategy. (right) Cross-Lingual Masking Strategy.

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

Pre-training SA-Graph B

Sg-

To enhance the representation power of the SA-Graph, we use translated parallel data to pre-train the graph. To be more specific, given a source sentence S and the translated sentence T, we construct the SA-graph in the same way as the fine-tuning stage. Therefore, we keep consistent between the pre-training stage and the fine-tuning stage. The difference is that in the pre-training stage, we randomly mask some nodes in the graph and use their neighbors to recover them, which is shown in Figure 4. To be more specific, for each masked node *i*, we aggregate the representations of its neighbors. The aggregated representation is then fed into a linear classifier, which outputs the probabilities over the whole vocabulary. Cross entropy is used to compute the recovery loss as follows:

$$\mathcal{L}_{SA}(i) = -\log P(i|\mathcal{N}(i)). \tag{12}$$

Mono-lingual Masking Given a source sentence S and the translated sentence T, the first masking strategy constrains all the masked tokens to be within only one language, i.e., either the source or the target language. For those masked tokens, since the corresponding words in the other language should not be masked according to the masking constraint, the model can learn from the alignment information to predict the masked ones. In other words, this masking strategy encourages the model to explore the semantic correlation from the alignment information. At each iteration in our implementation, we first choose a language, and then randomly mask 15% of nodes belonging to the chosen language in SA-Graph are masked at random.

Cross-Lingual Masking With the above monolingual masking strategy, the model learns to leverage the word-alignment edges to predict the

masked ones. However, this would make the model tend to ignore the word-dependency edges. To facilitate the model to leverage the syntax information, we further develop a cross-lingual masking strategy: whenever a node is masked, its aligned node must be masked together. In this way, we cut off the alignment information flow, and the model is forced to learn from word-dependency edges to recover the masked nodes.

916

917

918

919

921

922

925

927

928

930

932

933

934

935

937

938

939

941

942

947

948

951

952

954

957

958

962

963

During the pre-training stage, we adopt each masking strategy half of the time. Besides the above two masking strategies, we also employ translation language modeling (TLM) in our pre-training process, which has shown strong performance in XLM pre-trained model (Lample and Conneau, 2019). For each masked word *i*, we compute the recovery loss as follows:

$$\mathcal{L}_{TLM}(i) = -\log P(i|\boldsymbol{h}_i). \tag{13}$$

The final loss for the pre-training is the sum of the loss of translation language modeling and our graph masking tasks, i.e., $L(i) = \mathcal{L}_{SA}(i) + \mathcal{L}_{TLM}(i)$.

C Implementation Details

We implement on top of HuggingFace's Transformers (Wolf et al., 2019) and report results on two both base and large models, i.e., GraFusionMRCbase and GraFusionMRClarge. We initialize our base model by the pre-trained XLM-R base model released by HuggingFace⁴, which contains 12 layers; and use XLM-R large model for initializing our large model, which contains 24 layers. We set the number of intermediate Transformer layers, i.e., the Syntax-Enhanced and Alignment-Aware Transformer layers, to 10 in the base model and to 22 in the large model. The first bottom Transformer layer is used for encoding the raw input sentences and the top layer converts the joint representation of the sentences in the source and target languages back to individual language spaces. To make a fair comparison, we reproduce FILTER_{base} and FILTER_{large} models based on the XLM-R base and large model, respectively. During the fine-tuning stage, following Liu et al. (2020) and Conneau et al. (2020), we pair the passage and question for each language as $|\langle cls \rangle, question, \langle s \rangle, \langle s \rangle$, passage, </s>] as the input. We then concatenate the input of the target language with the translated input in the source language.





Figure 5: Further analysis on the choice of different aggregation functions using the *GraFusionMRC*+ad model.

To make a fair comparison in our experiments, all the methods use translation approaches only. As future work, we can combine the synthesized data, the translated data, as well as the original English training data to train an initial cross-lingual MRC model. Then the proposed GraFusionMRC approach can be further applied on top of this initial model. 964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

Please note that, since we employ individual Transformer layers to encode the inputs for each language at the bottom of our model, the supporting max sequence length of each language can be 512. However, we set the max sequence length to 384 for each language to balance efficiency. Then we concatenate the source and target input and feed a sequence of length 768 to the SA-Transformer to extract their correlation.

D Further Analysis

In this section, we further explore the choice of aggregation function in (8) and also report the performances of our model when pre-trained with different masking strategies. Note that due to the timeconsuming issue of training the large models under different settings, we conduct further analysis on our base models, i.e., *GraFusionMRC*_{base}.

D.1 Comparison of Different Combinations

In this subsection, we study the effect of different aggregation functions in our *GraFusionMRC*+ad model. For each node *i* in SA-Graph, to aggregate the f_i^a and f_i^d vector, we experiment with three different types of aggregation functions: *concatenation*, weighted sum, and max-pooling. Please note that, for concatenation operation, we need an extra trainable matrix $W_{cat} \in \mathbb{R}^{2d \times d}$ to project the vector back to the dimension of *d* for consistency. For better comparison, we also report the results of



Figure 6: Effectiveness of graph masking strategies to the target languages, *Arabic, German* and *Spanish*, on MLQA dataset.

the average operation, which is used in our model.

1000

1001

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1016

1017

1018

1019

1021 1022

1023

1024

1025

1026

1027

1028

1030

1031

1032

1034

1035

1036

1037

1038

As shown in Figure 5, the average function tends to show the best performance on both MLQA and TyDiQA-GoldP datasets, which indicates that the average operation is a simple but effective method to aggregate the cross-lingual and mono-lingual correlation. We can also see that the concatenation operation shows the worst. A possible explanation could be that the introduction of the matrix W_{cat} increases the computational complexity of the aggregation process, making it more difficult to find the global optimal solution. Another interesting observation is that the performance of the weighted sum operation outperforms the max-pooling operation. This may be because the weighted sum operation could capture the fine-grained correlation between vectors while the max-pooling operations only capture the significant information.

D.2 Comparison of Masking Strategies

Below, we provide a detailed analysis to better understand the effectiveness of graph masking strategies. Specifically, we further explore the impacts on *GraFusionMRC*+ad of different choices of graph masking strategies, including using either mono-lingual or cross-lingual masking strategies or employ them both. Figure 6 shows the F1 scores of target languages, *Arabic, German* and *Spanish*, on MLQA datasets with different choices of strategies and the pre-training steps.

It can be observed that, in all three languages, the cross-lingual masking strategy performs better, but converges slower than the mono-lingual masking strategy. This may be attributed to the reason that the cross-lingual masking strategy is more complex than mono-lingual, and thus needs more training steps to fit the training data, and provide stronger alignment capture capabilities. Meanwhile, combining both strategies can obtain the best performance, which indicates that those two strategies might focus on different alignment information and1039thus be complementary to each other.1040

E Fine-tuning Details

We select the hyperparameters from batch size: 1042 {16, 32, 64}, learning rate: {1e-5, 3e-5, 5e-5, 1e-1043 6, 5e-6, and warmup rate: $\{5\%, 10\%, 15\%\}$. 1044 All experiments are conducted on 4 16G NVIDIA 1045 P100 GPUs. Each experiment is repeated 10 times 1046 and the average results are reported. The number of parameters of FILTER, GraFusionMRC+a 1048 and GraFusionMRC+ad model are 1.02, 1.07 and 1049 1.15 times that of XLM-R. For MLQA dataset, the 1050 best performance of GraFusionMRC+adp model is achieved at batch size=32, learning rate=5e-5, 1052 warmup rate=10%. It takes about 8 hours to get 1053 the best result running on 4 16G P100. As for 1054 the TyDiQA-GoldP dataset, we achieve the best re-1055 sults at batch size=64, learning rate=3e-5, warmup rate=5%. It takes about 10 hours to get the best 1057 result running on 4 16G P100.