Transfer Learning from Semantic Role Labeling to Event Argument Extraction with Template-based Slot Querying

Anonymous ACL submission

Abstract

In this work, we investigate transfer learning from semantic role labeling (SRL) to event argument extraction (EAE), considering their 004 similar argument structures. We view the extraction task as a role querying problem, unifying various methods into a single framework. There are key discrepancies on role labels and 800 distance arguments between semantic role and event argument annotations. To mitigate these discrepancies, we specify natural language-like queries to tackle the label mismatch problem and devise argument augmentation to recover distant arguments. We show that SRL anno-013 tations can serve as a valuable resource for EAE, and a template-based slot querying strategy is especially effective for facilitating trans-017 fer. In extensive evaluations on two English EAE benchmarks, our proposed model obtains impressive zero-shot results by leveraging SRL annotations, reaching nearly 80% of the fullysupervised scores. It could further provides benefits in low-resource cases, where few EAE 023 annotations are available. Moreover, we show that our approach generalizes to cross-domain 024 and multilingual scenarios.

1 Introduction

027

034

040

Event argument extraction (EAE) is a key component in the task of event extraction (Ahn, 2006) that aims to identify the arguments that serve as roles for event frames. While recent developments in neural network models have enabled impressive improvements on this task in the fully-supervised setting (Wang et al., 2019b; Pouran Ben Veyseh et al., 2020; Ma et al., 2020; Li et al., 2021b), EAE still remains challenging when abundant annotations are not available. In particular, event schemes are usually *specific* to the target scenarios. For example, events in biomedical domains, like GENE-EXPRESSION in GENIA (Kim et al., 2008), can be quite different than the ones in ACE (LDC, 2005), such as ATTACK and CONTACT. It is costly and



Figure 1: Example annotations with ACE events (above) and PropBank semantic frames (below). Brown and red rectangles indicate predicate and argument words, respectively. Green lines denote argument links.

inefficient to annotate large amounts of data for every new application.

Compared with the *specific* event schemes, semantic role labeling (SRL) is a *general* linguistic task with the goal of extracting predicate-argument structures from text (Gildea and Jurafsky, 2002; Palmer et al., 2010). There are rich and carefully annotated SRL resources, such as PropBank (Palmer et al., 2005) and FrameNet (Baker et al., 1998), covering a wide range of semantic frame types. As shown in the example in Figure 1, the SRL task resembles EAE much: they both specify semantic frames triggered by predicate words and aim at finding arguments for participating roles. Therefore, it is natural to consider applying transfer learning (Pan and Yang, 2009; Ruder et al., 2019) to enhance EAE with general SRL resources.

Notwithstanding the similarities, there are two main discrepancies between SRL and EAE structures that should be managed in order to facilitate transfer between the tasks. The first is *label mismatch*. For example, ACE adopts role names with natural language words, such as BUYER and PLACE, whereas PropBank utilizes generalized labels like ARG0 and ARGM-LOC. Although FrameNet also adopts natural language role names, it is challenging to find clear direct mappings to the target event frames. Moreover, SRL resources do not typically

089

094

098

102

103

105

106

108

109

110

111

112

113

114

115

annotate *distant arguments*, where there are no explicit syntactic encodings expressing the argument relation.¹ For example, in the sentence depicted in Figure 1, though it can be understood that the "store" is very likely to be the place where the "buying" happens, SRL annotations do not include this semantically inferred link, whereas it is considered an argument in event annotations.

In this work, we provide a comprehensive investigation of transfer from SRL to EAE. We view the tasks as a role querying problem within a *unified* framework, which covers many argument extraction methods, including classification-based methods (Ouchi et al., 2018; Ebner et al., 2020), machine reading comprehension (MRC)-based methods (Liu et al., 2020; Du and Cardie, 2020; Li et al., 2020; Feng et al., 2020; Lyu et al., 2021; Liu et al., 2021a) as well as sequence-to-sequence generation based ones (Li et al., 2021b; Hsu et al., 2021; Lu et al., 2021). We further explore a templatebased slot querying strategy, by querying argument roles using contextualized representations of the corresponding role slots in the frame template. We tackle the label-mismatch problem by forming the queries in templated natural language, which allows for the same query representation to be shared across varied schemes. To mitigate the lack of distant argument annotations in SRL, we apply two argument augmentation techniques: Data augmentation by shuffling input texts, which reduces the model's reliance on syntax, and knowledge distillation from question answering (QA) data, which incorporates distant argument signals.

With experiments on the standard ACE and ERE English event benchmarks, we show that SRL annotations are valuable resources for EAE. With the template-based querying strategy, a model trained with SRL can reach nearly 80% of the fullysupervised F1 score in the zero-shot scenario, and an intermediate-training scheme provides further benefits in the low-resource setting. The model also obtains promising results in extensions to crossdomain and multi-lingual scenarios, demonstrating its generalizability. Our work highlights the utility of SRL annotations in the context of downstream applications with limited direct annotations.

2 Method

2.1 Querying Methods

For either semantic roles or event arguments, we can view the extraction task as a role querying problem. Specifically, we are given a sequence of words $\mathbf{s} = \{w_1, ..., w_n\}$ as input contexts as well as a predicate or event trigger word w_e and the semantic frame or event type t. Each type is associated with a list of participating roles to be filled and our aim is to extract arguments from the input contexts for each role. We adopt one specific modeling simplification, that is, our model only predicts the syntactic head word of an argument. For EAE, a heuristic method is further adopted to expand from head words to spans: We simply include the head word's child that is linked with a MWE dependency relation² and has an uppercase first letter. We find that this heuristic works well in practice, covering nearly 95% of the argument spans in the ACE and ERE event datasets. We take this approach to make it easier to transfer across different schemes, which may have different annotation criteria on span ranges.

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

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

In this way, we can view both SRL and EAE as a role querying problem over the input words. Specifically, the probability³ of a candidate word w to be the argument filling a role r is:

$$p_r(w) = \frac{\exp(\lambda \mathbf{h}_w^T \mathbf{q}_r)}{\sum_{w' \in \mathbf{s} \cup \{\epsilon\}} \exp(\lambda \mathbf{h}_{w'}^T \mathbf{q}_r)}$$
 143

Here, \mathbf{h}_w denotes the representation vector of the word w and \mathbf{q}_r indicates the querying vector of the role r. We further include a scaling factor λ , which is fixed to $\frac{1}{\sqrt{d}}$, where d is the dimension of \mathbf{h} and \mathbf{q} , following the attention calculation in Transformer (Vaswani et al., 2017). We specify a dummy token ϵ to handle the cases where no arguments can be found for a role. This modeling scheme is flexible and allows different argument extraction strategies to be viewed *in a unified way*. In this work, we explore four strategies, as illustrated in Figure 2. Since these strategies are not totally novel, we give brief descriptions in the main content and please refer to Appendix A.1 for more details.

1) **CLF.** We start with querying based on traditional classification, which assigns to each role a

¹These are also known as *implicit arguments* (O'Gorman, 2019). While there are more fine-grained linguistic criteria, we take a simplified approximate approach by checking the syntactic distances between triggers and arguments.

²Multi-word expressions: {"fixed", "flat", "compound"}.

³We also tried more complex scoring functions than dot product, such as multi-layer perceptron or bi-affine scorers, but found similar results. We thus choose this simplest one.



Figure 2: Illustrations of different role querying strategies (for the "artifact" role), based on 1) CLF: classification, 2) MRC: machine reading comprehension, 3) GEN: generation and 4) TSQ: template-based slot querying.

non-contextualized vector. To allow transfer to different role names, we initialize the role vectors with average-pooled representations obtained by passing the role names individually to a pre-trained language model. We call this strategy classificationbased since the role vectors can be viewed as weights in a linear classifier. This corresponds to more traditional argument extraction methods (Ouchi et al., 2018; Ebner et al., 2020). One shortcoming of this strategy is that the query vectors are constructed without access to input contexts, limiting their representation ability.

160

161

162

163

164

165

167

168

170

171

2) MRC. Recently, the strategy of casting NLP 172 tasks as machine reading comprehension problems 173 (Rajpurkar et al., 2016, 2018) has been applied to 174 EAE (Liu et al., 2020; Du and Cardie, 2020; Li 175 et al., 2020; Feng et al., 2020; Lyu et al., 2021; Liu et al., 2021a). In this strategy, each role is queried 177 with a *contextualized* question that is encoded to-178 gether with the context. Unless otherwise specified, 179 we form the role questions using the templates of Liu et al. (2021a), which can be automatically gen-181 erated from the role names. Since each question 182 queries only one role, this strategy requires a full pass through the encoder with respect to each role, 184 raising concerns regarding its computational efficiency,⁴ as compared to CLF.

3) GEN. More recently, many approaches ex-187 tract arguments by sequence-to-sequence generation (Paolini et al., 2021; Li et al., 2021b; Hsu et al., 2021; Lu et al., 2021; Du et al., 2021; Huang et al., 2022). Specifically, Li et al. (2021b) and Hsu et al. 191 (2021) adopt a template-based generation strategy, 192 which aggregates the queries of all roles for an 193 event into one template sentence. This strategy is 194 promising since the template can contain all roles 195 and query them in one pass. Since arguments come 196

from input contexts, we further adopt a pointer network (Vinyals et al., 2015) for argument selection rather than generation through output vocabularies, fitting our unified querying framework. Because of the auto-regressive decoding scheme, this strategy can also suffer lower efficiency compared to CLF. 197

199

200

201

202

203

204

205

206

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

227

228

229

230

231

4) **TSO.** We further explore a strategy that fully exploits the representative powerful and querying efficiency of templates. We do not fill the templates with actual words in the context but simply keep the role names as placeholders. We concatenate this template with the context, then pass that representation to the encoder for contextualization. Finally, the contextualized representations of the role slots in the template are adopted as role query vectors. We refer to this strategy as Template-based Slot Querying (TSQ). This approach is similar to the contemporaneous work of Ma et al. (2022). Our approach to template querying differs primarily in that: 1) We concatenate both the template and the context and feed them to the encoder, allowing for bidirectional modeling, and; 2) Our models predict argument head words rather than spans to facilitate transfer, since unlike Ma et al. (2022) the focus of our work is transfer learning. Appendix B.4 includes further comparisons.

2.2 SRL Templates

We take PropBank⁵ (Palmer et al., 2005) and FrameNet (Baker et al., 1998) as our main SRL resources. To allow transfer across different schemes, we need to specify extra information required by the role querying strategies. In particular, templates are not included in SRL frame definitions and it is infeasible to manually specify them for hundreds to thousands of SRL frames. We adopt a semi-

⁴Please refer to Appendix B.8 for speed comparisons.

⁵We also include NomBank (Meyers et al., 2004) and map its frames to PropBank frames.

316

317

318

319

320

321

323

324

325

326

284

automatic method to construct the templates, with extra information collected from data statistics:

233

234

235

240

241

243

245

246

247

248

253

254

259

261

265

266

267

270

273

274

275

278

279

281

- Role names. We directly take the role labels of FrameNet which are already in natural language form. We further train a role classifier with the FrameNet data and apply it to PropBank instances. For each frame-specific role in Prop-Bank, we choose the most frequently predicted FrameNet label as its role name.
 - **Role orders.** For each frame-specific role, we calculate its average relative distance to the predicate word (in non-passive verbal usages), and order all the roles with respect to this distance.
 - **Preposition words.** For each frame-specific role, we also count how frequent preposition words are used for the role's arguments. If the frequency of a certain preposition surpasses a threshold, we prepend the preposition word before the role name in the template.

With these three types of extra information, we construct the templates by concatenating all the corresponding ordered pieces. For example, the "buy.01" PropBank frame gets a template of "<u>buyer</u> **buy** goods for recipient from seller for money in place". Please refer to Appendix A.2 for additional details of this process and examples.

Note that this process can be noisy. Nevertheless, the above three pieces provide complementary information for role specification: the role names provide semantic information, the role orders include syntactic word order information, and the prepositions give further hints. In practice, we find that most of the the generated templates are reasonably close to natural language. In this way, we are able to form similar queries for both SRL and EAE, tackling the label mismatch problem between different tasks.

2.3 Argument Augmentation

In addition to label mismatch, another discrepancy between SRL and EAE is that arguments in traditional SRL are syntactically constrained whereas event arguments can be extracted from any place in the context. Therefore, SRL models will have difficulties in predicting syntactically distant arguments. To mitigate this problem, we apply data augmentation (Feng et al., 2021) and knowledge distillation (Hinton et al., 2015) to augment distant arguments for SRL instances.

Firstly, we apply a simple data augmentation method by shuffling the input contexts. Since the

SRL arguments are constrained by syntax, we hypothesize that by distorting syntax in some way, the model can be trained to focus more on the semantic relations between the predicates and arguments, thus allowing predicting more distant arguments. To distort syntax, we randomly chunk the input context sequence with sizes randomly chosen from one to three at each time. Then these text chunks are shuffled, re-concatenated and fed to the pretrained model for contextualized encoding. Since our model selects argument head words which are still tractable, there is no change to the later processing except for word position re-indexing. We only apply this procedure during training and simply mix vanilla unshuffled data with the shuffled ones by an 1:1 ratio.

Moreover, we seek signals of distant arguments from question answering $(QA)^{6}$ datasets, such as SQuAD (Rajpurkar et al., 2016, 2018). In QA annotations, the answers are not constrained by syntax and can be freely picked from the full context, providing valuable resources for distant arguments (Liu et al., 2021a). Motivated by this, we train a QA model with the MRC strategy and predict the missing arguments for SRL instances. Instead of hard predictions, we store a soft probabilistic distribution over the context words for each role and utilize these for SRL training. To avoid noises from the QA predictions, we adopt two filters. Firstly, we only apply distillation for the unfilled roles according to SRL annotations. This is intuitive since the filled roles already have gold annotations. Moreover, we apply distillation only when the prediction is confident enough. We perform calibration to the QA model by temperature scaling (Guo et al., 2017) and adopt a probability threshold of 0.5. In this way, we could borrow the signals of distant arguments from the QA datasets to enhance SRL instances with potential missing distant arguments.

3 Experiments

3.1 Settings

We conduct our main experiments⁷ with English ACE⁸ (Walker et al., 2006) and ERE (LDC, 2015) event datasets. We follow Lin et al. (2020) and utilize their pre-processing scripts. For the target

⁶Specifically we adopt the extractive QA-MRC data. To avoid confusion, we use "QA" when denoting data resources while using "MRC" for the querying strategy.

The link to our implementation is omitted for review.

 $^{^{8}}$ We adopt ACE05-E⁺ (Lin et al., 2020) which keeps pronouns as arguments.

| Method | P% | ACE R% | F1% | P% | ERE R% | F1% |
|--|---|---|--|---|---|--|
| Super. | $68.93_{\pm 1.07}$ | $68.94_{\pm0.95}$ | $68.93_{\pm0.95}$ | $72.75_{\pm 1.69}$ | $71.80_{\pm 1.29}$ | $72.24_{\pm0.34}$ |
| GPT-3 QA | 29.10 32.77 _{±3.70} | $\begin{array}{r} 34.25 \\ 47.43_{\pm 1.17} \end{array}$ | $\begin{array}{c} 31.47 \\ 38.62 _{\pm 2.58} \end{array}$ | $\begin{array}{c} 25.09 \\ 32.68_{\pm 2.78} \end{array}$ | $26.76 \\ 48.13_{\pm 4.08}$ | $\begin{array}{c} 25.90 \\ 38.74_{\pm 2.09} \end{array}$ |
| SRL _{CLF} SRL _{MRC} SRL _{GEN} | $\begin{array}{c} 47.97_{\pm 1.47} \\ 58.27_{\pm 0.75} \\ 55.77_{\pm 0.61} \end{array}$ | $\begin{array}{c} 25.37_{\pm 0.86} \\ 39.54_{\pm 1.60} \\ 45.31_{\pm 1.26} \end{array}$ | $\begin{array}{c} 33.18 _{\pm 0.92} \\ 47.08 _{\pm 0.89} \\ 49.99 _{\pm 0.93} \end{array}$ | $\begin{array}{c} 50.17_{\pm 1.72} \\ 62.02_{\pm 1.15} \\ 58.37_{\pm 0.66} \end{array}$ | $\begin{array}{c} 25.60 _{\pm 0.65} \\ 45.31 _{\pm 1.74} \\ 52.68 _{\pm 0.63} \end{array}$ | $\begin{array}{c} 33.89_{\pm 0.88} \\ 52.32_{\pm 0.83} \\ 55.38_{\pm 0.62} \end{array}$ |
| SRL _{TSQ} +shuf. +distill +both | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 49.61_{\pm 0.80} \\ 51.70_{\pm 0.52} \\ 55.85_{\pm 0.67} \\ 57.04_{\pm 0.93} \end{array}$ | $\begin{array}{c} 53.36 {\scriptstyle \pm 0.53} \\ 54.82 {\scriptstyle \pm 0.44} \\ 55.17 {\scriptstyle \pm 0.42} \\ \textbf{56.35} {\scriptstyle \pm 1.07} \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 55.84_{\pm 0.78} \\ 57.42_{\pm 1.26} \\ 60.90_{\pm 0.85} \\ 61.48_{\pm 0.18} \end{array}$ | $\begin{array}{c} 57.81 {\scriptstyle \pm 0.34} \\ 58.54 {\scriptstyle \pm 1.05} \\ 57.95 {\scriptstyle \pm 0.65} \\ \textbf{58.96} {\scriptstyle \pm 0.48} \end{array}$ |

Table 1: Zero-shot EAE results on event test sets. Except for GPT-3, all results are averaged over five runs.

event frames, we manually specify extra information such as templates, adopting those of Li et al. (2021b). Unless otherwise specified, we assume given gold event triggers and focus on the extraction of event arguments. We also provide results with predicted event triggers in Appendix B.3. We evaluate arguments by labeled F1 scores, which require both argument spans and roles to match the gold ones. We run with five random seeds and report averaged results.

327

328

331

332

333

338

341

342

345

347

348

351

353

354

357

361

For external data, we take PropBank, Nom-Bank 1.0 and FrameNet 1.7 as our main SRL resources. We prepare the SRL templates by the semi-automatic process described in §2.2. For QA datasets, we take SQuAD 2.0 (Rajpurkar et al., 2018), QA-SRL 2.1 (FitzGerald et al., 2018), QANom (Klein et al., 2020) and QAMR (Michael et al., 2018). For the training of SRL or QA models, we simply adopt the concatenation of all the corresponding datasets. Except for those that have manual syntactic annotations, we utilize Stanza (Qi et al., 2020) to parse the texts to obtain the syntactic head words of the arguments.

We adopt pre-trained language models for initialization and fine-tune the full models during training. Specifically, we use RoBERTa_{base} (Liu et al., 2019) for encoder-only models (CLF, MRC, TSQ) and BART_{base} (Lewis et al., 2020) for encoder-decoder models (GEN). Please refer to Appendix B.1 for more detailed experimental settings.

3.2 Main Transfer Experiments

We conduct our main experiments with English ACE and ERE datasets. Thanks to the unified querying framework, we are able to conduct experiments in a zero-shot setting (§3.2.1), where models trained on external data are directly evaluated on EAE. We also investigate low-resource settings where some amounts of EAE annotations are available for further fine-tuning (§3.2.2).

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

385

388

389

390

391

392

393

394

3.2.1 Zero-shot

In the zero-shot setting, we further compare with two methods in addition to SRL: 1) GPT-3 (Brown et al., 2020), where we form prompts⁹ for each role and use GPT-3 to generate the answers; 2) QA, where we train QA models¹⁰ with the QA datasets. We also provide the fully-supervised results¹¹ (Super.) as references.

Results The main results are shown in Table 1. Except for the one with the CLF strategy, SRL models perform generally better than QA and GPT-3, showing the effectiveness of utilizing SRL resources. Among the SRL models, the TSQ strategy generally performs the best, indicating the effectiveness of this in-context querying strategy. Further improvements can be obtained with the argument augmentation techniques. Interestingly, if only using shuffling augmentation (+shuf.), the precisions roughly keep the same while recalls increase. If only using distillation (+distill), the recalls get boosted with the sacrifice of precisions. Finally, if both are utilized (+both), the improvements on recalls get accumulated while the precisions slightly recover when compared with the distillation-only case. This leads to the overall best F1 scores, reaching around 80% of the supervised results.¹²

Analysis As shown in Figure 3, we further perform breakdowns on the syntactic distances between triggers and arguments. We especially com-

¹²Please refer to Appendix B.6 for manual analysis.

⁹Please refer to Appendix B.2 for more GPT-3 details.

¹⁰Notice that we can only use the MRC strategy for QA models because of the task-specific format.

¹¹We take those of the TSQ model. More details of the supervised results are provided in Appendix B.5.



Figure 3: Breakdowns on trigger-argument syntactic distances (on ACE dev set). Numbers in the parentheses denote the percentages in the gold annotations.

pare the QA model and the four SRL_{TSO} models. Firstly, the QA model performs worse than SRL models except for the long distant ones (" $d \ge 4$ "). This is due to that SRL annotations mainly capture syntactically local arguments while QA is not constrained by this. Within the SRL models, when adding shuffling ("+shuf.") or distillation ("+distill"), the middle-ranged arguments consistently obtain improvements. One interesting pattern is that shuffling benefits "d = 1" but hurts "d > 4", while distillation seems to have the opposite effects. This may indicate that shuffling enhances more robust predictions of short- and middle-ranged arguments while distillation encourages longer-ranged ones. Finally, when combining these two techniques ("+both"), the model can reach a good balance, achieving the best overall results. Due to its overall better performance, we will use the "SRL_{TSO}+both" strategy for our SRL models in the remaining of this work.

3.2.2 Low-resource

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

417

420

421

422

423

424

425

427

429

We further investigate scenarios where we have 416 some amounts of target EAE annotations. With target data, we can directly train an EAE model (from 418 pre-trained language models). We further apply a 419 simple intermediate-training scheme (Phang et al., 2018; Wang et al., 2019a) to transfer the knowledge from SRL. We take the SRL-trained model and further fine-tune it on the target event data. A similar scheme can also be adopted with the QA model. Figure 4 shows the results with different amounts of training instances. Generally, SRL intermedi-426 ate training is beneficial especially for middle- and low-resource cases, again showing that SRL an-428 notations can be valuable transfer sources for the extraction of event arguments. 430



Figure 4: Model performance with direct or intermediate training. Here x-axis (drawn in log scale) denotes the percentage of utilized training data. The shaded areas indicate the ranges of standard deviations.

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

3.3 Further Extensions

In the previous experiments, we take ACE and ERE as the targets, which are still relatively similar to the SRL annotations. In this sub-section, we further investigate scenarios where there are larger discrepancies between the source and the target. Specifically, we examine the transfer from SRL to EAE in cross-domain (§3.3.1), multi-lingual (§3.3.2) and multi-sentence $(\S3.3.3)$ cases.

3.3.1 Cross-domain

We first investigate the biomedical domain, utilizing the GENIA BioNLP-11 benchmark (Kim et al., 2011). The GENIA events are quite different than general SRL frames and mainly describe detailed bio-molecule behavior (Kim et al., 2008). Still focusing on the argument extraction step, we take the event triggers predicted by the supervised system BEESL (Ramponi et al., 2020). We perform zeroshot argument extraction and evaluate the QA and SRL models, with manually compiled role questions and templates. We adopt the official evaluation metric of approximate recursive matching. Please refer to Appendix C.1 for more details.

Our main comparison is between the QA and SRL models, while we also include the supervised results of BEESL as references. We further adopt a self-training approach where we predict SRL

| Types | QA | SRL | SRL+self | Super. |
|-----------------|-------|-------|----------|--------|
| Expression | 70.71 | 76.66 | 77.59 | 80.90 |
| Transcription | 63.64 | 55.63 | 59.72 | 69.46 |
| Catabolism | 62.07 | 66.67 | 66.67 | 74.07 |
| Phosphorylation | 75.95 | 78.98 | 83.64 | 89.52 |
| Localization | 53.28 | 66.89 | 67.33 | 69.51 |
| – Simple – | 68.26 | 73.63 | 75.27 | 79.31 |
| Binding | 39.41 | 34.90 | 35.10 | 50.19 |
| Regulation | 33.80 | 38.95 | 38.52 | 45.90 |
| Pos. regulation | 31.85 | 38.96 | 39.95 | 49.41 |
| Neg. regulation | 36.62 | 44.84 | 44.51 | 47.17 |
| – Complex – | 33.36 | 40.44 | 40.88 | 48.32 |
| – All – | 47.42 | 51.95 | 52.76 | 60.22 |

Table 2: BioNLP-11 event extraction results (F1%).

frames on the unlabeled texts from the original GENIA training set and include these predicted structures for SRL training.

The results on the test set are shown in Table 2. SRL generally outperforms QA for most of the types. This may due to the difficulty of asking proper questions. For example, for the "Regulation" event, we ask "What is regulated?" for the role of "Theme" and "What causes the regulation?" for "Cause". These questions may be unrelated to the actual contexts, while for the SRL models, extra hints from the query templates may be helpful. This may also explain the reason why QA is better on some of the types where it is relatively easy to ask questions. For example, for "Transcription", the question of "What is transcribed?" would be accurate for most contexts. For the SRL models, the self-training method is beneficial overall, showing the effectiveness of utilizing unlabeled corpus from the target domain.¹³ Finally, our best zero-shot model could recover more than 80% of the overall performance of the supervised model, showing that general SRL resources can still be helpful in the biomedical domain. The main gaps between the zero-shot and supervised systems are on the "Binding" and Complex events where there are complicated and even nested structures. One future direction is to investigate ways to better handle these complex structures.

3.3.2 Multi-lingual

We next explore a multi-lingual setting, taking ACE05 Arabic and Chinese datasets as our targets. We follow Huang et al. (2022) and utilize their pre-precessing scripts¹⁴ for data preparation.

| Model | Arabic | Chinese | | | | |
|---|------------------------|--------------------------|--|--|--|--|
| Zero-shot results without any EAE annotations. | | | | | | |
| QA _{en} | $ 22.56_{\pm 1.48} $ | $26.58_{\pm 2.61}$ | | | | |
| QA _{en+tgt} | 23.54 ± 1.43 | 27.08 ± 1.79 | | | | |
| SRLen | $37.75_{\pm 0.52}$ | $39.37_{\pm 1.45}$ | | | | |
| SRL _{en+tgt} | 40.64 ±1.49 | 41.50 ± 1.04 | | | | |
| Multi-lingual | results with | English EAE annotations. | | | | |
| $GATE^{\dagger}$ | 44.5 | 49.2 | | | | |
| X-Gear [†] | 44.8 | 54.0 | | | | |
| En _{MRC} | $37.44_{\pm 3.02}$ | $51.86_{\pm 0.92}$ | | | | |
| +QA _{en} | $39.06_{\pm 2.86}$ | $53.36_{\pm 1.06}$ | | | | |
| +QA _{en+tgt} | $44.27_{\pm 1.37}$ | $53.97_{\pm 1.41}$ | | | | |
| En _{TSQ} | 37.64±1.96 | $53.54_{\pm 0.65}$ | | | | |
| +SRL _{en} | $41.86_{\pm 0.92}$ | $53.96_{\pm 0.85}$ | | | | |
| +SRL _{en+tgt} | 51.51 ±1.32 | 58.90 ±0.76 | | | | |
| Supervised results with target EAE annotations. | | | | | | |
| Super. | 58.09 _{±1.51} | $65.11_{\pm 0.94}$ | | | | |

Table 3: Results (Argument F1%) on ACE05 Arabic and Chinese test sets. "†" denotes reported results from Huang et al. (2022).

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

We further include multi-lingual external resources. For SRL, we utilize Arabic and Chinese PropBank annotations from OntoNotes (Hovy et al., 2006; Weischedel et al., 2013). For the role names in SRL frames, we again adopt a statistical approach: predicting with a FrameNet classifier based on a multilingual pre-trained encoder and adopting the mostly predicted label for each role. Due to differences in word order and usage of prepositional words in non-English languages, we exclude preposition words and simply order the roles by their ARG numbers.¹⁵ We also include QA datasets for the target languages, adopting CMRC-2018 (Cui et al., 2019) for Chinese and the Arabic portion of TyDiQA (Clark et al., 2020) for Arabic. All our models in this experiment are based on the pretrained multilingual encoder of XLM-R_{base} (Conneau et al., 2020).

The results are shown in Table 3. In the first group, we compare zero-shot performance without any EAE training resources. Similar to the previous trends, SRL models are obviously better than QA models, while including annotations in the target language could provide further benefits. In the second group, we assume access to English EAE training data. Similar to §3.2.2, we adopt an intermediate-training scheme by further finetuning the QA or SRL model on the English EAE

488

489

490

¹³We also tried masked-language-model objectives using the unlabeled target data, but did not find improvements.

¹⁴https://github.com/PlusLabNLP/X-Gear

¹⁵The Arabic and Chinese frames adopt similar schemes as in English, specifying roles of {ARG0, ARG1, ...}. We find it reasonable by simply ordering them by the role numbers and forming templates of "ARG0 V ARG1 ARG2 ...".

| Model | Overall | Same-Sent. | Cross-Sent. |
|-------------------------|--|---|--|
| QA SRL SRL+pseudo | $ \begin{vmatrix} 28.23_{\pm 0.74} \\ 48.03_{\pm 0.30} \\ 48.00_{\pm 0.14} \end{vmatrix} $ | $\begin{array}{c c} 35.16_{\pm 1.42} \\ 53.36_{\pm 0.30} \\ 53.50_{\pm 0.16} \end{array}$ | $\begin{array}{c} 11.66_{\pm 0.69} \\ 2.81_{\pm 0.78} \\ 11.17_{\pm 1.88} \end{array}$ |
| Super. | $ 57.38_{\pm 0.84} $ | $63.45_{\pm 0.86}$ | $25.52_{\pm 1.31}$ |

Table 4: Argument head F1(%) on RAMS test set.

data. Compared with the results of directly training on English, intermediate-training with external resources could bring improvements. Again we see that models enhanced with SRL resources obtain the overall best results, which are quite promising when compared with the supervised ones.

3.3.3 Multi-sentence

520

521

522

523

525

526

528

530

535

537

540

541

542

543

545

546

547

551

552

553

557

559

560

Finally, we investigate multi-sentence event arguments, which are not constrained in the same sentence of the event trigger but can come from the document-level contexts. To investigate this phenomenon, we evaluate¹⁶ on the RAMS dataset (Ebner et al., 2020), which annotates event arguments within five-sentence windows around the triggers. We similarly extend contexts to fivesentence windows in our training of QA and SRL models for this experiment.

The zero-shot results are shown in the first group of Table 4. Consistent with our previous findings, SRL performs better than QA for same-sentence arguments. Nevertheless, it predicts very few crosssentence arguments. This is not surprising because there are no such signals in the SRL training data. Inspired by previous works on coreference and anaphora resolution (Varkel and Globerson, 2020; Konno et al., 2021), we create pseudo SRL data with cross-sentence arguments by surface-string matching. Specifically, for each nominal argument in an SRL instance, we search for words in nearby sentences that have the same lemma as the argument's head word. If there are, we delete the original true argument and add pseudo cross-sentence argument links to those matched words. Although deletion may create ungrammatical instances, we find it better than other schemes such as replacing the original argument with a "[MASK]" token. With the additional pseudo training data, the model could recover certain cross-sentence arguments while keeping similar same-sentence performance. Multi-sentence argument extraction is still a difficult task, where even the supervised system

could only obtain an F1 score of around 25%. This calls for more future investigations, and exploring how to better utilize external data resources such as SRL might be a promising direction.

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

583

584

585

586

588

589

590

591

592

593

594

595

596

597

598

600

601

602

603

604

605

606

607

609

4 Related Work

Utilizing shallow semantics for the event tasks has been explored previously. Liu et al. (2016) leverage the FrameNet frames to enhance event detection. Wang et al. (2021) conduct contrastive pre-training with AMR structures to enhance event extraction. Several works utilize predicted shallow semantic structures as inputs to help low-resource event extraction (Peng et al., 2016; Huang et al., 2018; Lyu et al., 2021) and event schema induction (Huang et al., 2016). This work focuses on the sub-task of EAE and shows that SRL can be a valuable direct training resource for EAE.

For the EAE task, most previous works adopt methods with a classification-based strategy where each role is assigned static querying parameters (Chen et al., 2015; Nguyen et al., 2016; Wang et al., 2019b; Pouran Ben Veyseh et al., 2020; Ma et al., 2020; Ebner et al., 2020). Recently, two interesting alternative strategies are explored to enable extraction in more flexible ways: MRC-based methods cast the problem as answering role questions (Liu et al., 2020; Du and Cardie, 2020; Li et al., 2020; Feng et al., 2020; Lyu et al., 2021; Liu et al., 2021a), while generation-based methods adopt sequence-tosequence generation schemes (Paolini et al., 2021; Li et al., 2021b; Hsu et al., 2021; Lu et al., 2021; Du et al., 2021; Huang et al., 2022). We cover all these strategies within a unified role querying framework and further explore a template-based role querying strategy. This strategy is also related with prompt-based learning (Liu et al., 2021b; Schick and Schütze, 2021; Li and Liang, 2021; Petroni et al., 2019), but differs in the extractiontargeted paradigm. Concurrently, Ma et al. (2022) adopt a very similar idea, while our work differs mainly on our focus upon transfer learning.

5 Conclusion

In this work, we explore transfer learning from semantic roles to event arguments. With unified role querying strategies, we show that SRL annotations could be valuable resources to help the extraction of event arguments. The SRL model could also obtain promising results when extended to new scenarios with domain and language differences.

¹⁶Since our head-expanding heurist does not cover the argument span annotation conventions on RAMS, for simplicity we only evaluate argument head words.

610 References

611

613

614

617

618

619

622

623

629

641

643

651

660

661

664

- David Ahn. 2006. The stages of event extraction. In Proceedings of the Workshop on Annotating and Reasoning about Time and Events, pages 1–8, Sydney, Australia. Association for Computational Linguistics.
- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. The Berkeley FrameNet project. In 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1, pages 86–90, Montreal, Quebec, Canada. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. Event extraction via dynamic multipooling convolutional neural networks. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 167–176, Beijing, China. Association for Computational Linguistics.
- Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: 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, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng Chen, Wentao Ma, Shijin Wang, and Guoping Hu.
 2019. A span-extraction dataset for Chinese 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 5883–5889, Hong Kong, China. Association for Computational Linguistics.
- Xinya Du and Claire Cardie. 2020. Event extraction by answering (almost) natural questions. In *Proceedings* of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 671–683, Online. Association for Computational Linguistics.

Xinya Du, Alexander Rush, and Claire Cardie. 2021. GRIT: Generative role-filler transformers for document-level event entity extraction. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 634–644, Online. Association for Computational Linguistics. 666

667

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

703

704

705

706

707

708

709

710

711

712

713

714

715

718

- Seth Ebner, Patrick Xia, Ryan Culkin, Kyle Rawlins, and Benjamin Van Durme. 2020. Multi-sentence argument linking. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8057–8077, Online. Association for Computational Linguistics.
- Rui Feng, Jie Yuan, and Chao Zhang. 2020. Probing and fine-tuning reading comprehension models for few-shot event extraction. *arXiv preprint arXiv:2010.11325*.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey of data augmentation approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- Nicholas FitzGerald, Julian Michael, Luheng He, and Luke Zettlemoyer. 2018. Large-scale QA-SRL parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2051–2060, Melbourne, Australia. Association for Computational Linguistics.
- Daniel Gildea and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245–288.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1321– 1330. PMLR.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. 2015. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2(7).
- Eduard Hovy, Mitchell Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2006. OntoNotes: The 90% solution. In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, pages 57–60, New York City, USA. Association for Computational Linguistics.
- I Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott Miller, Prem Natarajan, Kai-Wei Chang, Nanyun Peng, et al. 2021. Degree: A data-efficient generative event extraction model. *arXiv preprint arXiv:2108.12724*.

820

821

822

823

824

825

826

827

828

829

830

831

776

Kuan-Hao Huang, I Hsu, Premkumar Natarajan, Kai-Wei Chang, Nanyun Peng, et al. 2022. Multilingual generative language models for zero-shot crosslingual event argument extraction. *arXiv preprint arXiv:2203.08308*.

720

721

725

727

732

733

734

736

737 738

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

761

764

767

768

770

775

- Lifu Huang, Taylor Cassidy, Xiaocheng Feng, Heng Ji, Clare R. Voss, Jiawei Han, and Avirup Sil. 2016.
 Liberal event extraction and event schema induction. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 258–268, Berlin, Germany. Association for Computational Linguistics.
- Lifu Huang, Heng Ji, Kyunghyun Cho, Ido Dagan, Sebastian Riedel, and Clare Voss. 2018. Zero-shot transfer learning for event extraction. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2160–2170, Melbourne, Australia. Association for Computational Linguistics.
- Jin-Dong Kim, Tomoko Ohta, and Jun'ichi Tsujii. 2008. Corpus annotation for mining biomedical events from literature. *BMC bioinformatics*, 9(1):1–25.
- Jin-Dong Kim, Yue Wang, Toshihisa Takagi, and Akinori Yonezawa. 2011. Overview of Genia event task in BioNLP shared task 2011. In *Proceedings of BioNLP Shared Task 2011 Workshop*, pages 7–15, Portland, Oregon, USA. Association for Computational Linguistics.
- Ayal Klein, Jonathan Mamou, Valentina Pyatkin, Daniela Stepanov, Hangfeng He, Dan Roth, Luke Zettlemoyer, and Ido Dagan. 2020. QANom: Question-answer driven SRL for nominalizations. In Proceedings of the 28th International Conference on Computational Linguistics, pages 3069–3083, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Ryuto Konno, Shun Kiyono, Yuichiroh Matsubayashi, Hiroki Ouchi, and Kentaro Inui. 2021. Pseudo zero pronoun resolution improves zero anaphora resolution. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3790–3806, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- LDC. 2005. ACE (automatic content extraction) english annotation guidelines for events version 5.4.3. *Linguistic Data Consortium*.
- LDC. 2015. Deft Rich ERE annotation guidelines: Events version 3.0. *Linguistic Data Consortium*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020.
 BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,

pages 7871–7880, Online. Association for Computational Linguistics.

- Fayuan Li, Weihua Peng, Yuguang Chen, Quan Wang, Lu Pan, Yajuan Lyu, and Yong Zhu. 2020. Event extraction as multi-turn question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 829–838, Online. Association for Computational Linguistics.
- Manling Li, Sha Li, Zhenhailong Wang, Lifu Huang, Kyunghyun Cho, Heng Ji, Jiawei Han, and Clare Voss. 2021a. The future is not one-dimensional: Complex event schema induction by graph modeling for event prediction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5203–5215, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sha Li, Heng Ji, and Jiawei Han. 2021b. Documentlevel event argument extraction by conditional generation. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 894–908, Online. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. 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 4582– 4597, Online. Association for Computational Linguistics.
- Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. A joint neural model for information extraction with global features. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7999–8009, Online. Association for Computational Linguistics.
- Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. Event extraction as machine reading comprehension. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1641–1651, Online. Association for Computational Linguistics.
- Jian Liu, Yufeng Chen, and Jinan Xu. 2021a. Machine reading comprehension as data augmentation: A case study on implicit event argument extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2716– 2725, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.

- 832 833
- 83
- 83
- 839
- 841
- 8
- 8
- 8
- 849 850
- 8
- 8
- 8 8
- 8 8

86

- 8
- 8

87

01

- 874 875
- 8
- 8
- 8
- 8

8

- 884 885 886
- 88 88

887 888

- Shulin Liu, Yubo Chen, Shizhu He, Kang Liu, and Jun Zhao. 2016. Leveraging FrameNet to improve automatic event detection. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2134– 2143, Berlin, Germany. 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.
- Yaojie Lu, Hongyu Lin, Jin Xu, Xianpei Han, Jialong Tang, Annan Li, Le Sun, Meng Liao, and Shaoyi Chen. 2021. Text2Event: Controllable sequence-tostructure generation for end-to-end event extraction. 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 2795–2806, Online. Association for Computational Linguistics.
- Qing Lyu, Hongming Zhang, Elior Sulem, and Dan Roth. 2021. Zero-shot event extraction via transfer learning: Challenges and insights. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 322–332, Online. Association for Computational Linguistics.
- Jie Ma, Shuai Wang, Rishita Anubhai, Miguel Ballesteros, and Yaser Al-Onaizan. 2020. Resource-enhanced neural model for event argument extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3554–3559, Online. Association for Computational Linguistics.
- Yubo Ma, Zehao Wang, Yixin Cao, Mukai Li, Meiqi Chen, Kun Wang, and Jing Shao. 2022. Prompt for extraction? paie: Prompting argument interaction for event argument extraction. *arXiv preprint arXiv:2202.12109*.
- Adam Meyers, Ruth Reeves, Catherine Macleod, Rachel Szekely, Veronika Zielinska, Brian Young, and Ralph Grishman. 2004. The NomBank project: An interim report. In *Proceedings of the Workshop Frontiers in Corpus Annotation at HLT-NAACL 2004*, pages 24–31, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Julian Michael, Gabriel Stanovsky, Luheng He, Ido Dagan, and Luke Zettlemoyer. 2018. Crowdsourcing question-answer meaning representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 560–568, New Orleans, Louisiana. Association for Computational Linguistics.

Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016. Joint event extraction via recurrent neural networks. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 300–309, San Diego, California. Association for Computational Linguistics.

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

- Timothy J O'Gorman. 2019. *Bringing together computational and linguistic models of implicit role interpretation*. Ph.D. thesis, University of Colorado at Boulder.
- Hiroki Ouchi, Hiroyuki Shindo, and Yuji Matsumoto. 2018. A span selection model for semantic role labeling. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 1630–1642, Brussels, Belgium. Association for Computational Linguistics.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The Proposition Bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.
- Martha Palmer, Daniel Gildea, and Nianwen Xue. 2010. Semantic role labeling. *Synthesis Lectures on Human Language Technologies*, 3(1):1–103.
- Sinno Jialin Pan and Qiang Yang. 2009. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359.
- Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, RISHITA ANUBHAI, Cicero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. Structured prediction as translation between augmented natural languages. In *International Conference on Learning Representations*.
- Haoruo Peng, Yangqiu Song, and Dan Roth. 2016. Event detection and co-reference with minimal supervision. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 392–402, Austin, Texas. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Jason Phang, Thibault Févry, and Samuel R Bowman. 2018. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. *arXiv preprint arXiv:1811.01088*.
- Amir Pouran Ben Veyseh, Tuan Ngo Nguyen, and Thien Huu Nguyen. 2020. Graph transformer networks with syntactic and semantic structures for

- 945
- 94

955

960

961

962

965

968

969

970

971

978

979

987

994

995

998

event argument extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3651–3661, 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, pages 101–108, Online. Association for Computational Linguistics.
 - Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
 - 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.
 - Alan Ramponi, Rob van der Goot, Rosario Lombardo, and Barbara Plank. 2020. Biomedical event extraction as sequence labeling. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5357–5367, Online. Association for Computational Linguistics.
 - Sebastian Ruder, Matthew E. Peters, Swabha Swayamdipta, and Thomas Wolf. 2019. Transfer learning in natural language processing. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Tutorials, pages 15–18, Minneapolis, Minnesota. Association for Computational Linguistics.
 - Timo Schick and Hinrich Schütze. 2021. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the* 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 255–269, Online. Association for Computational Linguistics.
 - Yuval Varkel and Amir Globerson. 2020. Pre-training mention representations in coreference models. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8534–8540, Online. Association for Computational Linguistics.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. *Advances in neural information processing systems*, 28.

1000

1001

1003

1004

1005

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

- Christopher Walker, Stephanie Strassel, Julie Medero, and Kazuaki Maeda. 2006. ACE 2005 multilingual training corpus. *Linguistic Data Consortium*, 57.
- Alex Wang, Jan Hula, Patrick Xia, Raghavendra Pappagari, R. Thomas McCoy, Roma Patel, Najoung Kim, Ian Tenney, Yinghui Huang, Katherin Yu, Shuning Jin, Berlin Chen, Benjamin Van Durme, Edouard Grave, Ellie Pavlick, and Samuel R. Bowman. 2019a. Can you tell me how to get past sesame street? sentence-level pretraining beyond language modeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4465–4476, Florence, Italy. Association for Computational Linguistics.
- Xiaozhi Wang, Ziqi Wang, Xu Han, Zhiyuan Liu, Juanzi Li, Peng Li, Maosong Sun, Jie Zhou, and Xiang Ren. 2019b. HMEAE: Hierarchical modular event argument extraction. 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 5777–5783, Hong Kong, China. Association for Computational Linguistics.
- Ziqi Wang, Xiaozhi Wang, Xu Han, Yankai Lin, Lei Hou, Zhiyuan Liu, Peng Li, Juanzi Li, and Jie Zhou. 2021. CLEVE: Contrastive Pre-training for Event Extraction. 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 6283–6297, Online. Association for Computational Linguistics.
- Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2013. Ontonotes release 5.0. *Linguistic Data Consortium, Philadelphia, PA*, 23.

1055

1056

1057

1058

1061

1062

1063

1065

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1085

1086

1087

1088

1089

1040

A Details of Methods

A.1 Modeling Details

This sub-section provides more details of the models and the querying strategies that are briefly described in §2.1.

We first introduce some common modeling settings before diving into specific querying strategies. As described in the main context, we adopt a unified role querying model for argument extraction with representations of the role queries q_r and the candidate words \mathbf{h}_w . The construction of word representations follows common practice: we feed the input sequence to a contextualized encoder and utilize each word's output hidden vector. When a word is split into multiple sub-words, the first subword is taken. To encode the trigger word, the input embedding of the trigger is added with a specific trigger embedding, which is randomly initialized and tuned together with the model. There are cases when the context does not have mentions for some roles (no arguments), where we adopt an all-zero dummy \mathbf{h}_{ϵ} , which essentially fixes the no-argument scores to zero.

During training, we use the standard crossentropy loss function. When there are more than one gold arguments, we simply apply equal weights to them. In testing, for each role, we select the words whose score is larger than zero and ranks within the top-two among all candidate words. One important aspect that we do not explicitly consider in the output modeling is the interactions between arguments as well as frame-level global features, which have been shown effective for event extraction (Lin et al., 2020) and event schema induction (Li et al., 2021a). Incorporating these for the transfer scenarios would be an interesting future direction. The selected words are further expanded to argument spans using the dependency-tree based heuristic as described in the main content.

The main difference among the querying strategies is on the construction of the querying vectors, which is described in the following.

1) CLF

For the traditional classification based strategy, we allocate a specific vector for each role, which is included as model parameters. In the case where there are enough supervision, these vectors can be randomly initialized. To fit our goal of transfer learning, we take advantage of the natural language role names and encode them individually using a

| Role | Question Template |
|--------|--|
| Person | Who is the [] _{role_name} in the [] _{trigger_text} event? |
| Place | Where does the [] _{trigger_text} event take place? |
| Others | What is the [] _{role_name} of the [] _{trigger_text} event? |

Table 5: Question templates from Liu et al. (2021a).

vanilla pre-trained model, with an input format of: 1090

$$[CLS] \ role_name \ [SEP]$$
 1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

For example, for the role of "artifact", the input is simply "[CLS] artifact [SEP]". We take the averagely-pooled output representations. Notice that since each role name is encoded by itself without any other contexts, the representations are non-contextualized, making this strategy almost the same as using a classifier.

2) MRC

For the MRC based strategy, we form a question for each role and dynamically obtain the query vectors by encoding the question together with the context. We adopt the question templates from Liu et al. (2021a), as shown in Table 5. For example, to query the "artifact" role of the "bought" event in Figure 1, we ask: "What is the artifact of the bought event?". One advantage of this strategy is that we only need role types, role names and trigger texts to form a question, making it less difficult to extend to the SRL cases. In our preliminary experiments, we also tried role-specific questions for ACE utilizing those from Lyu et al. (2021), such as "What is bought?", and found similar results. Following standard MRC models, we concatenate the question and the context as the input sequence, which is fed to the encoder:

[CLS] role_question [SEP] context [SEP]

Furthermore, instead of introducing extra parameters with an extra answer selection head, we simply take the contextualized representations of the question word¹⁷ as the querying vector.

3) GEN

The template generation based strategy requires a template for each event or SRL frame, which specifies a canonical realization of this frame in a

¹⁷We choose the question word instead of the role name to allow easier transfer from QA datasets where there may be no specific querying roles. In preliminary experiments, we also tried average pooling over the question tokens to form querying vectors but did not find better results.

natural language sentence. For example, for the 1126 "TransferOwnership" event, we have a template 1127 of "seller give artifact to buyer for beneficiary in 1128 place", where each role occupies a placeholder slot. 1129 In this strategy, a sequence-to-sequence encoder-1130 decoder model is utilized. The context is encoded 1131 by the encoder while the filled template is gener-1132 ated by the decoder. We mostly follow Li et al. 1133 (2021b) but make some modifications to the out-1134 put modeling. Instead of directly replacing the 1135 slots with actual argument words, we keep the role 1136 names and insert the actual arguments after the role 1137 slots. For example, we output¹⁸ "seller [UNK] give 1138 artifact book to buyer he for beneficiary [UNK] 1139 in place store" instead of "[UNK] give book to he 1140 for [UNK] in store". We keep the role names for 1141 two reasons: firstly, the role names in the target 1142 sequence can act as a guidance of the to-be-filled 1143 arguments; moreover, since the arguments are re-1144 stricted to be words from the context, we can utilize 1145 the representations of the role names as queries to 1146 point to the context words. The second point allows 1147 us to form a pointer-network styled model, which 1148 directly selects arguments from the context word 1149 representations, fitting in our unified role-querying 1150 framework. 1151

4) TSQ

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

Actually we do not need to fill in the template with actual argument words, since our target task is an extraction task where we only need to find the argument mentions in the context. Moreover, if no generation is required, we could merge the context and the template together to allow bidirectional modeling. Motivated by this, we keep the unfilled but already natural language styled template as it is, concatenate it with the context and feed the full sequence to the encoder:

[CLS] template [SEP] context [SEP]

After the encoding, we take the output representations (first sub-token) of each role slot as its query vector and apply all the role queries parallely to select the corresponding argument words. This can be viewed as a combination of MRC and GEN, taking the advantages of both methods. As in MRC, we perform extraction for the role queries and no generation is needed, and as in GEN, the template allows us to embed all the role queries in one sequence rather than forwarding multiple times for differ-
ent roles. Moreover, since all the role queries are
performed parallelly without inter-dependencies,
this can be viewed as a non-autoregressive method1174understand
which is more efficient than GEN.1176

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

Extra Information Requirements

To allow transferring among different semantic frame schemes, all the above querying strategies require some extra information for the event or SRL frames. Firstly, the natural-language styled role names are needed for all the methods, since without this, we can hardly model what each role aims to query. With the pre-trained language models, we have a way to obtain role representations within a shared space by encoding the natural language styled queries. For MRC, we further need the role types, more specifically, whether the role is personrelated to decide using "What" or "Who".¹⁹ For GEN and TSQ, we need templates that embed all the roles within a natural language sentence.

For the target EAE tasks where there are only tens of event types, we manually specify all the required information. For SRL frames, we semiautomatically collect them with data statistics, as described in §2.2 and the next section.

A.2 Extra Information for SRL Frames

We adopt SRL frames from PropBank 3.1,²⁰ Nom-Bank 1.0²¹ and FrameNet 1.7²². Since many Nom-Bank frames are derived from PropBank frames, we simply map²³ them to the PropBank counterparts and ignore the ones that do not have such mappings. We filter event-related SRL frames by excluding the ones that do not have any verb realizations, which are judged by the POS sets provided in the frame files. Moreover, we only consider a subset²⁴ of non-core or modifier roles that are related to the target EAE task.

To apply the role querying strategies, we semiautomatically collect extra information from the SRL data. Specifically, for each role in a frame, we collect the following information.

propbank-frames/releases/tag/v3.1
²¹https://nlp.cs.nyu.edu/meyers/
nombank/nombank.1.0/

¹⁸We utilize a specific [UNK] symbol to denote the case when there are no arguments in the context.

¹⁹The "Where" question is designated to the role of "place".
²⁰https://github.com/propbank/

²²https://framenet.icsi.berkeley.edu/ fndrupal/frameIndex

²³We check the "source" attribute in a NomBank frame, which points to a PropBank frame if existing.

²⁴These include "ARGM-LOC" in PropBank and {"Place", "Instrument", "Weapon", "Vehicle"} in FrameNet.

| Scheme | Frame | Template |
|----------|---|---|
| PropBank | buy.01 forbid.01 rent.01 swim.01 | buyer buy goods for recipient from seller for money in place authority forbid protagonist action in place at lessor lessee rent goods from lessor for money in place from area self mover swim against area to goal in place |
| FrameNet | Abandonment Commerce_buy Employing Mention | agent abandon theme in place buyer buy goods in place employer employ employee field position task in place communicator mention specified content message in medium in place |

Table 6: Examples of the auto-generated templates. The predicate is **emboldened** and the roles are underlined. Note that some of the examples are especially picked to show typical problems of this semi-automatic process.

• Role names. We directly take the role label names of FrameNet since they are already in natural language forms. We further train a role label classifier²⁵ with the FrameNet data and apply it to the PropBank data. Then for each framespecific role, the most frequently predicted label will be the role name. For example, for the "Arg0" role of the "buy.01" frame, its arguments in the dataset are mostly predicted to the "buyer" label, which is thus assigned as its role name.

1214

1215

1216

1217

1218 1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234 1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

- Role orders in the templates. We construct a template for an SRL frame by concatenate its predicate word²⁶ and role names. The main thing to specify is their ordering. We again take a statistical approach and collect each role's relative distance to the predicate. For example, in the "buy.01" frame instance of "He bought a book in a store.", "Arg0" (He) gets a distance of -1, "Arg1" (book) gets a +1 and "ArgM-LOC" (store) gets a + 2. Finally, the role orders in the templates are decided by the roles' average relative distances. We aim to obtain a canonical verb-styled ordering in active voice, and thus we only consider frame instances that are realized by non-passive verbal predicates.
 - Preposition words in the templates. When realized in natural language sentences, many roles are accompanied by prepositions. We count the frequency that a role is filled by an argument that utilizes a preposition²⁷ and keep the prepositions that appear more frequent than 25%. When there

are such prepositions, we add the preposition before the role name and put them together into the slot. When there are multiple feasible prepositions, we randomly sample one in training and utilize the most frequent one in testing.

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1269

1270

1275

1280

1282

• Question words. For the MRC strategy, we also need to identify the role types to select question words. Since the "Where" question is specific to places ("ArgM-LOC" in PropBank and "place" in FrameNet), we only need to distinguish whether the role is person-related. We again take a counting-based method by checking how many times a role is filled by a personal pronoun. If this happens for a role with a frequency larger than 10%, we regard it as potentially person-related. Since there are cases where a role can be filled by either a person or an object, at training time we randomly pick "Who" or "What" questions for these potentially person-related roles, while only asking "Who" in testing.

Most of the above heuristics are decided by manually checking the generated outputs for the Prop-1266 Bank and FrameNet frames. We provide some 1267 examples of generated templates in Table 6. No-1268 tice that this semi-automatic approach is far from perfect and there can be noises and inconsistencies. For example, in the templates of "rent.01" 1271 and "swim.01", the role ordering is slightly strange 1272 and there are repeated role names, and in the "Em-1273 ploying" and "Mention" frames, prepositions are 1274 missing for some roles. Nevertheless, we find that most of the generated templates are close to natural 1276 language and they generally look reasonable for 1277 the aim of our usage. There is one more caveat, 1278 that is, we only allow one template for one frame, 1279 while there can be multiple templates with syntactic variations. We leave the investigations of better 1281 template constructing methods to future work.

²⁵This classifier is similar to our CLF querying model except that no extraction is needed. Its accuracy on the FrameNet dev set is around 0.7. Notice that even when the classifier does not hit the most suitable label, the predicted ones may still be reasonable for our usage.

²⁶For PropBank frames we simply take the predicate's lemma, while for FrameNet frames, we choose the lemma of the lexicon unit that overlaps mostly with the frame name.

²⁷The criterion is that the argument's head word has a dependency relation of "case" to a child whose POS is "ADP".

| Dataset | Split | Sent. | Event | Arg. | A/E |
|---------|-------|--------|-------|------|-----|
| ACE | Train | 192.2K | 4.4K | 6.6K | 1.5 |
| | Dev | 0.9K | 0.5K | 0.8K | 1.6 |
| | Test | 0.7K | 0.4K | 0.7K | 1.6 |
| ERE | Train | 147.3K | 6.2K | 8.9K | 1.4 |
| | Dev | 1.2K | 0.5K | 0.7K | 1.4 |
| | Test | 1.2K | 0.6K | 0.8K | 1.5 |

Table 7: Statistics of the ACE and ERE data. "A/E" denotes the averaged argument number per event.

| Туре | Dataset | Sent. | Inst. | Arg. | A/I |
|------|----------|--------|--------|--------|-----|
| SRL | PropBank | 77.7K | 256.1K | 374.5K | 1.5 |
| | NomBank | 28.2K | 56.9K | 86.8K | 1.5 |
| | FrameNet | 173.0K | 173.4K | 208.5K | 1.5 |
| QA | SQuAD | 62.2K | 130.3K | 86.8K | 0.7 |
| | QA-SRL | 64.0K | 299.3K | 299.3K | 1.0 |
| | QANom | 7.1K | 26.4K | 26.4K | 1.0 |
| | QAMR | 4.8K | 88.3K | 88.3K | 1.0 |

Table 8: Statistics of the SRL and QA data. "Inst." denotes the number of SRL or QA instances, where "A/I" denotes the averaged number of argument or answer per instance.

B **Details of Main Experiments**

B.1 Settings

1283

1284

1285

1286

1287

1288

1289

1290

1291

1294

1295

1296

1297

The main experiments are conducted with English ACE²⁸ (ACE05- E^+) and ERE²⁹ (ERE-EN) datasets. We adopt the preprocessing scripts³⁰ from ONEIE (v0.4.8) Lin et al. (2020). The statistics of the event data are shown in Table 7. For SRL data, we take those from the latest PropBank³¹ (EWT and OntoNotes), NomBank³² and FrameNet³³. For FrameNet, we utilize the lexicographic annotation sets since there are much more instances. As described in §A.2, we ignore SRL frames that do not have verbal predicates and only keep related non-core roles. For QA data, we include SQuAD³⁴, QA-SRL³⁵, QANom³⁶ and QAMR³⁷. For the QA

```
<sup>28</sup>https://catalog.ldc.upenn.edu/
LDC2006T06
  <sup>29</sup>LDC2015E29, LDC2015E68, and LDC2015E78.
  <sup>30</sup>http://blender.cs.illinois.edu/
software/oneie/
  <sup>31</sup>https://github.com/propbank/
propbank-release
  <sup>32</sup>https://nlp.cs.nyu.edu/meyers/NomBank.
html
  <sup>33</sup>https://framenet.icsi.berkeley.edu/
fndrupal/
  <sup>34</sup>https://rajpurkar.github.io/
SQuAD-explorer/
  <sup>35</sup>https://github.com/uwnlp/qasrl-bank
  <sup>36</sup>https://github.com/kleinay/QANom
  <sup>37</sup>https://github.com/uwnlp/qamr
```

instances, we follow Michael et al. (2018) and use a question-context alignment heuristic to find a predicate in the context for each question. Since the external data is mainly utilized as training resources, we simply concatenate the all the available data portions for training while splitting a small subset for development. Data statistics of SRL and QA are shown in Table 8.

1298

1299

1300

1301

1303

1304

1305

1307

1308

1314

1319

1320

1324

1325

1327

1329

1330

1331

1333

1334

1335

1336

1337

1338

1339

1341

1342

1343

1344

1345

1346

We utilize pre-trained language models 1306 (RoBERTabase for encoder-only models (CLF, MRC, TSQ) and BART_{base} for encoder-decoder models (GEN)) to initialize our models and fine-tune the full models in all the experiments. 1310 The model parameter numbers are 125M and 1311 139M, for those with RoBERTa and BART 1312 respectively. For the hyper-parameter settings, we 1313 mostly follow common practices. Adam is utilized for optimization. The learning rate is initially set 1315 to 2e-5 and linearly decayed to 2e-6 throughout the 1316 training process. The models are trained for 50K 1317 steps with a batch size of 16 for event and SRL and 1318 32 for QA. We pick models by the performance on the development set of each task. In low-resources cases, the original event development set is also down-sampled accordingly as the training set. All 1322 the experiments can be conducted with one 1080 1323 Ti GPU and the training can usually be finished within several hours.

B.2 Details of GPT-3 Prompting

To perform prompting with GPT-3, we utilize the OpenAI API.³⁸ We adopt the "Davinci" model and the Completion endpoint. We design the prompts with a strategy that is similar to MRC. The prompts consist of three parts: the context sentence, a question for the querying role and a partial answer to be completed. The context is simply the sentence where the event trigger appear, while the questions are those shown in Table 5 as in the MRC strategy. The to-be-completed answer sentence is a declarative repetition of the question. For example, we have the following prompt to query the "artifact" role with the context of "He went to the store and bought a book.": He went to the store and bought a

book. Q: What is the artifact of the bought event?

A: The artifact of the bought event is

³⁸https://openai.com/api/

| Mathad | ACE | | ERE | |
|--------------------|-------|-------|-------|-------|
| Method | gold | pred. | gold | pred. |
| Super. | 68.93 | 53.98 | 72.24 | 49.77 |
| QA | 38.62 | 29.94 | 38.74 | 26.02 |
| SRL _{CLF} | 33.18 | 26.65 | 33.89 | 25.78 |
| SRL _{MRC} | 47.08 | 25.50 | 52.32 | 38.39 |
| SRL _{GEN} | 49.99 | 37.01 | 55.38 | 38.77 |
| SRL _{TSQ} | 53.36 | 39.41 | 57.81 | 40.08 |
| +shuf. | 54.82 | 41.15 | 58.54 | 40.78 |
| +distill | 55.17 | 41.50 | 57.95 | 40.01 |
| +both | 56.35 | 42.16 | 58.96 | 41.15 |

Table 9: Zero-shot EAE results (F1%) on event test sets with gold or predicted event triggers.

| Model | ACE | ERE |
|-----------------|-------|-------|
| encoder-only | 56.35 | 58.96 |
| encoder-decoder | 55.36 | 58.44 |

Table 10: Comparisons between encoder-only and encoder-decoder TSQ models.

We let the GPT-3 model to greedily decode the remaining of the answer sentence and match the results to the tokens in the original context to obtain the arguments. When there are no matchings or the answer is "not specified", no arguments are predicted for the querying role. Since the answer should come from the context, we utilize the "logit_bias" parameter to constrain the model to adopt sub-tokens that appear in the context (or those from "not specified").

B.3 Results with Predicted Triggers

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1360

1362

1363

1364

1365

1366

1368

1369

1370

1372

1373

In our main experiments, we assume given gold event triggers. In this sub-section, we train a supervised sequence-labeling event detector and further utilize the predicted triggers to perform zero-shot argument extraction. The results are shown in Table 9. The event detectors could obtain labeled F1 score of 71.0 and 58.4 for ACE and ERE, respectively. With the predicted triggers, the results drop correspondingly against those with gold triggers. Nevertheless, the overall trends are similar. The TSQ strategy performs the best while the argument augmentation is also helpful with predicted triggers. One interesting direction to explore is full event extraction in the zero-shot and low-resource scenarios, which we leave to future work.

B.4 Model Choice for TSQ

1374 Concurrently, Ma et al. (2022) explore an idea 1375 that is similar to TSQ, while taking sequence-to-

| Method | Gold | Predicted |
|----------------------------|-------|-----------|
| OneIE (Lin et al., 2020) | - | 54.8 |
| EEQA (Du and Cardie, 2020) | 63.34 | - |
| GenIE (Li et al., 2021b) | 66.67 | 53.71 |
| CLF | 66.96 | 52.62 |
| MRC | 66.55 | 52.43 |
| GEN | 66.76 | 52.81 |
| TSQ | 68.93 | 53.98 |

Table 11: Comparisons of fully-supervised ACE05- E^+ test results (F1%) (with gold or predicted triggers).



Figure 5: Argument F1(%) scores on ACE and ERE test sets with different amounts of training data. Here *x*-axis (drawn in log scale) denotes the percentage of original training data sampled. The shaded areas indicate the ranges of standard deviations.

sequence encoder-decoder model to perform argument extraction. Specifically, they encode the contexts with the encoder while put the template at the decoder side. We also compare this encoderdecoder scheme with our encoder-only TSQ in the transfer scenario. The results are shown in Table 10, where the encoder-only model is slightly better. Therefore, we utilize the encoder-only model for the TSQ strategy. 1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

B.5 Supervised Results

Although our main focus is on the transfer scenar-
ios, we also conduct purely supervised experiments1386on the target EAE datasets. We first compare fully-
supervised results with previous works. As shown1389

| Category | Example | SRL | Super. |
|-------------|--|--------------|--------------|
| Correct | | 194 (50.52%) | 238 (63.47%) |
| Role | Actually, they paid _{TransferMoney} for [it] _{Beneficiary} . | 45 (11.72%) | 22 (5.87%) |
| Local | <u>My</u> _{Buyer} plan is to pay _{TransferOwnership} off my car. | 40 (10.42%) | 23 (6.13%) |
| Head | They fired _{Attack} mortars in the [direction] _{Target} of the <u>7th Cavalry</u> _{Target} . | 24 (6.25%) | 8 (2.13%) |
| Global | "We condemned the attack _{Attack} ," he said, adding that his messages to the <u>terrorists_{Attacker}</u> is: Their efforts will not be successful. | 24 (6.25%) | 17 (4.53%) |
| Others | | 14 (3.65%) | 10 (2.67%) |
| Ambiguous | At least four [policeman] _{Attacker Victim} were injured in clashes _{Attack} . | 18 (4.69%) | 20 (5.33%) |
| Span | The <u>1st [Brigade]</u> _{Attacker Attacker} took Karbala with a minimal fight _{Attack} . | 12 (3.12%) | 14 (3.73%) |
| Coreference | HeDefendant skipped bail during [his]Defendant trialHearing. | 13 (3.39%) | 23 (6.13%) |

Table 12: Examples of the categories and results of the manual error analysis. In the examples, the triggers are shown in **bold** texts with $_{brown}$ event types. The gold arguments are presented in <u>underlined spans</u> with $_{red}$ roles, while predicted ones are [bracketed] followed by $_{blue}$ roles. Results are denoted with number counts and (percentages). The rows of the error categories are sorted by the gap between SRL and supervised counts.

in Table 11, our results are generally comparable to those in previous works, which validates the quality of our implementation.

Furthermore, we compare the four querying strategies with different amounts of training data. The results are shown in Figure 5. The overall trend is similar in both datasets. In high-resource scenarios, different querying strategies could obtain similar results. In low-resource cases, generally the methods that capture more contextual information in the queries can perform better. The CLF strategy with non-contextualized queries obtain obviously worse results than the others, while TSQ is the overall best-performing strategy. This is also consistent with the results in the zero-shot transfer scenarios.

B.6 Manual Analysis

1390

1391

1392

1393

1394

1395

1396

1397

1398

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

We further perform a manual error analysis to investigate what the main error types are. We randomly take 100 event frames that contain prediction errors from the ACE development set and categorize the errors. We perform this analysis for both our best zero-shot SRL model and the supervised model to examine where the main gaps are. We specify eight error categories:

- **Ambiguous** cases, where there are annotation errors or ambiguities and the predictions could be regarded as correct in some way.
 - **Coreference**, where predicted and gold arguments are co-referenced in some way.
- **Span** mismatch, where the main contents are captured with non-crucial boundary mismatches.

• **Head** mismatch, where the main contents are roughly captured but not with the exact annotated words. This happens mostly in appositions or noun modifiers with more specific contents.

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

- **Role** misunderstanding, where the semantic meaning of a role is not correctly understood.
- Local inference, where correct predictions require semantic inference at the local clause.
- **Global** understanding, where correct predictions require global understanding of the full context.
- **Others**, where the error does not fall into any of the above categories.

Examples of these categories and the results are 1433 shown in Table 12. According to the statistics, 1434 the main gaps between the SRL and supervised 1435 models are on the categories of role misunderstand-1436 ing, lacking of semantic inference as well as head 1437 mismatch. Head mismatches are due to the discrep-1438 ancies between syntactic head and semantic core 1439 words, and might not cause severe problems. The 1440 first two are more semantic errors that are related 1441 to the essence of the EAE task. Role misunder-1442 standing may be related with template mismatches, 1443 where roles in the SRL templates are different than 1444 those in target event ones. Lacking of semantic in-1445 ference is mostly upon distant arguments. Though 1446 the argument augmentation techniques recover cer-1447 tain distant arguments for SRL frames, this prob-1448 lem is still far from being solved. Notice that these 1449 semantic errors reveal the main difficulties of the 1450 EAE task, which even supervised systems have not 1451 yet fully tackled. To solve these problems, more 1452 comprehensive semantic understanding is required. 1453

| Event | Template | Role questions |
|--|---|--|
| Expression Transcription Catabolism Phosphorylation Localization | agent express theme agent transcribe theme agent degrade theme agent phosphorylate theme agent localize theme | What is expressed? What is transcribed? What is degraded? What is phosphorylated? What is localized? |
| Binding | agent bind theme1 to theme2 | What is bound? What is something bound to? |
| Regulation | cause regulate theme | What causes the regulation? What is regulated? |

Table 13: Manually specified templates and role questions for GENIA events ("agent" is a dummy role introduced to form the templates in active voice).

| Role | QA | SRL | +shuf. | +distill | +both |
|-------------|-------|-------|--------|----------|-------|
| Place | 51.73 | 47.10 | 46.94 | 56.46 | 57.38 |
| Attacker | 34.59 | 56.52 | 59.03 | 57.88 | 58.19 |
| Entity | 38.02 | 40.96 | 43.97 | 41.59 | 43.36 |
| Target | 33.15 | 38.62 | 38.66 | 37.59 | 39.04 |
| Victim | 66.98 | 80.58 | 81.18 | 79.27 | 79.49 |
| Artifact | 13.95 | 49.52 | 62.82 | 46.47 | 60.54 |
| Person | 58.53 | 71.82 | 72.16 | 73.46 | 73.20 |
| Recipient | 45.75 | 46.68 | 47.51 | 50.64 | 49.37 |
| Destination | 66.11 | 65.97 | 66.14 | 68.02 | 66.00 |
| Instrument | 34.28 | 40.00 | 47.05 | 43.54 | 49.56 |

Table 14: F1% score breakdowns by argument roles.

| Method | Single-instance | Batched |
|--------|-----------------|---------|
| CLF | 184 | 316 |
| MRC | 106 | 146 |
| GEN | 28 | 144 |
| TSQ | 167 | 281 |

Table 15: Decoding speed (instances per second) comparisons of different role querying strategies. We evaluate both single-instance and batched decoding modes.

Role Breakdowns B.7

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

We perform breakdowns on argument roles on ACE with the zero-shot models. Table 14 shows the results of the top-ten frequent roles. Interestingly, distillation generally helps more on the non-core roles, such as PLACE and DESTINATION, while shuffling enhances core roles, like ATTACKER and VICTIM. Finally, applying both could lead to the overall best results.

B.8 Speed Comparisons

We also perform decoding speed comparisons to ex-1464 amine the efficiency of different querying strategies. 1465 The results are shown in Table 15. There are no sur-1466 prises that the simplest CLF strategy achieves the 1467 highest decoding speed, since its input sequences 1468 are the shortest and there is no further complex 1469

query encoding. TSQ is only around 10% slower, 1470 but still efficient compared with other two meth-1471 ods, where MRC suffers from multiple forwarding for different role queries and GEN requires autoregressive decoding at testing time.

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

С **Further Extensions**

C.1 GENIA Details

For the GENIA experiments, one more assumed input is the protein entities, following the settings of BioNLP shared task. Since our model-predicted argument head words might not match the protein entities, we perform a syntax-based post-processing heuristic. For a predicted argument word, we check its descendants in the dependency tree and relocate the argument to the highest node that belongs to an entity (or an event for the "Theme" of "Regulation"). If no such items can be found, the prediction is ignored. The evaluation metric is approximate recursive matching using the official online service.³⁹

For the GENIA events, we manually specify templates, which are shown in Table 13. We also manually specify role questions since the templated questions mostly fail in this scenario. For the three regulation events, we simply adopt the same specifications since no obvious differences are found when adding modifiers of "positively" or "negatively". Since our SRL templates are all formed in active voice, we introduce a dummy "agent" role to form non-passive GENIA templates. The prediction of this dummy role is ignored in testing.

C.2 Multi-lingual Analysis

One interesting aspect in the multi-lingual scenario is how the predictions are influenced by the word

³⁹http://bionlp-st.dbcls.jp/GE/2011/ eval-test/

| Language | Model | Pearson | Spearman |
|----------|---------|---------|----------|
| Arabic | w/o SRL | 0.6050 | 0.6727 |
| | w/ SRL | 0.5157 | 0.1394 |
| Chinese | w/o SRL | 0.6910 | 0.5636 |
| | w/ SRL | 0.5025 | 0.2727 |

Table 16: Correlations between relative role order differences and performance gaps to supervised systems for multi-lingual EAE (with top-10 frequent roles).

1503 order difference between the source and target languages. We analyze the influence by measuring 1504 the performance differences in different roles. We 1505 firstly calculate the directional statistics for each 1506 role in each language, specifically: for a role in 1508 a language, what percentage of its arguments appear after the trigger? For example, "Attacker" ap-1509 pears after the trigger in 26.9% of the times in 1510 English, while this percentage is 72.7% in Arabic. 1511 Then for each role, we have a source-target order 1512 difference metric, which is the absolute value of the frequency difference. We further calculate the 1514 performance differences between a transfer model 1515 1516 trained with English data and a supervised model directly trained on the target language. Finally, we measure the correlation between the order differ-1518 ences and performance differences for the top-ten 1519 frequent roles in each language. The results for the 1520 1521 transfer model with or without (multi-lingual) SRL intermediate-training are shown in Table 16. Inter-1522 estingly, if directly transferring from English to the 1523 target languages, there are at least moderate corre-1524 lations between the order differences and perfor-1525 mance gaps. While using SRL, the correlations de-1526 crease probably because of the extra signals about 1527 target language order in the SRL data. This shows 1528 that order difference may be a major factor influ-1529 encing the effectiveness of cross-lingual transfer. 1530 Currently our templates are all in English-styled 1531 and it would be an interesting future direction to 1532 explore the influences of template specifications 1533 such as role orders.

C.3 More Multi-sentence Results

1535

1536We also perform intermediate-training on the
RAMS dataset with different amounts of target
training instances. The test results are shown in
Figure 6. The same-sentence patterns are similar
to those in previous ACE experiments, while SRL
seems to be able to provide small but consistent
benefits for cross-sentence arguments.



Figure 6: RAMS results (same- or cross-sentence argument F1%) with or without SRL intermediate-training.

1543

D Discussion of Limitations

This work has several limitations. Firstly, we only 1544 focus on the event argument extraction step and 1545 assume given event triggers. Though the first step 1546 of event detection is also important for event extrac-1547 tion, we do not cover it in this work mainly due to 1548 two reasons: 1) the annotation of event triggers is 1549 generally less laborious than argument annotation 1550 since word-level tagging instead of pairwise linking 1551 is required; 2) Event detection is highly specific to 1552 the target scheme, which is different than argument 1553 extraction where there are more sharings between 1554 semantic roles and event arguments. Secondly, in 1555 this work, SRL templates are created heuristically 1556 and do not cover syntactic and language variations. 1557 For example, we only construct English-styled tem-1558 plates in active voice, which might not be ideal 1559 for all the cases. We mainly aim to show that the 1560 template-based method is a promising way to per-1561 form argument extraction especially in transferring 1562 scenarios, but surely there could be better ways to 1563 construct the querying templates. Finally, though 1564 the applying of argument augmentation recovers 1565 certain amounts of distance arguments, it is still far 1566 from an ideal solution to the problem. This calls 1567 for more future investigations on this direction, re-1568 searching towards deeper and more comprehensive 1569 semantic understanding of natural languages. 1570