## Instruction-Tuning LLMs for Event Extraction with Annotation Guidelines

**Anonymous ACL submission** 

#### Abstract

In this work, we study the effect of annotation guidelines-textual descriptions of event types and arguments, when instruction-tuning large language models for event extraction. We 005 conducted a series of experiments with both human-provided and machine-generated guidelines in both full- and low-data settings. Our results demonstrate the promise of annotation guidelines when there is a decent amount of training data and highlight its effectiveness in improving cross-schema generalization and low-frequency event-type performance.

#### 1 Introduction

007

011

013

017

025

031

036

Event Extraction (EE) aims to identify and structure what, who, when, where, and how of realworld events from given textual resources (Doddington et al., 2004; Ji and Grishman, 2008; Li et al., 2022; Xu et al., 2024). Translating this abstraction requires complex schema specifications that define event types, argument roles, and their interrelationships, yet being able to precisely capture the language nuances and distinguish between event types and argument roles, which posits the task as an inherently challenging problem.

Recently, large language models (LLMs) have transformed NLP research and practices dramatically, owing to the rich knowledge and other capabilities (e.g., reasoning) they have obtained from extensive pre-training (Wei et al., 2022; Chen et al., 2023; Shi and Lipani, 2023). This transformation has similarly impacted the broader research field of Information Extraction (IE). Existing applications of LLMs to IE can be categorized into two lines. The prompt engineering-based approaches, often based on proprietary LLMs, consider an LLM as a black box, querying it with task specifications via zero- or few-shot prompting and relying on its latent knowledge to extract interested information (Gao et al., 2023; Wang et al., 2023b; Li et al.,

2023a; Srivastava et al., 2023). However, these approaches not only lead to inferior performance but also incur prohibitive costs, especially when the task is complex.

040

041

042

045

046

047

048

051

052

054

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

079

Our work will thus focus on the second line of approach, namely, instruction-tuning open-weight LLMs. This line of approach adapts an LLM to specific IE tasks and schemas by directly training it to follow the task instructions, which offers a promising yet cost-effective solution. For example, Wang et al. (2023a) leverages natural language instructions to guide large language models for IE tasks; Li et al. (2024) proposed a two-phase learning framework that enhances schema understanding and following ability via automatically annotated data. More recently, Sainz et al. (2024) instruction-tuned LLaMA (Touvron et al., 2023) on multiple IE datasets and discovered annotation guidelines-textual descriptions of an event type and its argument roles used by human annotators when collecting the dataset, as effective components of an IE task's instruction. Despite the promise of the existing explorations, however, most of them have focused on the relatively simpler task of Named Entity Recognition, yet how to properly instruction-tune LLMs for the structured EE task is still understudied.

To fill this gap, we study instruction-tuning LLMs for EE, with a focus on the role of annotation guidelines in task instructions (Fig. 1). We conduct a systematic analysis using LLaMA-3.1-8B-Instruct on two EE datasets (ACE05 (Doddington et al., 2004) and RichERE (Song et al., 2015)) under varied training settings. Our key findings are organized around four themes:

1) Effect of Annotation Guidelines on Event Ex*traction* — We found that annotation guidelines improve performance by helping the model distinguish fine-grained event types. However, this advantage may diminish when negative sampling is introduced during training, which allows the



Figure 1: Overview of our exploration of automatically generating annotation guidelines to augment code-format instruction tuning for EE. Prompt template for Guideline-PN and the example outputs are shown.

model to learn event distinctions from additional contrastive examples instead.

2) Comparing Machine-Generated and Human-Written Guidelines — Prior work assumed access to human-authored guidelines, which may not hold in practice. We thus proposed 5 different ways to automatically generate annotation guidelines. We find that they outperform human-written ones by up to 11% and 7% in trigger and argument classifications, respectively.

3) Guidelines in Data-Scarce Scenarios — Our results show that with only 2000 training samples, guidelines allow LLMs to reach performance levels comparable to full-data training. However, when data is extremely scarce (100 samples), models tend to rely more on memorization than on guideline-driven schema constraints.

4) *Cross-Schema Generalization* — We assess whether structured guidelines help models generalize to different EE schemas. While models trained on RichERE transfer well to ACE (suggesting fine-to-coarse schema adaptation is feasible), the reverse scenario sees a performance drop due to RichERE's more complex event structures and expanded argument roles.

Finally, we confirmed a consistent effect of annotation guidelines on the smaller LLaMA-3.2-1B- Instruct model and conducted an in-depth analysis showing that the guidelines can help LLMs reduce common types of EE errors and are beneficial to event types with any frequency in the training set. 108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

#### 2 Approach

#### 2.1 Task Formulation

Given an input sentence X, the goal of EE is to extract the structured event information Y from the sentence, adhering to predefined schema constraints  $\mathcal{E}$ . The extraction task consists of (1) **Trigger Extraction**, which localizes an event trigger span and classifies its event type, and (2) **Argument Extraction**, where the task is to identify spans in X that serve as argument roles within the extracted event instance.

When an autoregressive LLM is tasked with EE, the extraction of event instances is formulated in a generative way, with the LLM generating a sentence describing the extracted event instances. Specifically, the prompt to the LLM is defined as  $P = [I \oplus E_e \oplus X]$ , where  $\oplus$  is the concatenation operation, I represents the task instruction, which specifies the structured output format and task definition, and  $E_e \in \mathcal{E}$  denotes the event schema of an interested type e from a predefined set  $\mathcal{E}$ .

081

Let  $\mathcal{D} = \{(e_i, X_i, Y_i)\}_{i=1}^N$  denote a dataset of annotated event examples, where each  $X_i$  corresponds to a prompt instance  $P_i$  for the interested event type  $e_i$ . The objective function of instruction tuning for EE is as follows:  $\mathcal{L}(\mathcal{D}; \theta) =$  $-\sum_{i}\sum_{j}\log p_{\theta}(Y_{ij} \mid P_i, Y_{i, < j})$ , where  $Y_{i, < j}$  represents previously generated tokens in the structured output sequence, ensuring an autoregressive formulation.

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

159

160

161

162

163

164

168

169

170

171

172

174

176

177

179

181

183

Existing work identified the structure of EE outputs to be critical (Jiao et al., 2023; Wang et al., 2023b). In particular, Wang et al. (2023b) found that formulating the EE output in a *code format* can take advantage of Programming Language features such as inheritance and type annotation to introduce external knowledge or add constraints. In our work, we follow the same formatting strategy and represent the EE task as a code generation problem. Specifically, the event schema  $E_e$  is represented as a Python class; accordingly, every extracted event instance is represented as a Python object of the corresponding event class. When there are multiple event instances implied in the input X, a list of Python objects will be generated; when there is no event specified in X, we expect an empty Python list to be the model output. An example is shown in Figure 1.

During training, we provide only the groundtruth event schema in the prompt; when the text input X does not include any event, a random event schema will be chosen. At inference time, given a test instance X, we pair the input with every possible event type in the schema set  $\mathcal{E}$ , prompt the LLM to extract any implied event instances, and perform model evaluation based on the aggregated extraction outputs. As such, a well-performing LLM needs both extract the complete event instances and avoid events that are not indicated in X.

#### 2.2 **Instruction-Tuning LLMs with** Annotation Guidelines

Recent work by Sainz et al. (2024) demon-173 strated the effectiveness of integrating annotation guidelines in the code-format instructions of IE 175 tasks. Specifically, when describing the event type schema  $E_e$ , a textual description is added to the event type and each of its argument roles (Figure 1). As such, the LLM is expected to more easily understand the meaning of the event type while be-180 ing instructed to extract any occurring events from the input X. While Sainz et al. (2024) evaluated annotation guidelines in the broad IE task, their

main focus has been on Named Entity Recognition, rather than the complicated EE task. Furthermore, their approach assumed the availability of pre-existing human-curated guidelines, an assumption that may not always hold in real-world applications. To bridge this gap, we explore methods to automatically generate annotation guidelines and assess their effectiveness in comparison to humanauthored ones.

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

To develop a scalable and cost-effective approach for guideline generation, we employ a reverse engineering strategy, leveraging both annotated event examples and the strong generative capabilities of LLMs. As illustrated in Figure X, we construct a guideline generation prompt for each event type e by providing a few annotated examples  $\{(X_i, Y_i)\}$  demonstrating the existence or non-existence of event instance of type e, and then prompt an LLM (GPT-40 in our experiment) to generate annotation guidelines for e. In total, we experimented with five variants of machinegenerated guidelines: (1) Guideline-P: We prompt the LLM with 10 positive annotated examples of type e to generate the annotation guidelines. Inspired by Sainz et al. (2024), we sample 5 distinct guidelines for each event type, which can be used during the model training to ensure that the model is exposed to multiple rephrasings of the guidelines rather than memorizing and overfitting to a specific one. (2) Guideline-PN (Positive + Negative Examples): In addition to 10 positive event annotations, we also provide 15 negative annotations where the input X does not imply event instances of type e. Similarly, we prompt the LLM to generate 5 distinct guidelines for each event type. (3) Guideline-PS (Positive + Sibling Events): Similar to Guideline-PN, we prompt the LLM with both positive and negative event annotations. However, the negative annotations are selected from the sibling event types of the target type e (e.g., Arrest vs. Jail), as defined by the event ontology. We hypothesize that the critical challenge for EE lies in distinguishing between sibling event types; hence, an instructed LLM can benefit from following annotation guidelines that particularly emphasize the difference between sibling event types. As in the earlier variants, we generate 5 guideline samples per event type. (4) Guideline-PN-Int and (5) Guideline-PS-Int: Finally, we create two more variants that Integrate the 5 diverse guideline samples from Guideline-PN and Guideline-PS into a comprehensive one, respectively. Examples of the

	Examples of Annotation Guidelines for Event Type: Extradite (ACE05)
GUIDELINE-H Avg. Length - 107.67 tokens	<ul> <li>Event Type: An EXTRADITE Event occurs whenever a PERSON is sent by a state actor from one PLACE to another place for the purposes of legal proceedings there.</li> <li>Arguments:</li> <li>AGENT: The extraditing agent.</li> <li>PERSON: The person being extradited.</li> </ul>
GUIDELINE-P Avg. Length - 163.87 tokens	<ul> <li>Event Type: The Extradition event refers to the formal process where one jurisdiction delivers a person accused () The event can be triggered by terms such as 'extradition' () Edge cases include situations where the term 'extradition' is used metaphorically or in a non-legal context.</li> <li>Arguments: <ul> <li>AGENT: () the agent is the organization or authority (). Examples include 'court', 'government', ()</li> <li>PERSON: () individual who is being transferred to another jurisdiction. Examples are 'she', ()</li> </ul> </li> </ul>
GUIDELINE-PN Avg. Length - 285.24 tokens	<ul> <li>Event Type: The event is triggered by the formal request () for legal reasons. Triggers such as 'extradition' are indicative of this event type, not 'Transport' which involves general movement without legal context.</li> <li>Arguments: <ul> <li>AGENT: The agent is responsible for the legal and procedural aspects of the extradition, (). An example is 'the original court' ()</li> <li>PERSON: () one who is being moved from one place to another under legal authority. For example, 'he' ()</li> </ul> </li> </ul>
GUIDELINE-PS Avg. Length - 159.79 tokens	<ul> <li>Event Type: () person being moved to a new jurisdiction (). This differs from events like 'TrialHearing' or 'Convict', which focus on the legal proceedings and outcomes within a single jurisdiction.</li> <li>Arguments: <ul> <li>AGENT: () Edge cases may include international organizations or coalitions () such as the U.N. ()</li> <li>PERSON: Unlike the 'defendant' in events like 'TrialHearing' or 'Convict', the person in the 'Extradite' event is specifically being transferred for legal proceedings or punishment.</li> </ul> </li> </ul>
GUIDELINE-PN-INT Avg. Length - 439.94 tokens	<ul> <li>() Key triggers include terms like 'extradite', 'extradition', and 'extraditing'. It is distinct from events like 'ArrestJail' and 'ReleaseParole', as it specifically involves ()</li> <li>Arguments: <ul> <li>AGENT: The agent () typically a legal or governmental body. Examples include 'court', 'government'()</li> <li>PERSON: The person is the individual being extradited, the subject of the legal transfer. Examples include 'she', 'him', and 'her'.</li> </ul> </li> </ul>
GUIDELINE-PS-INT Avg. Length - 434.64 tokens	The 'Extradite' event involves the legal transfer of a person (). It is distinct from events like 'ArrestJail', (), and 'ReleaseParole' or 'Pardon', () Arguments: - AGENT: The agent is the entity () such as a court, government, or police department. This entity ensures the transfer is conducted according to legal protocols () - PERSON: () They are the central figure in the extradition process, distinct from a 'defendant' in other legal events, () This may include high-profile individuals or groups.

Table 1: Examples of annotation guidelines for the event type Extradite from ACE05. Due to space limits, only agent and person were shown for arguments, and only 1 out of the 5 guideline samples were shown for **P**, **PN**, and **PS**. We highlight distinctions from other event types, example mentions, and edge cases in guidelines.

5 guideline variants are shown in Table 1. The prompt templates used for generating guidelines are provided in Appendix A.3.

#### **3** Experiments

239

241

242

243

245

247

248

251

#### 3.1 Experimental Setup

**Datasets.** We perform experiments on two standard EE datasets: **ACE05** (Doddington et al., 2004) and **RichERE** (Song et al., 2015). Both of them exhibit fine-grained event distinctions, and RichERE includes sparser event annotations (i.e., fewer event-labeled sentences), which makes it more challenging. Moreover, RichERE does not come with human-written annotation guidelines. Datasets were split following the TextEE benchmark (Huang et al., 2024) and then converted to code format automatically by our scripts. **Evaluation.** Following prior work (Huang et al., 2024), we evaluate the model on four F1 metrics: (1) **Trigger Identification** (**TI**), which measures correct trigger span extraction, (2) **Trigger Classification** (**TC**), which additionally requires event-type correctness, (3) **Argument Identification** (**AI**), which ensures correct argument role association with the predicted trigger, and (4) **Argument Classification** (**AC**), which further requires role-type correctness and is thus the most comprehensive metric on a model's EE performance. When evaluating the model on the Guideline-P, PN, and PS variants, one guideline is randomly selected each time.

As a side benefit of representing events in a structured code format, we can easily evaluate an extracted event instance by directly instantiating its corresponding Python object based on the event schema's Python class definitions, checking if the 252

Experiments	ACE w/o NS				ACE w/ NS				RichERE w/o NS				RichERE w/ NS			
	TI	ТС	AI	AC	TI	ТС	AI	AC	TI	ТС	AI	AC	ТІ	ТС	AI	AC
NoGuideline	39.57	39.57	31.05	29.73	84.15	84.15	64.99	61.96	<u>35.11</u>	<u>35.11</u>	27.16	25.32	42.27	42.27	32.38	31.56
Guideline-H	40.71	40.71	30.76	28.64	56.30	56.30	44.82	43.13	_	-	-	-	-	-	-	-
Guideline-P	51.46	51.46	37.82	35.20	72.86	72.86	55.01	53.73	34.38	34.38	<u>28.04</u>	<u>26.35</u>	67.92	67.92	52.29	44.93
<b>Guideline-PN</b>	<u>49.60</u>	<u>49.60</u>	35.80	32.81	<u>80.77</u>	<u>80.77</u>	<u>63.20</u>	<u>60.34</u>	40.89	40.89	30.04	27.18	<u>75.35</u>	<u>75.35</u>	60.85	57.10
Guideline-PS	47.93	47.93	<u>37.19</u>	<u>34.88</u>	79.23	79.23	59.00	56.88	32.41	32.41	24.63	22.78	76.45	76.45	<u>60.42</u>	<u>56.26</u>
Guideline-PN-Int	40.17	40.17	30.46	28.34	51.95	51.95	41.09	39.32	27.11	27.11	21.93	20.81	42.40	42.40	33.22	31.67
Guideline-PS-Int	39.51	39.51	31.27	30.26	53.70	53.70	42.62	41.10	31.61	31.61	26.70	24.96	52.60	52.60	41.06	39.46

Table 2: Evaluation results (%) for end-to-end EE tasks trained on complete train data. Models trained with Negative Samples (w/ NS) include negative example augmentation. (Best and Second Best performances)

object is valid (e.g., missing arguments or including hallucinated arguments) and comparing it with the ground truth. This code-based evaluation thus prevents the tedious string-matching process adopted in prior work (Li et al., 2021).

271

272

274

275

276

277

278

279

280

285

290

291

296

297

298

299

301

304

307

Model Training. We experimented with the LLaMA-3.1-8B-Instruct model (Grattafiori et al., 2024), selected for its demonstrated proficiency in processing structured code-based inputs and generating coherent outputs. When instruction-tuning the model under the Guideline-P, PN, and PS variants, we randomly sample one of the generated guidelines, a strategy found to prevent the model from memorizing specific guidelines in Sainz et al. (2024). For parameter-efficient training, we implemented rsLoRA (Kalajdzievski, 2023) using the Unsloth framework (Daniel Han and team, 2023).

We include all details about datasets, evaluation, and model training in Appendix A-B.

### 3.2 RQ1: Do the annotation guidelines allow an LLM to more precisely extract occurring events?

To assess the impact of incorporating annotation guidelines in the EE instructions, we compare instruction-tuning an LLM with and without guidelines. We hypothesize that including the annotation guidelines can help the LLM more easily distinguish between similar event types. To understand its impact, we also compare this approach with a "negative sampling (NS)" approach. Specifically, we instruction-tune the LLM on an augmented training set, where each training example is supplemented with 15 randomly selected negative samples, i.e., triplets of  $(e_{neg}, X, \phi)$  with non-existing event type  $e_{neg}$  yielding empty extraction output. We note that annotation guidelines and negative sampling are two complementary approaches for an LLM to learn to distinguish between event types. In our experiments, we thus evaluated the effect of annotation guidelines in two independent settings: (1) training on the original training set (**w/o NS**) and (2) training on the negative sample-augmented training set (**w/ NS**).

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

340

341

342

343

344

345

346

Table 2 shows the results. In the **w/o NS** setting, including annotation guidelines (**Guideline-P**, **PN**, and **PS**) consistently improves performance across both datasets. Our analysis in Section 3.6 further validated that the guidelines indeed enable the LLM to understand the nuanced differences between event types. On ACE w/o NS, **Guideline-P** achieves the highest scores across all four metrics, leading to around 10% TC and 5% AC gains over **NoGuideline**. Similarly, on RichERE w/o NS, **Guideline-PN** outperforms **NoGuideline** by about around 5% TC and 2% AC.

Training the LLM with augmented negative samples, as we expected, helps the model better distinguish between event types; for example, NoGuideline in the w/ NS setting achieves 30% higher AC on ACE and 6% higher AC on RichERE, compared to its counterparts in the w/o NS setting. However, the effects of annotation guidelines in the w/ NS setting diverge between the two datasets. For ACE, adding the guidelines in the instruction does not offer a further advantage, where NoGuideline and Guideline-PN achieved a comparable, the best performance, while all other guideline variants do not show to help. On RichERE, however, the benefit of annotation guidelines complements the negative samples', where Guideline-PN and Guideline-PS achieve around 25% gain on AC over NoGuideline. We notice that RichERE is annotated with a smaller training set but defines more fine-grained event schemas than ACE; for example, the coursergrained Transport event type in ACE is represented by two event types, i.e., TransportPerson

Experiments	ACE w/o NS				ACE w/ NS				RichERE w/o NS				RichERE w/ NS			
	TI	ТС	AI	AC	TI	ТС	AI	AC	TI	TC	AI	AC	TI	TC	AI	AC
NoGuideline	10.60	10.60	5.19	3.68	31.64	31.64	25.91	24.22	19.87	19.87	13.34	11.69	36.29	36.29	28.15	25.58
Guideline-H	29.01	29.01	16.37	14.78	32.62	32.62	25.35	22.87	-	-	-	-	-	-	-	-
Guideline-P	<u>36.91</u>	<u>36.91</u>	<u>24.17</u>	<u>21.24</u>	<u>56.99</u>	<u>56.99</u>	43.44	40.51	40.28	40.28	21.97	18.33	<u>62.04</u>	<u>62.04</u>	<u>46.33</u>	<u>42.03</u>
<b>Guideline-PN</b>	30.94	30.94	19.27	17.64	60.29	60.29	<u>42.88</u>	<u>39.95</u>	<u>31.23</u>	<u>31.23</u>	<u>19.48</u>	<u>17.51</u>	67.16	67.16	47.85	43.39
Guideline-PS	40.53	40.53	28.03	26.12	55.1	55.1	41.57	38.91	26.16	26.16	16.64	15.19	58.95	58.95	42.79	38.1
Guideline-PN-Int	34.11	34.11	22.73	21.18	28.31	28.31	23.82	22.37	25.73	25.73	16.75	14.6	33.59	33.59	28.06	26.0
Guideline-PS-Int	30.04	30.04	19.69	16.9	27.96	27.96	21.55	20.37	23.33	23.33	15.35	13.38	34.92	34.92	27.31	25.04

Table 3: Evaluation results (%) on full test data, for end-to-end EE tasks, trained on 2000 train data samples.

		ACE	w/ NS		RichERE w/ NS								
	TI	тс	AI	AC	TI	тс	AI	AC					
NoGuide	37.08	37.08	21.53	19.18	24.98	24.98	15.05	13.15					
Н	29.00	29.00	17.93	16.34	-	-	-	-					
Р	27.95	27.95	15.94	14.21	23.93	23.93	13.56	12.71					
PN	29.60	29.60	17.87	15.92	27.43	<u>27.43</u>	17.10	15.28					
PS	<u>29.85</u>	<u>29.85</u>	<u>19.49</u>	<u>17.04</u>	19.61	19.61	11.77	10.48					
PN-Int	24.34	24.34	14.08	12.56	27.59	27.59	16.21	<u>14.47</u>					
PS-Int	22.51	22.51	13.59	12.48	18.99	18.99	10.67	9.56					

Table 4: Evaluation results (%) for end-to-end EE tasks on full test data, averaged over three runs using 100 training samples. We did not experiment with the "w/o NS" setting because the model performance with 100 training samples is negligible for all variants.

and TransportArtifact. As the guideline provides not only a detailed description of an event type but also a comparison with similar ones (Table 1), the LLM can leverage this information for better EE performance.

347

348

351

352

354

361

366

367

371

# 3.3 RQ2: Are machine-generated annotation guidelines effective?

Interestingly, from Table 2, we noticed that the guidelines provided by the ACE annotators do not yield a performance gain and that the machinegenerated guideline variants are not equally effective. Specifically, **Guideline-H** achieves a comparable performance in w/o NS and an inferior one in w/ NS on ACE; **Guideline-PN-Int** and **Guideline-PS-Int** provide either no or limited performance gain in both w/o NS and w/ NS settings, while **Guideline-P** and **Guideline-PS** are not consistently better than **NoGuideline**. **Guideline-PN** shows to be the most stable, outperforming **NoGuideline** on RichERE and performing comparably to the best model on ACE.

Qualitatively, as shown in Table 1, the humanwritten guidelines (**Guideline-H**) lack explicit contrasts, making event boundaries ambiguous—for instance, Transport (a movement event) and Extradite (a justice event) both involve relocation, yet the fact that only the latter is legally enforced is not clarified in the guidelines. Guideline-**P** provides examples and edge cases of the target event, but these may not be sufficient for the model to distinguish between similar event types. While both Guideline-PS and Guideline-PN have supplied this comparison, -PS shows to be limited by focusing on only sibling differentiations (e.g., Extradite vs. Convict). Finally, surprisingly, the two -Int variants, despite being comprehensive, lead to mixed results. We observed that models tend to overfit to these comprehensive instructions. In contrast, training the models with 5 diverse guidelines per event type as in -PN and -PS avoids this issue, which shares a similar finding as Cai et al. (2024); Sainz et al. (2024).

373

374

375

376

377

378

379

381

382

385

386

387

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

## **3.4 RQ3:** Are the annotation guidelines helpful when there is only a small amount of training data?

With 2000 samples (Table 3), Guideline-P, Guideline-PN and Guideline-P improve NoGuideline on ACE and RichERE w/o NS by up to 30% TC and 20% AC. Unlike our observation on the fulltraining setting, this trend also holds in ACE w/ NS, where guidelines provide a similar advantage. Excitingly, the results also show that annotation guidelines can compensate for limited training data, enabling models trained with only 2000 samples to achieve performance comparable to full-data training. For example, on ACE, Guideline-P w/ NS (2k) outperforms NoGuideline w/o NS (full) by 10% AC; on RichERE, Guideline-PN w/ NS (2k) outperforms NoGuideline (full) by 18% AC in "w/o NS " and 12% AC in "w/ NS".

However, when training data is reduced to 100 samples (Table 4), the benefits become datasetdependent. In ACE w/ NS, NoGuideline slightly outperforms guideline-based models, suggesting

Experiments	RichE	RE w/	o NS —	ACE	RichERE w/ NS $\rightarrow$ ACE			$\Big  \textbf{ ACE w/o NS} \rightarrow \textbf{ RichERE} \\$				ACE w/ NS $\rightarrow$ RichERE				
	TI	тс	AI	AC	TI	тс	AI	AC	ТІ	тс	AI	AC	TI	тс	AI	AC
NoGuideline	29.55	29.55	21.34	16.60	44.10	44.10	33.91	25.17	33.41	33.41	24.34	22.68	37.19	37.19	27.74	25.87
Guideline-P	31.78	31.78	22.51	15.90	61.69	61.69	39.83	27.93	42.95	42.95	31.61	27.79	54.72	54.72	38.63	35.00
Guideline-PN	40.12	40.12	27.78	19.77	<u>63.97</u>	<u>63.97</u>	48.74	36.24	41.72	41.72	29.54	26.10	<u>64.87</u>	<u>64.87</u>	48.25	44.51
Guideline-PS	29.28	29.28	20.13	15.38	64.23	64.23	<u>44.12</u>	<u>32.84</u>	<u>42.33</u>	<u>42.33</u>	<u>29.93</u>	<u>26.73</u>	65.54	65.54	<u>45.57</u>	<u>41.68</u>
Guideline-PN-Int	27.00	27.00	18.91	14.66	35.35	35.35	28.07	21.82	28.65	28.65	22.13	19.87	38.60	38.60	27.46	26.02
Guideline-PS-Int	<u>31.96</u>	<u>31.96</u>	<u>23.60</u>	<u>19.00</u>	51.71	51.71	39.36	31.34	34.33	34.33	26.65	24.24	36.85	36.85	27.69	26.19
In-Distribution	39.57	39.57	31.05	29.73	84.15	84.15	64.99	61.96	35.11	35.11	27.16	25.32	42.27	42.27	32.38	31.56

Table 5: Evaluation of models (%) in cross-schema generalization. **In-Distribution** represents the NoGuideline performance when trained and tested on the same dataset and the same setting (w/o or w/ NS). We did not experiment with Guideline-H as RichERE does not come with human-annotated guidelines.

that with extremely limited data, the model resorts 411 to memorization rather than learning schema con-412 413 straints. In contrast, in RichERE w/ NS, which has more diverse and fine-grained event structures, 414 guidelines remain beneficial-Guideline-PN sur-415 passes NoGuideline by 2% AC, indicating that 416 guidelines help in settings where direct memoriza-417 tion is insufficient. 418

# 3.5 RQ4: Do annotation guidelines improve cross-schema generalization?

419

420

In Table 5, we evaluate different variants' gen-421 eralizability to a new schema in EE. Notably, 422 while ACE and RichERE share the same domain, 423 RichERE has a finer schema design. In **RichERE** 424 w/o NS  $\rightarrow$  ACE, performance remains below 425 the in-distribution baseline. While Guideline-426 PN achieves 40% TC, nearly matching the in-427 distribution score, its AC drops by nearly 10%, 428 429 likely due to RichERE's expanded argument roles that do not always align well with ACE's simpler 430 schema. This suggests that fine-to-coarse schema 431 migration is partially feasible but still faces chal-432 lenges in argument mapping. Contrastive learn-433 ing helps mitigate some of this gap, as seen in 434 Guideline-PS (w/NS), which improves TC to 64% 435 and AC to 32%, highlighting the benefits of struc-436 tured alignment. In contrast,  $ACE \rightarrow RichERE$ 437 generalizes even better, with Guideline-PN (w/ 438 NS) achieving 64% TC and 44% AC, surpass-439 ing the in-distribution baseline by over 22% TC 440 and 12% AC. This suggests that training on ACE, 441 442 which has well-defined event boundaries, provides a stronger foundation for adapting to RichERE's 443 more detailed schema. Since RichERE introduces 444 additional argument roles for certain events in ACE, 445 structured guidelines play a key role in prevent-446

		A	CE		RichERE						
	TI	тс	AI	AC	TI	тс	AI	AC			
NoGuide w/o NS	29.90	29.90	20.70	19.44	32.74	32.74	24.18	22.35			
PN w/o NS	30.88	30.88	21.82	20.15	33.72	33.72	25.24	24.48			
NoGuide w/ NS	79.81	79.81	56.41	53.85	45.70	45.70	35.68	32.69			
PN w/ NS	77.95	77.95	57.30	54.21	69.10	69.10	49.26	44.10			

Table 6: Evaluation results (%) of LLaMA-3.2-1B-Instruct trained on full ACE and RichERE.



Figure 2: Error categorization: CA (Context Ambiguity), PE (Parsing Errors), MAE (Missing Arguments/Events), AE (Argument Errors), TTE (Type/Trigger Errors), and LN (Label Noise).

ing role confusion and ensuring more consistent schema adaptation.

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

#### 3.6 Further Analysis

**Generalization to a Smaller LLM** We experimented with LLaMA-3.2-1B-Instruct for **NoGuideline** and the best-performing guideline variant **Guideline-PN**. Results in Table 6 display a consistent observation compared to experiments with the larger LLaMA-3.1-8B model (Table 2). That is, **Guideline-PN** achieves a comparable or better result than **NoGuideline** and shows the advantage of guidelines, particularly on **RichERE w/ NS**.

Error Analysis We randomly selected 100 examples on each dataset where NoGuideline w/o NS made mistakes and compared them with errors

made by other variants. The results in Figure 2 462 show that, on ACE w/o NS, including the annota-463 tion guidelines leads to increasing ungrammatical 464 code outputs and parsing errors (PE), although it 465 dramatically reduces the event type and trigger errors (TTE). In the case of w/ NS, guidelines help 467 in almost all aspects, with the majority of remain-468 ing errors being caused by missing arguments or 469 events (MAE) and label noise (LN). On RichERE, 470 however, we observe that for both w/o and w/ NS 471 cases, the annotation guidelines enhance the model 472 performance in all dimensions. 473

Effectiveness of Guidelines per Event Type Frequency As shown in Figure 3, frequent event types show consistent improvements with annotation guidelines, as indicated by the green bars, suggesting that even well-represented events benefit from enhanced annotations. Only a few declines (red bars) occurred with mid- to low-frequency event types, whereas most event types still benefit from the guidelines. In fact, these event types, especially the very rare ones (top of the figure), generally present larger gains (i.e., longer green bars) than the more frequent types, which demonstrates a unique advantage of annotation guidelines in low-resource settings.

#### 4 Related Work

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

503

507

510

511

LLMs for IE and EE With the growing capabilities of LLMs, recent efforts have explored their potential in IE (Xu et al., 2024) and studied EE as an auxiliary task. Existing LLM-based IE methods generally fall into two categories: In-Context Learning (ICL) and Supervised Fine-Tuning (SFT). ICL-based approaches (Li et al., 2023b; Guo et al., 2023; Ashok and Lipton, 2023; Wang et al., 2023b) rely on providing a few-shot context within prompts, enabling LLMs to infer structured information without explicit parameter updates. While being data-efficient, they were found to misinterpret the task specifications (Gao et al., 2024) and suffer from brittle sensitivity to prompt phrasing and example ordering (Gao et al., 2023). In addition, they also incur prohibitive costs due to the lengthy reasoning chains especially for complex tasks. In contrast, SFT-based methods (Lu et al., 2023; Wang et al., 2023a; Gui et al., 2024; Zhou et al., 2024) fine-tune LLMs on annotated datasets, which can significantly improve their EE performance. Our work deepens this line of research and particularly explores the inclusion of



Figure 3: Impact of guidelines on AC scores per ET, sorted by frequency in the full training set. Smaller index indicate a higher frequency. Green/red bars indicate improvements/declines. Dashed/solid lines denote average AC scores without/with guidelines.

annotation guidelines in instructions. While there have been existing works on similar topics, they did not focus on EE (Sainz et al., 2024) or instruction tuning (Pang et al., 2023).

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

Code Prompts for EE While EE tasks are typically represented in texts, code-based prompting has emerged as a promising alternative, leveraging structured representations to enhance schema adherence. Early works have applied code-style prompts to event argument extraction (Wang et al., 2023b) and other IE tasks (Li et al., 2023b), demonstrating potential but often underperforming compared to SFT-based models due to the absence of fine-tuning. EventRL (Gao et al., 2024) utilizes outcome supervision with specific reward functions to reduce information mismatch and hallucination. KnowCoder (Sainz et al., 2024; Li et al., 2024) addresses this limitation by introducing a comprehensive schema representation in code format, integrating taxonomies, constraints, and structured definitions. Complementary to these works, we study generating annotation guidelines to enhance the instruction tuning of LLMs for code-formatted EE and demonstrate their effectiveness.

### 5 Conclusion

We demonstrate that incorporating structured annotation guidelines improves the instruction-tuning of LLMs for EE, bridges the data gap when only a limited amount of training data is available, and enhances the model's cross-schema generalization. Our explorations of guideline generation also highlight the promise of automatically generating effective instructions.

## 6 Limitations

While our study demonstrates the benefits of structured annotation guidelines for event extraction, 547 several limitations remain. First, our evaluation is 548 limited to two datasets (ACE and RichERE), both within the news domain, which may not fully capture how guidelines generalize to other domains 551 such as biomedical or legal text. Future work should assess whether schema differences in other domains exhibit similar trends. Second, while we analyze guideline length and diversity, we do not 555 explicitly optimize guideline generation, leaving open the question of how to best balance concise-557 ness and informativeness. Exploring adaptive methods that retrieve or refine guidelines dynamically during training and inference could further improve 560 561 efficiency. Lastly, our study primarily focuses on instruction-tuning an LLM with predefined event schemas; however, real-world applications often require handling previously unseen event types. Investigating how structured guidelines can aid zero-565 566 shot or few-shot event extraction remains an important avenue for future research. 567

#### References

568

571

572

573

574

575

576

577

580

582

588

591

595

- Dhananjay Ashok and Zachary C. Lipton. 2023. Promptner: Prompting for named entity recognition. *Preprint*, arXiv:2305.15444.
- Zefan Cai, Po-Nien Kung, Ashima Suvarna, Mingyu Ma, Hritik Bansal, Baobao Chang, P. Jeffrey Brantingham, Wei Wang, and Nanyun Peng. 2024. Improving event definition following for zero-shot event detection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2842–2863, Bangkok, Thailand. Association for Computational Linguistics.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Michael Han Daniel Han and Unsloth team. 2023. Unsloth.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- George Doddington, Alexis Mitchell, Mark Przybocki, Lance Ramshaw, Stephanie Strassel, and Ralph Weischedel. 2004. The automatic content extraction (ACE) program – tasks, data, and evaluation. In *Proceedings of the Fourth International Conference*

*on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA). 596

597

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

- Jun Gao, Huan Zhao, Wei Wang, Changlong Yu, and Ruifeng Xu. 2024. Eventrl: Enhancing event extraction with outcome supervision for large language models. *Preprint*, arXiv:2402.11430.
- Jun Gao, Huan Zhao, Changlong Yu, and Ruifeng Xu. 2023. Exploring the feasibility of chatgpt for event extraction. *Preprint*, arXiv:2303.03836.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Honghao Gui, Shuofei Qiao, Jintian Zhang, Hongbin Ye, Mengshu Sun, Lei Liang, Jeff Z. Pan, Huajun Chen, and Ningyu Zhang. 2024. Instructie: A bilingual instruction-based information extraction dataset. *Preprint*, arXiv:2305.11527.
- Yucan Guo, Zixuan Li, Xiaolong Jin, Yantao Liu, Yutao Zeng, Wenxuan Liu, Xiang Li, Pan Yang, Long Bai, Jiafeng Guo, and Xueqi Cheng. 2023. Retrieval-augmented code generation for universal information extraction. *Preprint*, arXiv:2311.02962.
- Kuan-Hao Huang, I-Hung Hsu, Tanmay Parekh, Zhiyu Xie, Zixuan Zhang, Prem Natarajan, Kai-Wei Chang, Nanyun Peng, and Heng Ji. 2024. Textee: Benchmark, reevaluation, reflections, and future challenges in event extraction. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 12804–12825.
- Heng Ji and Ralph Grishman. 2008. Refining event extraction through cross-document inference. In *Proceedings of ACL-08: HLT*, pages 254–262, Columbus, Ohio. Association for Computational Linguistics.
- Yizhu Jiao, Ming Zhong, Sha Li, Ruining Zhao, Siru Ouyang, Heng Ji, and Jiawei Han. 2023. Instruct and extract: Instruction tuning for on-demand information extraction. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 10030–10051, Singapore. Association for Computational Linguistics.
- Damjan Kalajdzievski. 2023. A rank stabilization scaling factor for fine-tuning with lora. *Preprint*, arXiv:2312.03732.
- Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023a. Evaluating chatgpt's information extraction capabilities: An assessment of performance, explainability, calibration, and faithfulness. *arXiv preprint arXiv:2304.11633*.
- Peng Li, Tianxiang Sun, Qiong Tang, Hang Yan, Yuanbin Wu, Xuanjing Huang, and Xipeng Qiu. 2023b. CodeIE: Large code generation models are better few-shot information extractors. In *Proceedings* of the 61st Annual Meeting of the Association for

Computational Linguistics (Volume 1: Long Papers),

pages 15339–15353, Toronto, Canada. Association

Qian Li, Jianxin Li, Jiawei Sheng, Shiyao Cui, Jia Wu,

Yiming Hei, Hao Peng, Shu Guo, Lihong Wang,

Amin Beheshti, et al. 2022. A survey on deep learn-

ing event extraction: Approaches and applications.

IEEE Transactions on Neural Networks and Learning

Sha Li, Heng Ji, and Jiawei Han. 2021. Document-level

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 Computa-

Zixuan Li, Yutao Zeng, Yuxin Zuo, Weicheng Ren,

Wenxuan Liu, Miao Su, Yucan Guo, Yantao Liu, Lix-

iang Lixiang, Zhilei Hu, Long Bai, Wei Li, Yidan

Liu, Pan Yang, Xiaolong Jin, Jiafeng Guo, and Xueqi

Cheng. 2024. KnowCoder: Coding structured knowl-

edge into LLMs for universal information extraction.

In Proceedings of the 62nd Annual Meeting of the

Association for Computational Linguistics (Volume 1:

Long Papers), pages 8758–8779, Bangkok, Thailand.

Keming Lu, Xiaoman Pan, Kaiqiang Song, Hongming

Zhang, Dong Yu, and Jianshu Chen. 2023. PIVOINE:

Instruction tuning for open-world entity profiling. In

Findings of the Association for Computational Lin-

guistics: EMNLP 2023, pages 15108-15127, Singapore. Association for Computational Linguistics.

Chaoxu Pang, Yixuan Cao, Qiang Ding, and Ping Luo.

2023. Guideline learning for in-context information

extraction. In Proceedings of the 2023 Conference

on Empirical Methods in Natural Language Process-

ing, pages 15372-15389, Singapore. Association for

Oscar Sainz, Iker García-Ferrero, Rodrigo Agerri,

prove zero-shot information-extraction.

Information Processing Systems.

Oier Lopez de Lacalle, German Rigau, and Eneko Agirre. 2024. GoLLIE: Annotation guidelines im-

Twelfth International Conference on Learning Repre-

Zhengxiang Shi and Aldo Lipani. 2023. Don't stop pretraining? make prompt-based fine-tuning powerful learner. In Thirty-seventh Conference on Neural

Zhiyi Song, Ann Bies, Stephanie Strassel, Tom Riese, Justin Mott, Joe Ellis, Jonathan Wright, Seth Kulick,

Neville Ryant, and Xiaoyi Ma. 2015. From light to rich ERE: Annotation of entities, relations, and

events. In Proceedings of the 3rd Workshop on EVENTS: Definition, Detection, Coreference, and Representation, pages 89-98, Denver, Colorado. As-

sociation for Computational Linguistics.

Computational Linguistics.

sentations.

Association for Computational Linguistics.

for Computational Linguistics.

Systems.

tional Linguistics.

- 670
- 671 672 673
- 674 675 676
- 677
- 679
- 683

689 690

691

697 698

701

- 703
- 706

Saurabh Srivastava, Gaurav Singh, Shou Matsumoto, Ali Raz, Paulo Costa, Joshua Poore, and Ziyu Yao. 2023. MailEx: Email event and argument extraction. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12964–12987, Singapore. Association for Computational Linguistics.

707

708

710

714

716

717

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. Preprint, arXiv:2302.13971.
- Xiao Wang, Weikang Zhou, Can Zu, Han Xia, Tianze Chen, Yuansen Zhang, Rui Zheng, Junjie Ye, Qi Zhang, Tao Gui, Jihua Kang, Jingsheng Yang, Siyuan Li, and Chunsai Du. 2023a. Instructuie: Multi-task instruction tuning for unified information extraction. Preprint, arXiv:2304.08085.
- Xingyao Wang, Sha Li, and Heng Ji. 2023b. Code4struct: Code generation for few-shot event structure prediction. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3640-3663.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, Yang Wang, and Enhong Chen. 2024. Large language models for generative information extraction: A survey. Frontiers of Computer Science, 18(6):186357.
- Wenxuan Zhou, Sheng Zhang, Yu Gu, Muhao Chen, and Hoifung Poon. 2024. Universalner: Targeted distillation from large language models for open named entity recognition. Preprint, arXiv:2308.03279.

10

In The

## 747 A Pr

748

749

750

751

753

757

758

759

761

764

765

## A Preprocessing and Data Sampling

For both datasets, ACE and RichERE, we follow the TextEE standardization (Huang et al., 2024) and formulate them as sentence-level EE tasks. We use the "split 1" data split of TextEE, but only sample a subset of 100 examples from its development (dev) set for better training efficiency. Specifically, we ensure that for each event type, two event instances will be included in our dev set, prioritizing those with larger coverages of arguments, with the remaining being examples with no event occurrences. The datasets are then converted to the code format shown in Figure 1. Table 7 summarizes dataset statistics.

Dataset	#Event Types	#Role Types	Instances (train/dev/test)
ACE05 (Dodding- ton et al., 2004)	33	22	16531/1870/2519
RichERE (Song et al., 2015)	38	35	9105/973/1163

Table 7: Dataset statistics. For efficiency purposes, in our experiments, we curated a subset of 100 examples as our development (dev) set.

To perform the low-data experiments (RQ3), we additionally create the following subsets of the full training set for each dataset. **Train2k** includes uniformly sampled 2,000 examples from the full training set. **Train100-1/2/3** are three distinct subsets including 100 examples from the full training set, each of which was selected following the same procedure as how we prepare the dev set, ensuring all event types are included and prioritizing instances covering more arguments.

#### **B** Evaluation Methodology and Metrics

Evaluation Methodology. Our methodology contrasts with GoLLIE (Sainz et al., 2024), which 773 follows a pipeline-based structure and selectively includes only parent event types in its prompts, lim-775 iting granularity in event representation. For argument extraction, GoLLIE further restricts schema inclusion to sibling event types, introducing man-778 ual design choices that reduce automation and scal-779 ability. To ensure fair and comprehensive evaluation, we adopt a methodology that enumerates 782 all possible event types for each test and development sample during prompt construction. Unlike setups where only the gold-standard event schema is included in the prompt, we avoid implicit event detection bias-if the correct event type were pro-786

vided, the model would not need to identify the event type itself and could directly extract arguments, which would not reflect its real performance on real-world data. Due to these fundamental differences in methodology, we do not compare our results with GoLLIE. 787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

## C Prompt Design and Model Training

Model. We conducted experiments on an instruction-tuned LLaMA-3-8B model, selected for its demonstrated proficiency in processing structured code-based inputs and generating coherent outputs. For parameter-efficient adaptation, we implement RSLoRA (Kalajdzievski, 2023), applying LoRA transformations to all linear layers in the transformer blocks following the methodology of Dettmers et al. (2024). Key hyperparameters-including LoRA rank (64), scaling factor  $\alpha$  (128), and batch size (32)—were determined through preliminary experiments to balance computational efficiency with model performance. The models were trained for 10 epochs using a single NVIDIA A100 GPU (80GB VRAM), with early stopping triggered after three consecutive validation steps without improvement. We adopt a cosine learning rate scheduler with an initial rate of 1e-5 and a warmup period of 350 steps. Input sequences are padded to 3,000 tokens to maintain consistency while accommodating long-form code structures. To ensure reproducibility and minimize memory fragmentation, we implement deterministic padding and truncation strategies.

**Prompt Design.** We adopt a structured prompt format consisting of four components: (i) task instruction,(ii) event schema, (iii) input text, and (iv) expected output, formatted as a structured event representation. Our approach follows a schemafirst prompting strategy, where event definitions are explicitly encoded in a structured format to enhance model comprehension of event relations and argument constraints. For each input instance, a randomly sampled guideline definition is used to annotate the event schema, ensuring that the model is exposed to multiple rephrasings rather than memorizing and overfitting on a static definition. Formally, we prepare the input sequence as follows: "[BoS]  $-task_instruction$  (I) -annotated event schema  $(E_e)$  \$-input\_sample  $(X_i)$  [EoS]" where the event schema  $E_e$  for an event e is annotated with one of the generated guideline definitions.

```
837
838
            You are an expert in annotating NLP datasets for event extraction. Your task is to
               generate "detailed" annotation guidelines for the event type Acquit which is a child event type of super class JusticeEvent.
839
841
842
            Input Format will be as following
843
844
            Event Schema:
            Event Name and its parent class
845
            Arguments:
847
            Arguments separated by new lines. If there are no arguments None will be given.
848
849
            Examples
851
            Instructions:
852
            1) Identify and list all unique arguments related to the event type.
853
            2) Define the event type and each argument. You can take help of examples below to
               understand the events and their arguments.
854
855
            3) Please remember that the examples may not cover all the arguments in the list. In
                 some cases, you may not have arguments at all, in such cases, you can have an
857
               empty list for arguments.
            4) For each definition, provide 5 illustrative definitions in JSON format. For
859
               events you can add example triggers and the explanation of the events such as
               edge cases and other critical details starting with "The event can be triggered
861
               by \ldots ". Similarly for arguments also you can add examples, and detailed
862
                information for them including any edge case or domain knowledge starting with "
863
               Examples are ... ".
864
            5) Remember to not generate any additional information such as examples, etc. and
865
               strictly follow the output format shown below.
            6) Remember also to add detailed information for the events and arguments so that
867
               the annotators who are not familiar with machine learning and NLP can still
                solve the task. Remember to add required domain knowledge and please cover the
               edge cases when possible.
870
            7) Remember that while generating examples for the event or attributes you should
871
               generate diverse set of triggers or argument values rather than picking them
               from the examples I have provided for each of the 5 generated guidelines.
872
873
            Output Format:
874
875
            {
              "Event Definition": [
876
                "Definition 1",
877
                "Definition 2",
878
879
                "Definition 3"
                "Definition 4",
880
                "Definition 5"
881
              ٦.
              "Arguments Definitions": {
                "Argument1": [
884
                  "Definition 1"
885
                  "Definition 2",
                  "Definition 3",
                  "Definition 4"
                  "Definition 5"
890
                ],
891
                "Argument2": [
                  "Definition 1",
"Definition 2",
"Definition 3",
893
894
                  "Definition 4"
                  "Definition 5"
896
897
                ٦
                // Add additional arguments as necessary
              }
900
            }
901
            Event Schema:
902
903
            Acquit which is a child event type of super class JusticeEvent
904
            Arguments:
905
            Argument 1 -> adjudicator
           Argument 2 -> defendant
906
```

```
Example 1
### Input Text ###
Sentence 1.
### Event Trigger ###
[event trigger]
### Event Arguments ###
For argument "defendant" extracted spans ['x']
For argument "adjudicator" extracted spans ['y']
Example 2
### Input Text ###
Sentence 2.
### Event Trigger ###
[event trigger]
### Event Arguments ###
For argument "defendant" extracted spans ['a']
(...)
```

Listing 1: Prompt example for generating Guideline-P, Guideline-PN, and Guideline-PS.

907

908

909

910

911 912

913

914

915 916

917

918

919 920

921 922

923 924

928

927

928 929 930

931 932

933

934

935

936 937

938

939 940

941 942

943 944

945

946

947

948

949

950

951

952

953

954

955

956

957

958 959

960

961

962 963 964

965

966

967

968

969 970

971 972

**Prompt for generating consolidated guidelines.** The exact prompts used for generating consolidated guidelines - Guideline-PN-Int, and Guideline-PS-Int is shared below

```
You are an expert in summarizing NLP event extraction guidelines. Your goal is to
   consolidate multiple detailed descriptions into a single concise, comprehensive
    "Intergrated" guideline.
### Input Format ###
Event Type: Event Type Name
  `json
{
  "Event Definition": [
    "Definition 1",
"Definition 2",
    "Definition 3",
"Definition 4",
    "Definition 5"
  ],
  "Arguments Definitions": {
    "mention": [
      "Definition 1"
      "Definition 2",
      "Definition 3",
      "Definition 4"
      "Definition 5"
   ],
    "Argument1": [
      "Definition 1",
      "Definition 2"
      "Definition 2",
      "Definition 4"
      "Definition 5"
    ٦.
    // Add additional arguments as necessary
  }
}
### Task ###
1. Integrated the 5 definitions under "Event Definition" into a single definition:
   - Highlight all critical points and examples from the five definitions.
   - Ensure the description is concise, comprehensive, and clear, using formal
       language that non-experts can understand.
2. Do the same for each argument under "Arguments Definitions," producing a single
  intergrated definition for each.
```

```
### Output Format ###
···json
{
  "Event Definition": "Consolidated intergrated guideline for the event type.",
  "Arguments Definitions": {
    "mention": "Consolidated intergrated guideline for the mention argument.",
    "Argument1": "Consolidated intergrated guideline for Argument1."
    "Argument2": "Consolidated intergrated guideline for Argument2."
    // Add additional arguments as necessary
  }
}
### Guidelines to Summarize ###
Event Type: prompt_Acquit(JusticeEvent)
```ison
{
    "Acquit(JusticeEvent)": {
        "description": [
            "Definition 1"
            "Definition 2",
            "Definition 3"
            "Definition 4"
            "Definition 5"
        ]
    },
    "attributes": {
        "mention": "The text span that triggers the event."
        "adjudicator": [
            "Definition 1",
            "Definition 2"
            "Definition 2",
            "Definition 4"
            "Definition 5"
        ],
        "defendant": [
            "Definition 1"
            "Definition 1",
"Definition 2",
            "Definition 3",
            "Definition 4"
            "Definition 5"
        ]
    }
}
```

975 976

977

978

981 982

983

984

986 987

991

993

994 995

996

997

999

1000

1002

1003 1004

1005

1006

1008

1009

1011

1012 1013

1014 1015

1016

1017

1019

1020

1021

1022

1023

1024

1025 1026

1028

1030

1031

1032

1033

1034 1035

1036

Listing 2: Prompt example for generating consolidated guidelines: Guideline-PN-Int, and Guideline-PS-Int.

## D Dataset examples across multiple guideline settings

The below JSON example illustrates an event extraction task from the ACE dataset under the No Guideline setting. It defines how structured events are extracted from text, specifying event triggers, types, arguments, and roles. The instruction explains the task, the input provides a natural language sentence and its conversion into a structured Python-style format. The output presents the extracted event, including its trigger ("extradited") and associated arguments (e.g., "government" as the agent, "him" as the person).

```
{
    "doc_id": "APW_ENG_20030306.0191",
    "wnd_id": "APW_ENG_20030306.0191-6",
    "instance_id": "821",
    "dataset_name": "ace05-en",
    "task_type": "E2E",
    "is_auth": "0",
    "instruction": "# This is an event extraction task where the goal is to extract
    structured events from the text. A structured event contains an event trigger
    word, an event type, the arguments participating in the event, and their roles
```

in the event. For each different event type, please output the extracted	1037
information from the text into python-style dictionaries where the first key	1038
will be 'mention' with the value of the event trigger. Next, please output the	1039
arguments and their roles following the same format. The event type	1040
definitions and their argument roles are defined next.",	1041
"input": "# The following lines describe the task definition\n\n@dataclass\nclass	1042
Extradite(JusticeEvent):\n mention: str\n agent: List\n destination:	1043
List $n$ origin: List $n$ person: List $n \#$ This is the text to analyze $h$	1044
ntext = $\$ The post-Milosevic government later extradited him to the U.N. war	1045
crimes tribunal in The Hague, the Netherlands.\"\n\n# The list called result	1046
should contain the instances for the following events according to the	1047
guidelines above:\nresult = \n".	1048
"output": "[Extradite(\n mention=\"extradited\".\n person=[\"him\"]. \n	1049
destination=[\"Hague\"]. \n agent=[\"government\"].\n origin=[]\n)]"	1050
}	1854

1054

1094 1095

1097

1099

1100

Listing 3: Prompt example for generating consolidated guidelines: Guideline-PN-Int, and Guideline-PS-Int.

**NoGuideline** Shown below is an example from the NoGuideline setting in python code format with no doc string and argument definitions.

```
1055
#Task Instruction
# This is an event extraction task where the goal is to extract structured events
  1057
    from the text. A structured event contains an event trigger word, an event type,
   1058
    the arguments participating in the event, and their roles in the event. For
   1059
    each different event type, please output the extracted information from the text
   1060
   1061
    into python-style dictionaries where the first key will be 'mention' with the
    value of the event trigger. Next, please output the arguments and their roles
  1062
    following the same format. The event type definitions and their argument roles
   1063
   are defined next.
   1065
#Input
# The following lines describe the task definition
   1067
   1068
@dataclass
   1069
class Extradite(JusticeEvent):
  1070
    mention: str
  1071
   1072
    agent: List
    destination: List
  1073
  1074
    origin: List
    person: List
  1075
  1076
# This is the text to analyze
  1077
text = "The post-Milosevic government later extradited him to the U.N. war crimes
   1078
    tribunal in The Hague, the Netherlands.'
  1079
# The list called result should contain the instances for the following events
   1081
    according to the guidelines above:
   1082
  1083
result =
  1084
#Output
   1085
  1086
[Extradite(
    mention="extradited",
    person=["him"],
destination=["Hague"],
  1088
    agent=["government"],
   1090
    origin=[]
   1091
)]
```

Guideline-PN Shown below is an example from the Guideline-PN setting in python code format.

**#Task Instruction** # This is an event extraction task where the goal is to extract structured events from the text. A structured event contains an event trigger word, an event type, 1098 the arguments participating in the event, and their roles in the event. For each different event type, please output the extracted information from the text into python-style dictionaries where the first key will be 'mention' with the 1101 value of the event trigger. Next, please output the arguments and their roles 1102 1103 following the same format. The event type definitions and their argument roles

```
are defined next.
#Input
# The following lines describe the task definition
@dataclass
class Extradite(JusticeEvent):
     ""The event is triggered by the act of transferring a person from one
        jurisdiction to another for legal proceedings. Example triggers include '
       \ensuremath{\mathsf{extradite'}} , 'extradition', and 'extraditing'.'
    mention: str # The text span that triggers the event.
    agent: List # Examples are 'court', 'government', 'police department'. The
       agent is the authority or entity responsible for initiating or carrying out
        the extradition process.
    destination: List # Examples are 'jurisdiction', 'Hague', 'state'. The
       destination is the place to which the person is being extradited.
    origin: List # Examples are 'state', 'headquarters'. The origin is the place
       from which the person is being extradited.
    person: List # Examples are 'she', 'him', 'her'. The person is the individual
       being extradited.
# This is the text to analyze
text = "The post-Milosevic government later extradited him to the U.N. war crimes
    tribunal in The Hague, the Netherlands.'
# The list called result should contain the instances for the following events
   according to the guidelines above:
result =
#Output
[Extradite(
    mention="extradited",
    person=["him"],
    destination=["Hague"],
    agent=["government"],
    origin=[]
)]
```

1105

1107

1108 1109

1110

1111 1112

1113

1114

1115 1116

1117

1118

1119

1120 1121

1122

1123

1124 1125

1126 1127

1128

1129 1130

1131

1132 1133

1134

1135

1136

1137

1138 1139

1149

1142 1143

1144 1145

1146 1147

1148

1149 1150

1151

1152

1153

1154 1155

1156

1157

1158 1159

1160

1161

1162

1163

1164

1165

1166

1167 1168

1169 1170

1171 1172

1173

**Guideline-PN-Int** Similarly, shown below is an example from the Guideline-PN-Int setting in python code format.

```
# The following lines describe the task definition
@dataclass
class Extradite(JusticeEvent):
    """The Extradite event is triggered by the legal process of transferring a
       person from one jurisdiction to another for legal proceedings, such as
       facing charges or serving a sentence. This event involves formal actions by legal authorities and the movement of the individual across jurisdictions.
       Key triggers include terms like 'extradite', 'extradition', and 'extraditing
        '. It is distinct from events like 'ArrestJail' and 'ReleaseParole', as it
        specifically involves cross-jurisdictional transfer rather than initial
       detention or release from custody."""
   mention: str # The text span that triggers the event.
    agent: List # The agent is the authority or entity responsible for initiating
        or carrying out the extradition process, typically a legal or governmental
       body. Examples include 'court', 'government', and 'police department'.
   destination: List # The destination is the place to which the person is being
       extradited, where they will face legal proceedings or serve a sentence.
       Examples include 'jurisdiction', 'Hague', and 'state'.
   origin: List # The origin is the place from which the person is being
        extradited, where they are currently held or from where they are being
        transferred. Examples include 'state' and 'headquarters'
    person: List # The person is the individual being extradited, the subject of
        the legal transfer. Examples include 'she', 'him', and 'her'.
# This is the text to analyze
text = "The post-Milosevic government later extradited him to the U.N. war crimes
   tribunal in The Hague, the Netherlands."
```

# The list called result should contain the instances for the following events	1174
according to the guidelines above:	1175
result =	1176
	1177
#Output	1178
[Extradite(	1179
<pre>mention="extradited",</pre>	1180
person=["him"],	1181
<pre>destination=["Hague"],</pre>	1182
<pre>agent=["government"],</pre>	1183
origin=[]	1184
)]	1185