# MAVEN-ARG: Completing the Puzzle of All-in-One Event Understanding **Dataset with Event Argument Annotation**

Anonymous ACL submission

### Abstract

001

011

012

022

034

042

Understanding events in texts is a core objective of natural language understanding, which requires detecting event occurrences, extracting event arguments, and analyzing inter-event 005 relationships. However, due to the annotation challenges brought by task complexity, a largescale dataset covering the full process of event understanding has long been absent. In this paper, we introduce MAVEN-ARG, which augments MAVEN datasets with event argument annotations, making the first all-in-one dataset supporting event detection, event argument extraction (EAE), and event relation extraction. As an EAE benchmark, MAVEN-ARG offers 015 three main advantages: (1) a comprehensive schema covering 162 event types and 612 argument roles, all with expert-written definitions and examples; (2) a large data scale, containing 98, 591 events and 290, 613 arguments obtained with laborious human annotation; (3) the exhaustive annotation supporting all task variants of EAE, which annotates both entity and non-entity event arguments in document level. Experiments indicate that MAVEN-ARG is quite challenging for both fine-tuned EAE models and proprietary large language models (LLMs). Furthermore, to demonstrate the benefits of an all-in-one dataset, we preliminarily explore a potential application, future event prediction, with LLMs. MAVEN-ARG and our code will be publicly released.

#### 1 Introduction

Conveying information about events is a core function of human languages (Levelt, 1993; Pinker, 2013; Miller and Johnson-Laird, 2013), which highlights event understanding as a major objective for natural language understanding and a foundation for various downstream applications (Ding et al., 2015; Li et al., 2018a; Goldfarb-Tarrant et al., 2019; Huang et al., 2019; Wang et al., 2021a). As illustrated in Figure 1, event understanding is typically organized as three information extraction tasks (Ma



Figure 1: Illustration for the overall event understanding, consisting of event detection, event argument extraction, and event relation extraction tasks.

et al., 2021; Peng et al., 2023b): event detection (ED), which detects event occurrences by identifying event triggers and classifying event types; event argument extraction (EAE), which extracts event arguments and classifies their argument roles; event relation extraction (ERE), which analyzes the coreference, temporal, causal, and hierarchical relationships among events.

043

045

046

047

049

055

061

062

063

064

065

066

067

068

069

Despite the importance of event understanding, a large-scale dataset covering all the event understanding tasks has long been absent. Established sentence-level event extraction (ED and EAE) datasets like ACE 2005 (Walker et al., 2006) and TAC KBP (Ellis et al., 2015, 2016; Getman et al., 2017) do not involve event relation types besides the basic coreferences. RAMS (Ebner et al., 2020) and WikiEvents (Li et al., 2021) extend EAE to the document level but do not involve event relations. ERE datasets are mostly developed independently for coreference (Cybulska and Vossen, 2014), temporal (Chambers et al., 2014; Ning et al., 2018), causal (Mirza et al., 2014; Mostafazadeh et al., 2016b; Caselli and Vossen, 2017), and subevent (Hovy et al., 2013; Glavaš and Snajder, 2014) relationships and do not cover event arguments. Given annotation challenges from task complexity, these datasets often cover only

095

100

101

102

103 104

105

106

108

110

111

112

113

114

115

116

117

118

119

121

thousands of events. Due to the inconsistent event schemata and data, these datasets cannot be unified.
This status quo hinders the development of end-to-end event understanding methods and limits the potential for event-based downstream applications.

MAVEN (Wang et al., 2020) is the largest humanannotated ED dataset, with a high-coverage event schema for general-domain events. Based on it, Wang et al. (2022) further annotates the first unified ERE dataset MAVEN-ERE, which covers all four types of event relationships and has a massive scale with more than one million event relations. Building on the sustained efforts of these works over years, we complete the puzzle of an all-in-one event understanding dataset in this work. We construct MAVEN-ARG, which provides exhaustive event argument annotations based on MAVEN.

Beyond finishing an all-in-one event understanding dataset, three main advantages of MAVEN-ARG make it a valuable EAE benchmark. (1) Comprehensive Event Schema. The original MAVEN schema only defines event types but without argument roles. We engage experts to enhance MAVEN schema with argument roles and to write detailed definitions for them, which help annotators and can also serve as task instructions for prompting large language models. The resulting event schema contains 162 event types, 612 argument roles, and 14,655 words of definitions, which well cover general-domain events. (2) Large Data Scale. MAVEN-ARG comprises 107, 507 event mentions, 290, 613 event arguments, and 129, 126 entity mentions, all of which are human annotated. To our knowledge, this makes it the largest EAE dataset currently available. (3) Exhaustive Annotation. The development of EAE has seen many variations in task settings, including annotating only the topic event (Ebner et al., 2020; Tong et al., 2022) of a document or all fine-grained events (Walker et al., 2006), annotating event arguments at the sentence level (Walker et al., 2006) or document level (Ebner et al., 2020; Li et al., 2021), and limiting event arguments to entities (Walker et al., 2006; Li et al., 2021) or including non-entity arguments (Grishman and Sundheim, 1996; Parekh et al., 2023). MAVEN-ARG adopts the most exhaustive annotation. We annotate event arguments for all finegrained events at the document level, covering both entity and non-entity arguments. This enhances the dataset's utility for benchmarking and developing a wide range of EAE methods.

In the experiments, we reproduce several recent

state-of-the-art EAE models as baselines and also 122 evaluate large language models with in-context 123 learning. Experimental results show that they can 124 only achieve at most 40% F1 scores, which is far 125 from promising. It indicates that MAVEN-ARG 126 is quite challenging and more research efforts are 127 needed to develop practical EAE methods. Further-128 more, to demonstrate the advantage of an all-in-one 129 event understanding dataset for enabling sophisti-130 cated event-based applications, we conduct a pre-131 liminary exploration of *future event prediction*. We 132 sample causally related event chains from MAVEN-133 ARG and prompt LLMs to predict future events, 134 including their types and arguments. Experiments 135 show that while most of the predictions are reason-136 able, they seldom align with the actual future. We 137 encourage future work to further explore this appli-138 cation and hope MAVEN-ARG can help improve 139 EAE and develop diverse event-based applications. 140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

158

159

160

161

163

164

165

166

167

168

169

171

# 2 Dataset Construction

# 2.1 Event Schema Creation

The event schema of MAVEN (Wang et al., 2020) covers a broad range of general-domain events and has a well-defined hierarchical structure. To enable event argument annotation based on MAVEN, one author and two engaged linguistic experts devoted three years to manually designing argument roles for MAVEN schema. Each argument role is accompanied by informative text definitions that are easy to understand, and each event type is provided with detailed annotation examples. An example in shown in appendix A.1. This not only helps annotators understand their tasks but also can prompt LLMs to perform EAE via in-context learning. To ensure quality, the argument role design for each event type is reviewed by at least one expert.

Our event schema creation involves the following steps: (1) Initially, to reduce annotation difficulty, we invite ten ordinary annotators to review the event type schema and a portion of the data. Based on their feedback, we deleted 6 event types that are similar to others and renamed 4 event types for clarity. (2) The basic schema is constructed from a simplification and modification of FrameNet (Baker et al., 1998). The *frame elements* in FrameNet are widely considered akin to argument roles (Aguilar et al., 2014; Parekh et al., 2023), but they are often too complex for ordinary annotators to comprehend since FrameNet is primarily constructed for linguistic experts (Aguilar

Dataset	#Event Type	#Argument Role
ACE 2005	33	36
DocEE	59	356
WikiEvents	50	59
RAMS	139	65
MEE	16	23
GENEVA	115	220
MAVEN-ARG	162	612

Table 1: Event schema statistics of MAVEN-ARG compared with other datasets.

et al., 2014). Therefore, for each event type, we 172 manually select frame elements related to describ-173 ing events and suitable for annotation as MAVEN-174 ARG argument roles from their FrameNet equiva-175 lents, and we rewrite the definitions and examples. 176 (3) Extending argument roles based on the 5W1H 177 analysis (What, Where, When, Why, Who, How) 178 for describing events (Karaman et al., 2017; Ham-179 borg et al., 2019). Temporal and causal relations from event relation extraction describe When and 181 Why, while the event type describes What. We primarily refer to Who (participants), Where (locations), and How (manners, instruments, etc.) to 184 design argument roles. (4) Considering the hierar-185 chical structure. When designing subordinate types, we inherit and refine the argument roles of their su-188 perordinate types. (5) Sampling data to check if any event argument is missing. 189

Schema Statistics After the schema design, the 190 final MAVEN-ARG schema contains 162 event 191 types, 612 unique argument roles, and 14,655 words of definitions. Taking inspiration from se-193 mantic role labeling (Fillmore, 1976; Banarescu 194 et al., 2013), we tend to let the argument roles 195 sharing the same semantic role use the same name but distinguish them with different textual defi-197 nitions. For instance, we do not use Killer for the Killing event type and use Attacker for the 199 Attack event type. Instead, we use Agent to denote them both but write different definitions for 201 them. This is to encourage the knowledge transfer 202 between EAE for different event types. Therefore, 612 is the number of argument roles with unique definitions, and there are 143 unique names for all 206 the argument roles. Table 1 compares the event schema size of MAVEN-ARG with existing EAE datasets, including ACE 2005 (Walker et al., 2006), DocEE (Tong et al., 2022), WikiEvents (Li et al., 2021), RAMS (Ebner et al., 2020), MEE (Pouran 210

Ben Veyseh et al., 2022), and GENEVA<sup>1</sup> (Parekh et al., 2023). We can observe that MAVEN-ARG has the largest event schema, which more comprehensively covers the broad range of diverse events and will help develop more generalizable methods.

211

212

213

214

215

216

217

218

219

220

221

224

225

226

227

228

229

230

231

232

233

234

236

237

239

240

241

242

243

245

246

247

248

249

250

251

252

253

254

255

256

257

258

# 2.2 Entity Annotation

The mainstream task setting for EAE (Walker et al., 2006; Li et al., 2021) confines event arguments to entities, which reduces the task's complexity to some extent and provides more definite and standardized extraction results. Hence, before annotating event arguments, we annotate entities for the 4,480 MAVEN documents. We follow the task definition and guidelines of a recent named entity recognition benchmark Few-NERD (Ding et al., 2021), but we only annotate coarse-grained entity types, including Person, Organization, Location, Building, Product, Art, and MISC. To deliver more unambiguous EAE results and reduce the argument annotation difficulty, we also annotate entity coreference, which means judging whether multiple entity mentions refer to the same entity. During entity annotation, we engage 47 annotators, including 8 senior annotators selected during the annotation training. Each document is annotated by three independent annotators and further checked by one senior annotator. The final annotation results are aggregated via majority voting. If the senior annotator judged the accuracy of a document's annotation to be below 90%, the document will be returned to the three first-stage annotators for re-annotation. To check data quality, we calculate Fleiss' kappa (Fleiss, 1971) to measure the inter-annotator agreements. The result for entity recognition is 73.2%, and for entity coreference is 78.4%, both indicating high consistency.

### 2.3 Event Argument Annotation

Based on the event detection annotations of MAVEN and event coreferences of MAVEN-ERE, we conduct event argument annotations. For multiple coreferent event mentions (triggers), only one of them is displayed during annotation to reduce annotation overhead. Once the annotator selects an event trigger, the corresponding argument roles for its event type are displayed on the annotation interface, along with definitions and examples. This ensures that annotators do not have to memorize the lengthy event schema. To annotate an event

<sup>&</sup>lt;sup>1</sup>GENEVA has a larger "full ontology" but is without data. Here we compare with its schema actually used in dataset.

Dataset	#Doc.	#Event	#Trigger	#Arg.	#Entity	#Entity Mention	Fine-grained Event	Doc. Level	Entity Arg.	Non-Entity Arg.
ACE 2005	599	4,090	5,349	9,683	45,486	59,430	$\checkmark$	×	9,683	×
DocEE	27,485	27,485	-	180, 528	-	-	×	$\checkmark$	×	180, 528
WikiEvents	246	3,951	-	5,536	13,937	33,225	$\checkmark$	$\checkmark$	5,536	×
RAMS	3,993	9,124	-	21,237	-	-	×	$\checkmark$	×	21,237
MEE	13,000	17,642	-	13,548	-	190, 592	$\checkmark$	$\checkmark$	13,548	×
GENEVA	-	7,505	-	12,269	-	36,390	$\checkmark$	×	8,544	3,725
MAVEN-ARG	4,480	98,591	107, 507	290, 613	83,645	129, 126	$\checkmark$	$\checkmark$	116,024	174,589

Table 2: Statistics of MAVEN-ARG compared to existing widely-used EAE datasets. "Doc." is short for "Document" and "Arg." is short for "Argument". "-" denotes not applicable due to lack of document structure or corresponding annotations. "Fine-grained Event" means annotating all the events rather than only one topic event for a document. "Doc. Level" means annotating arguments within the whole document rather than only the sentence containing the trigger. For multilingual datasets, we only compare with its English subset.

argument, annotators can either choose an entity from the whole document or select a continuous 260 textual span; once an entity mention is selected, all 261 of its coreferent entity mentions are automatically 262 selected. Annotators have the option to report errors in the event type annotation of a trigger, which allows for the discarding of that trigger. In the annotation process, approximately 4% of triggers are discarded. We employ 202 annotators, including 71 senior annotators selected during annotation 268 training and 33 experts with rich annotation ex-269 periences. The annotation is divided into three 270 phases. Each document is first annotated by an ordinary annotator, and then modified by a senior 272 annotator. Finally, an expert will check whether the 273 annotation accuracy reaches 90%. If not, the doc-274 ument's annotation will be returned to the second 275 phase. To measure data quality, we randomly sam-276 ple 100 documents and conduct the three-phrase annotation for them twice with different annotator groups. The Fleiss' kappa is 68.6%, which indicates a satisfactory level of annotation agreement. More annotation details are shown in appendix A. 281

# 3 Data Analysis

### 3.1 Data Statistics

283

284Table 2 shows the main statistics of MAVEN-ARG285compared with various existing EAE datasets. Ap-286pendix B further shows the statistics of different287splits. We can observe that MAVEN-ARG has two288advantages: (1) MAVEN-ARG has the largest data289scale, surpassing previous datasets by several times.290This ensures that even for long-tail event types,291MAVEN-ARG has sufficient data to fully train and292stably evaluate EAE models. (2) The exhaustive an-293notation of MAVEN-ARG makes it the only dataset

that covers all settings of EAE task. MAVEN-ARG includes complete annotations of entity and event coreference and annotates both entity and non-entity arguments for all fine-grained events at the document level. This allows MAVEN-ARG to support the evaluation of all variants of EAE methods and the development of comprehensive event understanding applications. 294

295

296

297

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

# 3.2 Data Distribution

We present the distributions of the annotated entity and event arguments of MAVEN-ARG in Figure 2. Argument roles with the same name across different event types are merged for presentation clarity. We observe that: (1) The distribution of entity types is generally similar to that of Few-NERD (Ding et al., 2021), demonstrating sufficient diversity. (2) The three most frequent basic argument roles (Agent, Patient, and Location) account for over 60% of event arguments. This highlights their ubiquity and encourages knowledge transfer among different event types in EAE methods. (3) Event arguments exhibit a highly longtailed distribution. The 136 argument roles counted as "Others", each constituting less than 1.5%, collectively accounts for 27.8% of event arguments. The long-tailed distribution of MAVEN-ARG poses a significant challenge to model generalizability.

### 3.3 Trigger-argument Distance

We analyze the distribution of trigger-argument distances in Figure 3. For events with multiple coreferent triggers and entity arguments with multiple entity mentions, the distance is calculated between the nearest trigger-argument pairs. The overall average trigger-argument distance is 37.8. From Figure 3, we observe that while the majority of event



Figure 2: MAVEN-ARG entity and event argument distributions. For clarity, only the top event argument roles are shown and the others are summed up in "Others".

arguments are located near their triggers, which is natural for human writing, a substantial number of arguments are situated far from their triggers, with the furthest exceeding 800 words. This data characteristic challenges the ability of EAE methods to capture long-distance dependencies.

# 4 Experiment

331

332

336

337

338

340

341

345

347

352

361

365

### 4.1 Experimental Setup

Models To assess the challenge of MAVEN-ARG, we evaluate multiple advanced methods. For finetunned EAE models, we implement several stateof-the-art models, including DMBERT (Wang et al., 2019), CLEVE (Wang et al., 2021b), BERT+CRF (Wang et al., 2020), EEQA (Li et al., 2020), Text2Event (Lu et al., 2021), and PAIE (Ma et al., 2022). These methods cover all the mainstream EAE modeling paradigms (Peng et al., 2023c). Their detailed descriptions and implementations are introduced in appendix C.1.

We also evaluate large language models (LLMs) with in-context learning on MAVEN-ARG. Specifically, we select two advanced LLMs, **GPT-3.5** (OpenAI, 2022) and **GPT-4** (OpenAI, 2023), and evaluate them with 2-shot in-context learning. Here 2-shot means using full annotations of two documents as demonstrations. Considering time and cost constraints, we sample 50 documents from the test set for experimentation. We employ the gold trigger evaluation approach (Peng et al., 2023c) to directly assess their EAE performance.

**Evaluation Metric** Considering that MAVEN-ARG covers non-entity argument annotations, traditional evaluation metrics (Peng et al., 2023c) designed only for entity arguments are no longer applicable. By taking each argument role as a question to the document, we propose to view EAE as a **multi-answer question answering** task<sup>2</sup> and



Figure 3: Distribution of distances between triggers and arguments in MAVEN-ARG.

adopt its evaluation metrics (Rajpurkar et al., 2016; Amouyal et al., 2022; Yao et al., 2023), including **bag-of-words F1** and **exact match (EM)**.

Conventional evaluation calculates the micro average over all the entity and event mentions, which we dub it as **mention-level** evaluation. Considering that real-world applications only require the accurate prediction for one of all the coreferent mentions, we propose to consider entity (Li et al., 2021) and event coreference in evaluation. Specifically, for **entity coreference level** evaluation, an entity argument is considered as predicted correctly if one of its mentions is predicted correctly. For **event coreference level** evaluation, an argument is considered as predicted correctly if it is predicted correctly for one of the coreferent triggers.

### 4.2 Experiment Results of Fine-tuned Models

The results of fine-tuned EAE models are shown in Table 3, and we have the following observations:

(1) Existing state-of-the-art EAE models exhibit moderate performance on MAVEN-ARG, which is significantly worse than their results on existing datasets (Peng et al., 2023c). This indicates that MAVEN-ARG is challenging and there is a need for increased efforts in developing practical event understanding models. (2) The BERT+CRF and PAIE models exhibit the best performance, potentially attributable to their ability to model rich interactions between different event arguments. (3) The previous top-performing classification-based models (DMBERT and CLEVE) (Peng et al., 2023c) perform poorly on MAVEN-ARG, which is due to their inability to handle non-entity arguments. Therefore, future research necessitates more flexible approaches to tackle the complex and real-world scenario in MAVEN-ARG. (4) Text2Event notably underperforms. This is potentially due to the intensive annotations of MAVEN-ARG, i.e., a high volume of events and argument annotations within a single document, making generating all events

<sup>&</sup>lt;sup>2</sup>A single role may correspond to multiple argument spans.

Madal	#Domonic	Mention Level					Entity Coref Level				<b>Event Coref Level</b>			
Model	#F al allis	P	R	F1	EM	Р	Ř	F1	EM	Р	R	F1	EM	
DMBERT	110M	19.7	19.7	19.7	19.5	12.5	12.4	12.4	12.3	11.8	11.8	11.8	11.6	
CLEVE	355M	22.1	22.1	22.1	22.0	13.2	13.2	13.2	13.0	12.3	12.2	12.2	12.1	
BERT+CRF	110 <b>M</b>	31.7	31.4	30.9	27.0	33.5	32.8	32.2	27.1	32.3	31.8	31.2	26.3	
EEQA	110 <b>M</b>	21.4	19.5	19.6	15.8	24.5	22.9	22.8	18.8	23.7	22.2	22.1	18.1	
Text2Event	770M	12.9	12.9	12.7	11.3	12.5	12.4	12.1	10.4	10.8	10.7	10.5	9.0	
PAIE	406M	37.2	36.2	35.6	<b>30.3</b>	42.3	<b>41.1</b>	40.5	34.4	42.1	<b>41.0</b>	40.3	34.3	

Table 3: Experimental results (%) of existing state-of-the-art fine-tuned EAE models on MAVEN-ARG.

and arguments at once difficult. It indicates that generating complex structured outputs remains a major challenge for generation models (Peng et al., 2023a), requiring further exploration.

### 4.3 Experiment Results of LLMs

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420 421

422

423

424

425

426

497

428

429

430

431

432

433

434

435

436

437

438

The results of LLMs with in-context learning are presented in Table 4, revealing that while LLMs with in-context learning are competitive compared to some fine-tuned EAE models, they still fall significantly short of the state-of-the-art. This is consistent with previous findings, suggesting that existing LLMs with in-context learning perform notably worse on specification-heavy information extraction tasks (Peng et al., 2023a; Li et al., 2023; Han et al., 2023). The LLMs' bag-of-words F1 scores are notably higher than their exact match scores, suggesting that the LLMs' predictions tend to be free-format and do not strictly match human annotations (Han et al., 2023).

One possible reason for the suboptimal performance is that LLMs cannot easily understand the schema form their names. Therefore, we conduct experiments with more informative prompts by incorporating definitions for each used argument role into the prompt, which are high-quality instructions used for guiding human annotators during data annotation. The results of these enhanced prompts are also shown in Table 4 (w/ definition). There is an obvious but marginal improvement after adding definitions, possibly due to the LLMs' limitations in understanding long contexts (Shaham et al., 2022; Peng et al., 2023a; Liu et al., 2023).

### 4.4 Analysis on Trigger-Argument Distance

As shown in Figure 3, MAVEN-ARG provides
document-level annotations, covering data with
varying trigger-argument distances. We conduct
an analytical experiment on the impact of triggerargument distance to model performance. Specifically, we break down the predictions and annotations in the test set by their trigger-argument dis-



Figure 4: Mention-level F1 (%) of models on data with varying trigger-argument distances, i.e., the number of words between an event argument and its trigger.

tances and evaluate how the performance changes along with different distances. The experimental results are shown in Figure 4, which demonstrate that models perform significantly worse on samples with longer trigger-argument distances. This aligns with previous findings in document-level relation extraction regarding the distance between entity pairs (Ru et al., 2021). It suggests that modeling long-distance dependencies between triggers and arguments remains a challenge for existing EAE models. Future research can leverage MAVEN-ARG to explore advanced methods for handling long-distance trigger-argument instances. 446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

# 4.5 Analysis on Entity and Non-Entity Arguments

MAVEN-ARG provides comprehensive annotations, including both entity and non-entity arguments. We analyze the performance breakdown of investigated EAE models on these two types of arguments. The results are presented in Table 5, which reveals that EAE models generally perform better on non-entity arguments. The possible reason may be that there are more non-entity arguments in MAVEN-ARG and non-entity arguments are often presented in a

Model	Mention Level				E	Entity Co	oref Leve	el	Event Coref Level			
	Р	R	F1	EM	Р	Ŕ	F1	EM	Р	R	F1	EM
GPT-3.5 w/ definition	$21.3 \\ 21.8$	$20.9 \\ 21.7$	$\begin{array}{c} 19.9 \\ 20.6 \end{array}$	$14.3 \\ 15.2$	$24.5 \\ 25.0$	$25.1 \\ 25.8$	$23.4 \\ 24.1$	$\begin{array}{c} 16.8 \\ 17.8 \end{array}$	$24.4 \\ 24.9$	$24.8 \\ 25.4$	$23.2 \\ 23.9$	$\begin{array}{c} 16.9 \\ 17.9 \end{array}$
GPT-4 w/ definition	25.6 <b>27.2</b>	27.2 <b>28.7</b>	25.1 <b>26.6</b>	17.9 <b>19.1</b>	28.9 <b>30.5</b>	31.7 <b>33.3</b>	28.7 <b>30.3</b>	20.2 <b>21.3</b>	27.9 <b>29.8</b>	30.5 <b>32.3</b>	27.6 <b>29.5</b>	19.5 <b>21.1</b>

Table 4: Experimental results (%) of LLMs with 2-shot in-context learning on MAVEN-ARG.

Model	En F1	tity EM	Non-l F1	Entity EM
DMBERT	19.7	19.5	-	_
CLEVE	22.1	22.0	_	_
BERT+CRF	17.8	18.5	19.4	24.0
EEQA	6.2	5.6	17.5	13.9
Text2Event	5.5	5.2	1.6	1.1
PAIE	20.3	19.2	37.6	30.4

Table 5: Mention-level results (%) of EAE models on entity and non-entity arguments. Classification-based models, e.g., DMBERT and CLEVE, are not applicable to non-entity arguments.

looser form, making it easier for the models to learn 470 the patterns and extract them. An exception is ob-471 served for the generation-based model Text2Event, 472 which exhibits poorer performance on non-entity 473 arguments. This may be because non-entity ar-474 guments are typically longer, which are harder to 475 generate at once. It suggests that further explo-476 ration is needed to investigate how to well handle 477 EAE with generation methods. 478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

# **5** Future Event Prediction Demonstration

MAVEN-ARG, in conjunction with MAVEN and MAVEN-ERE, creates the first all-in-one event understanding benchmark, which covers the full process of ED, EAE, and ERE. Beyond serving as an evaluation benchmark for these tasks, an allin-one event dataset naturally enables a variety of event-based applications, especially considering the recent advances brought by LLMs. Here we preliminarily explore an application case, future event prediction, as a demonstration.

Predicting future events based on causality can help decision-making, which is of self-evident importance. Therefore, since the early script learning (Schank and Abelson, 1975; Mooney and De-Jong, 1985), future event prediction has continually attracted research interest (Chambers and Jurafsky, 2008; Jans et al., 2012; Granroth-Wilding and Clark, 2016; Hu et al., 2017; Chaturvedi et al., 2017; Li et al., 2018b; Lee and Goldwasser, 2019;

Model	Reasonable (%)	Matched (%)
GPT-3.5	92.7	7.8
GPT-4	95.2	12.2

Table 6: Future event prediction results (%), averaged over 2 evaluators and 3 prompts. **Reasonable** denotes the rate of predictions judged as reasonable to happen next. **Matched** denotes the rate of predictions matched with the actual future events.

Zhao, 2021). However, due to the lack of highquality event resources, the evaluation of future event prediction often compromises by merely predicting verbs and subjects (Chambers et al., 2014), predicting according to textual order (Jans et al., 2012), or selecting story endings (Mostafazadeh et al., 2016a; Chaturvedi et al., 2017). The MAVEN series of datasets, with annotations of complete event structures and rich causal relations, may aid in predicting future events in real-world scenarios. 499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

**Experiment Setup** We sample 100 event chains, each consisting of 3-5 events, from the training and validation sets. In each chain, preceding events cause the subsequent ones. Events are described in a structured JSON format, containing event type, event trigger, and event arguments. For each event chain, we hold out the last event and input the remaining incomplete chain into two proprietary LLMs, GPT-3.5 and GPT-4 (OpenAI, 2023), requiring them to predict the next occurring event. These LLMs are prompted with detailed task instructions and 5 demonstration event chains. To minimize the influence of the demonstrations, predictions are made independently three times under different demonstrations. More experimental details are shown in appendix D. We employ manual evaluation, with two experts engaged to judge (1) whether the prediction is reasonable, and (2) whether the prediction matches the actual future event.

**Experimental Results** Experimental results are shown in Table 6. From these, we can see that the powerful LLMs can produce highly reasonable

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

581

event predictions. However, their predictions 531 seldom align with the actual future, making them 532 not directly helpful. These observations suggest that using LLMs for future event prediction is promising, but there remain topics to explore on how to build practical future event prediction systems with LLMs. For instance, using retrieval-537 augmented methods may help LLMs access more timely evidence when making future predictions. As a preliminary attempt, the experiments demon-540 strate how our all-in-one event understanding 541 dataset can assist in conveniently building and 542 evaluating event-based applications. We hope that 543 future works can explore using the MAVEN series 544 datasets to build diverse applications. 545

# 6 Related Work

**Event Argument Extraction Datasets** Since the 547 early MUC datasets (Grishman and Sundheim, 1996), event argument extraction (EAE) as a part of event extraction has received widespread attention. To reduce task complexity and provide stan-551 dardized extraction results, the ACE datasets (Dod-553 dington et al., 2004) are designed with a schema covering 33 event types, limiting event argument annotation to entities within the same sentence as the trigger. ACE 2005 (Walker et al., 2006) has been the most widely used dataset for a long time, 557 and the practice of ACE has been broadly adopted. Rich ERE (Song et al., 2015) expands ACE schema 559 to 38 event types and constructs the TAC KBP datasets (Ellis et al., 2014, 2015, 2016; Getman et al., 2017). MEE (Pouran Ben Veyseh et al., 2022) 563 follows the ACE schema to build a multilingual dataset. With the advancement of NLP methods, 564 some works break some of the constraints of ACE 565 task definition to construct more practical datasets. RAMS (Ebner et al., 2020), WikiEvents (Li et al., 567 2021), and DocEE (Tong et al., 2022) extends the 568 annotation scope to the whole documents. How-569 ever, RAMS and DocEE only annotate one topic event per document, ignoring fine-grained events within documents. MAVEN (Wang et al., 2020) and GENEVA (Parekh et al., 2023) both construct 573 high-coverage general event schemata with over 574 100 event types. MAVEN supports only event 576 detection. GENEVA extends event arguments to cover non-entity spans but focuses on testing 577 the generalizability rather than developing practical EAE methods. Its data are repurposed from FrameNet (Baker et al., 1998) examples, which are 580

individual sentences without document structure. MAVEN-ARG meticulously designs 612 unique argument roles for MAVEN schema and conducts large-scale exhaustive annotation, which annotates both entity and non-entity arguments for finegrained events at the document level.

**Event Argument Extraction Methods** Traditional EAE methods primarily involve (1) Classification-based methods (Chen et al., 2015a, 2017; Sha et al., 2018; Wadden et al., 2019; Wang et al., 2019; Lin et al., 2020; Wang et al., 2021b; Zhou and Mao, 2022): employing text encoders like CNN (Krizhevsky et al., 2012) and BERT (Devlin et al., 2019), followed by an information aggregator, such as dynamic multipooling mechanism (Chen et al., 2015a), to obtain role-specific representations for classification. (2) Sequence labeling methods (Nguyen et al., 2016; Yang and Mitchell, 2017; Nguyen et al., 2021; Peng et al., 2023c): mainly adopting the conditional random field (CRF) (Lafferty et al., 2001) as the output layer to model structured dependencies between different arguments. Recently, increasing attention has been paid to transforming EAE into a question-answering task, transferring question-answering capabilities to boost EAE (Liu et al., 2020; Du and Cardie, 2020; Li et al., 2020; Ma et al., 2022; Lu et al., 2023). Additionally, some research focuses on using generation models to directly generate structured outputs containing events and their arguments (Lu et al., 2021; Li et al., 2021; Hsu et al., 2022; Lu et al., 2022; Ren et al., 2023; Hsu et al., 2023; Zhang et al., 2023), which has been becoming increasingly important with the advance of large language models.

# 7 Conclusion and Future Work

We introduce MAVEN-ARG, an event argument extraction dataset with comprehensive schema, large data scale, and exhaustive annotation. Experiments indicate that MAVEN-ARG is quite challenging for both fine-tuned EAE models and proprietary large language models. Together with MAVEN and MAVEN-ERE, MAVEN-ARG completes an all-inone dataset covering the entire process of event understanding. An application case of future event prediction demonstrates how an all-in-one dataset can enable broad event-based applications. In the future, we will explore constructing multilingual resources under this framework and developing practical EAE methods with MAVEN-ARG.

#### Limitations 631

(1) MAVEN-ARG currently includes only English 632 corpus, which limits its potential applications and coverage for diverse linguistic phenomena. In 634 future work, we will try to support more languages 636 under our framework and we also encourage community efforts in developing multilingual event understanding benchmarks. (2) MAVEN-ARG, along with MAVEN (Wang et al., 2020) and MAVEN-ERE (Wang et al., 2022), exclusively supports mainstream event understanding tasks. 641 However, these datasets do not cover more 642 broad event-related tasks such as event factuality identification (Qian et al., 2019, 2022) and event salience identification (Liu et al., 2018). We 645 encourage future explorations in building more challenging and diverse tasks and applications 647 on top of MAVEN data. (3) While previous research has found that LLMs perform poorly on specification-heavy tasks (Peng et al., 2023c; Han et al., 2023; Li et al., 2023) including the EAE task, there is no in-depth exploration of effective LLM-based approaches addressing the EAE task 653 in this paper. We leave the exploration of how to better leverage LLMs for EAE tasks in future work.

# **Ethical Considerations**

657

671

672

674

675

In this section, we discuss the ethical considerations of this work: (1) Intellectual property. The MAVEN dataset is released under the CC BY-SA 4.0 license<sup>3</sup>. The MAVEN-ERE is shared under GPLv3<sup>4</sup> license and the original Wikipedia corpus is shared under the CC BY-SA 3.0 license<sup>5</sup>. The usage of these data in this work strictly adheres to the corresponding licenses and intended use. (2) Intended use. MAVEN-ARG is an event argument extraction dataset. Researchers and practitioners can utilize MAVEN-ARG to train and evaluate models for event argument extraction, thereby advancing the field of event understanding. (3) Potential risk control. MAVEN-ARG is constructed based on publicly available data. We believe that the underlying public data has been adequately desensitized and anonymized. The event argument annotation does not involve judgments about social issues and thus we believe MAVEN-ARG will not involve additional risks. The MAVEN-ARG test set will not 676 be publicly released to prevent potential unfair performance comparisons. (4) Worker Treatments are discussed in appendix A.2. (5) AI assistant. 679 The writing of this paper was assisted by ChatGPT 680 in rephrasing some sentences. 681

677

678

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

# References

- Jacqueline Aguilar, Charley Beller, Paul McNamee, Benjamin Van Durme, Stephanie Strassel, Zhiyi Song, and Joe Ellis. 2014. A comparison of the events and relations across ACE, ERE, TAC-KBP, and FrameNet annotation standards. In Proceedings of the Second Workshop on EVENTS: Definition, Detection, Coreference, and Representation, pages 45-53.
- Samuel Joseph Amouyal, Ohad Rubin, Ori Yoran, Tomer Wolfson, Jonathan Herzig, and Jonathan Berant. 2022. QAMPARI: : An open-domain question answering benchmark for questions with many answers from multiple paragraphs. CoRR, abs/2205.12665.
- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. The Berkeley FrameNet project. In Proceedings of ACL-COLING, pages 86–90.
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract Meaning Representation for sembanking. In Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse, pages 178-186.
- Tommaso Caselli and Piek Vossen. 2017. The event StoryLine corpus: A new benchmark for causal and temporal relation extraction. In Proceedings of the Events and Stories in the News Workshop, pages 77-86.
- Nathanael Chambers, Taylor Cassidy, Bill McDowell, and Steven Bethard. 2014. Dense event ordering with a multi-pass architecture. Transactions of the Association for Computational Linguistics, 2:273– 284.
- Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised learning of narrative event chains. In Proceedings of ACL-HLT, pages 789–797.
- Snigdha Chaturvedi, Haoruo Peng, and Dan Roth. 2017. Story comprehension for predicting what happens next. In Proceedings of EMNLP, pages 1603–1614.
- Yubo Chen, Shulin Liu, Xiang Zhang, Kang Liu, and Jun Zhao. 2017. Automatically Labeled Data Generation for Large Scale Event Extraction. In Proceedings of ACL, pages 409-419.

<sup>&</sup>lt;sup>3</sup>https://creativecommons.org/licenses/by-sa/4. 0/

<sup>&</sup>lt;sup>4</sup>https://www.gnu.org/licenses/gpl-3.0.html

<sup>&</sup>lt;sup>5</sup>https://creativecommons.org/licenses/by-sa/3. 0/

727

- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015a. Event extraction via dynamic multipooling convolutional neural networks. In Proceedings of ACL-IJCNLP, pages 167-176.
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015b. Event extraction via dynamic multipooling convolutional neural networks. In Proceedings of ACL, pages 167-176.
- Agata Cybulska and Piek Vossen. 2014. Using a sledgehammer to crack a nut? lexical diversity and event coreference resolution. In Proceedings of LREC, pages 4545-4552.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT, pages 4171-4186.
- Ning Ding, Guangwei Xu, Yulin Chen, Xiaobin Wang, Xu Han, Pengjun Xie, Haitao Zheng, and Zhiyuan Liu. 2021. Few-NERD: A few-shot named entity recognition dataset. In Proceedings of ACL, pages 3198-3213.
- Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2015. Deep learning for event-driven stock prediction. In Proceedings of IJCAI.
- 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 LREC.
- Xinya Du and Claire Cardie. 2020. Event extraction by answering (almost) natural questions. In Proceedings of EMNLP, pages 671-683.
- Seth Ebner, Patrick Xia, Ryan Culkin, Kyle Rawlins, and Benjamin Van Durme. 2020. Multi-sentence argument linking. In Proceedings of ACL, pages 8057-8077.
- Joe Ellis, Jeremy Getman, Dana Fore, Neil Kuster, Zhiyi Song, Ann Bies, and Stephanie M Strassel. 2015. Overview of linguistic resources for the TAC KBP 2015 evaluations: Methodologies and results. In TAC.
- Joe Ellis, Jeremy Getman, Dana Fore, Neil Kuster, Zhiyi Song, Ann Bies, and Stephanie M Strassel. 2016. Overview of Linguistic Resources for the TAC KBP 2016 Evaluations: Methodologies and Results. In TAC.
- Joe Ellis, Jeremy Getman, and Stephanie M Strassel. 2014. Overview of linguistic resources for the TAC KBP 2014 evaluations: Planning, execution, and results. In TAC.
- Charles J Fillmore. 1976. Frame semantics and the nature of language. In Annals of the New York Academy of Sciences: Conference on the origin and development of language and speech, volume 280, pages 20-32.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. Psychological bulletin, 76(5):378.

783

784

785

786

787

788

789

790

791

792

793

794

795

797

798

799

800

801

802

803

804

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

- Jeremy Getman, Joe Ellis, Zhiyi Song, Jennifer Tracey, and Stephanie Strassel. 2017. Overview of linguistic resources for the tac kbp 2017 evaluations: Methodologies and results. In TAC.
- Goran Glavaš and Jan Šnajder. 2014. Event graphs for information retrieval and multi-document summarization. Expert systems with applications, 41(15):6904– 6916.
- Seraphina Goldfarb-Tarrant, Haining Feng, and Nanyun Peng. 2019. Plan, write, and revise: an interactive system for open-domain story generation. In Proceedings of NAACL: Demonstrations, pages 89-97.
- Mark Granroth-Wilding and Stephen Clark. 2016. What happens next? event prediction using a compositional neural network model. In Proceedings of the AAAI, volume 30.
- Ralph Grishman and Beth Sundheim. 1996. Message Understanding Conference- 6: A brief history. In Proceedings of COLING.
- Felix Hamborg, Corinna Breitinger, and Bela Gipp. 2019. Giveme5w1h: A universal system for extracting main events from news articles. In Proceedings of INRA@RecSys, CEUR Workshop Proceedings, pages 35-43.
- Ridong Han, Tao Peng, Chaohao Yang, Benyou Wang, Lu Liu, and Xiang Wan. 2023. Is information extraction solved by ChatGPT? An analysis of performance, evaluation criteria, robustness and errors. arXiv preprint arXiv:2305.14450.
- Eduard Hovy, Teruko Mitamura, Felisa Verdejo, Jun Araki, and Andrew Philpot. 2013. Events are not simple: Identity, non-identity, and quasi-identity. In Proceedings of Workshop on Events: Definition, Detection, Coreference, and Representation, pages 21-28.
- I-Hung Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott Miller, Prem Natarajan, Kai-Wei Chang, and Nanyun Peng. 2022. Degree: A data-efficient generation-based event extraction model. In Proceedings of NAACL-HLT, pages 1890–1908.
- I-Hung Hsu, Zhiyu Xie, Kuan-Hao Huang, Prem Natarajan, and Nanyun Peng. 2023. AMPERE: AMR-aware prefix for generation-based event argument extraction model. In Proceedings of ACL, pages 10976–10993.
- Linmei Hu, Juanzi Li, Liqiang Nie, Xiao-Li Li, and Chao Shao. 2017. What happens next? future subevent prediction using contextual hierarchical lstm. In Proceedings of AAAI, volume 31.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and 834 Yejin Choi. 2019. Cosmos QA: Machine reading 835

840	
841	
842	
843	(
844	
845	
846	
847	
848	
849	
850	
000	
851	
852	
853	
957	
054	
000	
856	1
857	Ĩ
959	
0.00	
859	1
860	
961	
000	
862	
863	,
267	
004	
865	1
266	1
000	
007	
808	
869	
870	
071	1
871	1
872	
873	
874	
875	,
876	
877	
070	,
878	4
879	
088	
881	
~~~	,
882	4
883	
884	
885	
886	
997	

837

839

comprehension with contextual commonsense reasoning. In *Proceedings of EMNLP-IJCNLP*, pages 2391–2401.

- Bram Jans, Steven Bethard, Ivan Vulic, and Marie-Francine Moens. 2012. Skip n-grams and ranking functions for predicting script events. In *Proceedings EACL*, pages 336–344.
- Çagla Çig Karaman, Serkan Yaliman, and Salih Atilay Oto. 2017. Event detection from social media: 5w1h analysis on big data. In *Proceedings of SIU*, pages 1–4. IEEE.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Proceedings of NeurIPs*, pages 1106–1114.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML*, pages 282– 289.
- I-Ta Lee and Dan Goldwasser. 2019. Multi-relational script learning for discourse relations. In *Proceedings of ACL*, pages 4214–4226.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of EMNLP*, pages 3045– 3059.
- Willem JM Levelt. 1993. Speaking: From intention to articulation. MIT press.
- Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023. Evaluating ChatGPT's information extraction capabilities: An assessment of performance, explainability, calibration, and faithfulness. *arXiv preprint arXiv:2304.11633*.
- 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 EMNLP*, pages 829–838.
- Sha Li, Heng Ji, and Jiawei Han. 2021. Document-level event argument extraction by conditional generation. In *Proceedings of NAACL-HLT*, pages 894–908.
- Zhongyang Li, Xiao Ding, and Ting Liu. 2018a. Constructing narrative event evolutionary graph for script event prediction. In *Proceedings of IJCAI*, pages 4201–4207.
- Zhongyang Li, Xiao Ding, and Ting Liu. 2018b. Constructing narrative event evolutionary graph for script event prediction. In *Proceedings of IJCAI*.
  - Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. A joint neural model for information extraction with global features. In *Proceedings of ACL*, pages 7999– 8009.

Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. Event extraction as machine reading comprehension. In *Proceedings of EMNLP*, pages 1641–1651. 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

- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*.
- Zhengzhong Liu, Chenyan Xiong, Teruko Mitamura, and Eduard Hovy. 2018. Automatic event salience identification. In *Proceedings of EMNLP*, pages 1226–1236.
- Di Lu, Shihao Ran, Joel Tetreault, and Alejandro Jaimes. 2023. Event extraction as question generation and answering. In *Proceedings of ACL*, pages 1666–1688.
- 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 ACL-IJCNLP*, pages 2795–2806.
- Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu Lin, Xianpei Han, Le Sun, and Hua Wu. 2022. Unified structure generation for universal information extraction. In *Proceedings of ACL*, pages 5755–5772.
- Mingyu Derek Ma, Jiao Sun, Mu Yang, Kung-Hsiang Huang, Nuan Wen, Shikhar Singh, Rujun Han, and Nanyun Peng. 2021. EventPlus: A temporal event understanding pipeline. In *Proceedings of NAACL: Demonstrations*, pages 56–65.
- 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. In *Proceedings of ACL*, pages 6759–6774.
- George A Miller and Philip N Johnson-Laird. 2013. Language and perception. In *Language and Perception*. Harvard University Press.
- Paramita Mirza, Rachele Sprugnoli, Sara Tonelli, and Manuela Speranza. 2014. Annotating causality in the TempEval-3 corpus. In *Proceedings of the EACL* 2014 Workshop on Computational Approaches to Causality in Language (CAtoCL), pages 10–19.
- Raymond J Mooney and Gerald DeJong. 1985. Learning schemata for natural language processing. In *Proceedings of IJCAI*, pages 681–687.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016a. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of NAACL-HLT*, pages 839–849.

Nasrin Mostafazadeh, Alyson Grealish, Nathanael Zhong Qian, Heng Zhang, Peifeng Li, Qiaoming Zhu, 992 and Guodong Zhou. 2022. Document-level event Chambers, James Allen, and Lucy Vanderwende. 993 CaTeRS: Causal and temporal relation factuality identification via machine reading com-2016b. 994 scheme for semantic annotation of event structures. prehension frameworks with transfer learning. In 995 In Proceedings of the Fourth Workshop on Events, Proceedings of COLING, pages 2622–2632. 996 pages 51-61. Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and 997 Minh Van Nguyen, Tuan Ngo Nguyen, Bonan Min, Percy Liang. 2016. SQuAD: 100,000+ questions for 998 and Thien Huu Nguyen. 2021. Crosslingual transfer machine comprehension of text. In Proceedings of 999 learning for relation and event extraction via word EMNLP, pages 2383-2392. 1000 category and class alignments. In Proceedings of EMNLP, pages 5414-5426. Yubing Ren, Yanan Cao, Ping Guo, Fang Fang, Wei 1001 Ma, and Zheng Lin. 2023. Retrieve-and-sample: 1002 Thien Huu Nguyen, Kyunghyun Cho, and Ralph Gr-Document-level event argument extraction via hy-1003 ishman. 2016. Joint event extraction via recurrent brid retrieval augmentation. In Proceedings of ACL, 1004 neural networks. In Proceedings of NAACL-HLT, pages 293-306. 1005 pages 300-309. Dongyu Ru, Changzhi Sun, Jiangtao Feng, Lin Qiu, 1006 Qiang Ning, Hao Wu, and Dan Roth. 2018. A multi-Hao Zhou, Weinan Zhang, Yong Yu, and Lei Li. axis annotation scheme for event temporal relations. 2021. Learning logic rules for document-level re-1008 In Proceedings of ACL, pages 1318–1328. lation extraction. In Proceedings of EMNLP, pages 1009 1239-1250. 1010 OpenAI. 2022. Introducing ChatGPT. Roger C Schank and Robert P Abelson. 1975. Scripts, 1011 OpenAI. 2023. GPT-4 technical report. arXiv preprint plans, and knowledge. In Proceedings of IJCAI, vol-1012 arXiv:2303.08774. 1013 ume 75, pages 151–157. Tanmay Parekh, I-Hung Hsu, Kuan-Hao Huang, Kai-Lei Sha, Feng Qian, Baobao Chang, and Zhifang Sui. 1014 Wei Chang, and Nanyun Peng. 2023. GENEVA: 2018. Jointly extracting event triggers and arguments 1015 1016 Benchmarking generalizability for event argument by dependency-bridge RNN and tensor-based argument interaction. In Proceedings of AAAI, pages extraction with hundreds of event types and argument 1017 roles. In Proceedings of ACL, pages 3664-3686. 5916-5923. 1018 Hao Peng, Xiaozhi Wang, Jianhui Chen, Weikai Li, Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori 1019 Yunjia Qi, Zimu Wang, Zhili Wu, Kaisheng Zeng, Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong, 1020 Bin Xu, Lei Hou, et al. 2023a. When does in-context Mor Geva, Jonathan Berant, and Omer Levy. 2022. 1021 learning fall short and why? a study on specification-SCROLLS: Standardized CompaRison over long lan-1022 heavy tasks. arXiv preprint arXiv:2311.08993. guage sequences. In Proceedings of EMNLP, pages 1023 12007-12021. 1024 Hao Peng, Xiaozhi Wang, Feng Yao, Zimu Wang, Chuzhao Zhu, Kaisheng Zeng, Lei Hou, and Juanzi Zhiyi Song, Ann Bies, Stephanie Strassel, Tom Riese, 1025 Li. 2023b. Omnievent: A comprehensive, fair, and Justin Mott, Joe Ellis, Jonathan Wright, Seth Kulick, 1026 easy-to-use toolkit for event understanding. pages Neville Ryant, and Xiaoyi Ma. 2015. From light 1027 508-517. to rich ERE: Annotation of entities, relations, and 1028 events. In Proceedings of the The 3rd Workshop on 1029 Hao Peng, Xiaozhi Wang, Feng Yao, Kaisheng Zeng, EVENTS: Definition, Detection, Coreference, and 1030 Lei Hou, Juanzi Li, Zhiyuan Liu, and Weixing Shen. Representation, pages 89–98. 1031 2023c. The devil is in the details: On the pitfalls of event extraction evaluation. In Findings of ACL. MeiHan Tong, Bin Xu, Shuai Wang, Meihuan Han, 1032 Yixin Cao, Jiangqi Zhu, Siyu Chen, Lei Hou, and 1033 Steven Pinker. 2013. Learnability and Cognition, new Juanzi Li. 2022. DocEE: A large-scale and fine-1034 edition: The Acquisition of Argument Structure. MIT grained benchmark for document-level event extraction. In Proceedings of NAACL-HLT, pages 3970press. 1036 3982. 1037 Amir Pouran Ben Veyseh, Javid Ebrahimi, Franck Dernoncourt, and Thien Nguyen. 2022. MEE: A novel David Wadden, Ulme Wennberg, Yi Luan, and Han-1038 naneh Hajishirzi. 2019. Entity, relation, and event exmultilingual event extraction dataset. In Proceedings 1039 of the EMNLP, pages 9603-9613. traction with contextualized span representations. In 1040 Proceedings of EMNLP-IJCNLP, pages 5783-5788. 1041 Zhong Qian, Peifeng Li, Qiaoming Zhu, and Guodong Zhou. 2019. Document-level event factuality identifi-Christopher Walker, Stephanie Strassel, Julie Medero, 1042 cation via adversarial neural network. In Proceedings and Kazuaki Maeda. 2006. ACE 2005 multilingual 1043 of NAACL-HLT, pages 2799-2809. training corpus. Linguistic Data Consortium, 57. 1044 12

941

942

947

951

955

957

958

962

965

966

967

968

969

970

971

972 973

974

975

977

985

Shichao Wang, Xiangrui Cai, Hongbin Wang, and Xiaojie Yuan. 2021a. Incorporating circumstances into narrative event prediction. In *Findings of EMNLP*, pages 4840–4849.

1045 1046

1047 1048

1049

1052

1054

1057 1058

1059

1060

1061 1062

1063

1064

1065

1066

1067 1068

1069

1070

1071

1072

1074

1075

1076 1077

1078

1079

1082

1084

1085

1088

1090

1091

1092

1093

1094

1095

1096 1097

1098

- Xiaozhi Wang, Yulin Chen, Ning Ding, Hao Peng, Zimu Wang, Yankai Lin, Xu Han, Lei Hou, Juanzi Li, Zhiyuan Liu, et al. 2022. MAVEN-ERE: A unified large-scale dataset for event coreference, temporal, causal, and subevent relation extraction. In *Proceedings of EMNLP*, pages 926–941.
- Xiaozhi Wang, Ziqi Wang, Xu Han, Wangyi Jiang, Rong Han, Zhiyuan Liu, Juanzi Li, Peng Li, Yankai Lin, and Jie Zhou. 2020. MAVEN: A Massive General Domain Event Detection Dataset. In Proceedings of EMNLP, pages 1652–1671.
- Xiaozhi Wang, Ziqi Wang, Xu Han, Zhiyuan Liu, Juanzi Li, Peng Li, Maosong Sun, Jie Zhou, and Xiang Ren. 2019. HMEAE: Hierarchical modular event argument extraction. In *Proceedings of EMNLP-IJCNLP*, pages 5777–5783.
- Ziqi Wang, Xiaozhi Wang, Xu Han, Yankai Lin, Lei Hou, Zhiyuan Liu, Peng Li, Juanzi Li, and Jie Zhou.
  2021b. CLEVE: Contrastive Pre-training for Event Extraction. In *Proceedings of ACL-IJCNLP*, pages 6283–6297.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of EMNLP: System Demonstrations*, pages 38–45.
- Bishan Yang and Tom M. Mitchell. 2017. Leveraging knowledge bases in lstms for improving machine reading. In *Proceedings of ACL*, pages 1436–1446.
- Zijun Yao, Yantao Liu, Xin Lv, Shulin Cao, Jifan Yu, Juanzi Li, and Lei Hou. 2023. KoRC: Knowledge oriented reading comprehension benchmark for deep text understanding. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11689– 11707.
- Kaihang Zhang, Kai Shuang, Xinyue Yang, Xuyang Yao, and Jinyu Guo. 2023. What is overlap knowledge in event argument extraction? APE: A crossdatasets transfer learning model for EAE. In *Proceedings of ACL*, pages 393–409.
- Liang Zhao. 2021. Event prediction in the big data era: A systematic survey. *ACM Comput. Surv.*, 54(5).
- Hanzhang Zhou and Kezhi Mao. 2022. Document-level event argument extraction by leveraging redundant information and closed boundary loss. In *Proceedings* of NAACL-HLT, pages 3041–3052.

#### Appendices 1100

1101

1102

1103

1104

1107

1111

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1135

1137

1141 1142

1143

1144

1145

#### **Data Collection Details** A

#### Annotation Instruction A.1

As introduced in § 2.1, we create a detailed event schema for both defining the task and instructing the annotators. We present the annotation instruc-1105 tions for the event type Incident in Table 7, includ-1106 ing its argument schema and annotation examples. The overall event schema will be released along 1108 with the dataset. To support the highly customized 1109 annotation process designed for us, we developed 1110 a new online annotation platform. A screenshot for the annotation platform is shown in Figure 5 to 1112 help understand the annotation operations.

# A.2 Annotation Coordination

We employ annotators (including senior annotators and expert annotators) from multiple commercial data annotation companies. 61% of them are female and 39% of them are male. The experts involved in schema creation are invited by the authors through personal connections. All the workers are fairly paid with agreed salaries and workloads. All employment is under contract and in compliance with local regulations. The overall annotation cost, including annotating entities and event arguments as well as developing and maintaining annotation platforms, is about 85,000 USD.

#### B Data Split Statistics

The detailed statistics of different data splits of MAVEN-ARG are shown in Table 8.

#### **EAE Experimental Details** С

# C.1 Fine-tuning Implementation Details

Here we provide brief descriptions of the finetuning-based models involved in our experiments. (1) DMBERT (Wang et al., 2019) utilizes 1134 BERT (Devlin et al., 2019) as the text encoder and a dynamic multi-pooling mechanism (Chen et al., 1136 2015b) on top of BERT to aggregate argumentspecific features and map them onto the distribution 1138 in the label space. (2) CLEVE (Wang et al., 2021b) 1139 is an event-oriented pre-trained language model, 1140 which is pre-trained using contrastive pre-training objectives on large-scale unsupervised data and their semantic structures. (3) BERT+CRF (Wang et al., 2020) is a sequence labeling model, which leverages BERT as the backbone and the conditional random field (CRF) (Lafferty et al., 2001) as 1146

the output layer to model the structural dependen-1147 cies of predictions. (4) EEQA (Li et al., 2020) is 1148 a span prediction model, which formulates event 1149 extraction as a question-answering task and out-1150 puts start and end positions to indicate triggers and 1151 arguments. (5) Text2Event (Lu et al., 2021) is 1152 a conditional generation model, which proposes 1153 a sequence-to-structure paradigm and generates 1154 structured outputs containing triggers and corre-1155 sponding arguments with constrained decoding. 1156 (6) PAIE (Ma et al., 2022) adopts prompt tun-1157 ing (Lester et al., 2021) to train two span selec-1158 tors for each argument role in the provided prompt 1159 and conduct joint optimization to find optimal role-1160 span assignments. We adopt the same backbones 1161 with their original papers for all EAE models in 1162 our experiments. We employ pipeline evaluation as 1163 suggested by Peng et al. (2023c). Specifically, for 1164 PAIE, we conduct EAE experiments based on the 1165 triggers predicted by CLEVE. For the other mod-1166 els, the EAE experiments are based on the triggers 1167 extracted by corresponding models. 1168

We implement the EAE models using code from the official repositories of OmniEvent (Peng et al., 2023b), PAIE (Ma et al., 2022), and Text2Event (Lu et al., 2021). The numbers of parameters of the EAE models are shown in Table 3. All open-source models are downloaded from the HuggingFace Transformers community (Wolf et al., 2020). Each of our fine-tuning experiments is conducted only once, on Nvidia A100 GPUs, consuming approximately 800 GPU hours in total. The hyper-parameters of the model are set based on prior experience and references from previous papers (Lu et al., 2021; Ma et al., 2022; Peng et al., 2023b). All hyper-parameters are shown in Table 9.

# C.2 LLM Experimental Details

1183

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

We access ChatGPT and GPT-4 through the offi-1184 cial OpenAI interfaces, namely gpt-3.5-turbo 1185 and gpt-4, respectively. The API access period 1186 spans from October 1 to October 31, 2023. The 1187 decoding sampling temperature for both models is 1188 set to 0. An example of the prompt, input, output, 1189 and ground-truth of this experiment are presented 1190 in Table 10. Model outputs are automatically ex-1191 tracted and evaluated using the evaluation approach 1192 mentioned in  $\S$  4.1. 1193

### [Incident] Accident, unfortunate event

# **Event Arguments:**

1. Participant : Entities involved in the accident (individuals, institutions, organizations, and even trains, ships, etc.). They can be the ones causing the accident or the ones affected by it. Similar to the combination of Agent and Patient in previous events, but due to the difficulty in distinguishing between Agent and Patient in accidents, they are uniformly labeled as Participants.

2. Location : The location or position where the incident occurred. If the incident involves multiple locations during the process, they should be marked separately.

3. Content : In general, only one annotation is needed, which accurately indicates the content and type of the accident.

4. Loss : The losses caused by accidents can include the number of deaths and injuries, property damage, and so on.

**Annotation Examples:** 

1. British losses were confined to a single man wounded by an accident aboard "Crescent".

2. On 6 June 1982, during the Falkland's war, the British Royal Navy type 42 destroyer engaged and destroyed a [British army

gazelle helicopter, serial number "XX377"] Participant + Loss , in a friendly fire incident, killing all four occupants

Table 7: Example annotation instructions for event type Incident. Different argument roles are denoted by different background colors. **Triggers** are bolded in red.

In the morning Izuku was quick to realize it was going to be a busy day .	Span Selection
	Trigger headed
There wasn 't anyone that told him that , nor was there a particular busy schedule .	Event Type Motion Movement Creek Details
As Shota and he headed to U . A to what was supposed to be a regular day , all around Izuku could hear about the Hosu incident .	Trigger mismatchs event type
Talks about Yamiyo were all around .	Agent     Izuku ×     Shota ×  The subject of overall displacement. It can be individuals, armles, or transportation vehicles such as airplanes and ca
On a big billboard was going on on a a live program going on about him .	2 Location_original The startice point of Apart's prelifion movement
Izuku watched it as they stood at a red light .	3 Location_final U.A×
Izuku continued to watch it as the car took off .	he end point of the location movement of the Agent, can be one or more
A vibration of his phone was what made him blink and tear his eyes away from it []	
A text from Kenta.	
That was still his contact number , though $\boxed{\text{tzuku}}$ briefly wondered to change it to $\boxed{\text{Habiki}}$ , he decided against it .	
After all , he himself told izuku hw was fine with being called either .	
> So , what have you been up to ?	
What exactly went on yesterday and the day before that ?' I got arrested	

Figure 5: Screenshot for the annotation platform. The trigger "headed" is selected for annotation (in the right panel) and entities are highlighted in green as the options for annotating event arguments.

Dataset	#Doc.	#Event	#Trigger	#Arg.	#Entity	#Entity Mention	Entity Arg.	Non-Entity Arg.
Train	2,913	64,923	70,775	190,479	55,421	86,969	76,882	113, 597
Dev	710	15,556	16,996	46,458	12,927	18,806	18,040	28,418
Test	857	18, 112	19,736	53,676	15,297	23,351	21,102	32,574

Table 8: Statistics of the data splits of MAVEN-ARG. "Doc." is short for "Document" and "Arg." is short for "Argument".

	DMBERT	CLEVE	BERT+CRF	EEQA	PAIE	Text2Event
Learning Rate	$5 \times 10^{-5}$	$1 \times 10^{-5}$	$5 \times 10^{-5}$	$5 \times 10^{-5}$	$2 \times 10^{-5}$	$5 \times 10^{-5}$
Weight Decay	$1 \times 10^{-5}$	$1 \times 10^{-2}$				
Batch Size	32	128	64	32	16	8
Epoch	6	5	10	10	_	30
Warmup Rate	0.1	0.1	0.1	0.1	0	0.1

Table 9: Hyper-parameters of fine-tuning EAE models on MAVEN-ARG. PAIE utilizes 10,000 gradient update steps to optimize the parameters.

**PROMPT:** Please extract event argument roles and the corresponding mentions for the events marked with <event> and </event> in the text, the possible roles must be chosen from the Roleset. If there are no roles for the event, place output NA.

**INPUT:** The February 2010 Australian **<event>** cyberattacks **</event>** were a series of denial-of-service attacks conducted by the Anonymous online community against the Australian government in response to proposed web censorship regulations . Operation Titstorm was the name given to the cyber attacks by the perpetrators . They resulted in lapses of access to government websites on 10 and 11 February 2010 . This was accompanied by emails , faxes , and phone calls harassing government offices . The actual size of the attack and number of perpetrators involved is unknown but it was estimated that the number of systems involved ranged from the hundreds to the thousands . The amount of traffic caused disruption on multiple government websites . Australian Telecommunications Minister Stephen Conroy proposed the regulations that would mainly filter sites with pornographic content . Various groups advocating uncensored access to the Internet , along with companies like Google and Yahoo ! , object to the proposed filter . A spokesperson for Conroy said that the actions were not a legitimate form of protest and called it irresponsible . The attacks also drew criticism from other filter protest groups . The initial stage was followed by small in-person protests on 20 February that were called " Project Freeweb " . Roleset: [Agent, Location, Patient, Consequence]. The event trigger is: **<event>** cyberattacks **</event>**.

**GPT-3.5 OUTPUT:** Answer: Patient: the Australian government; Consequence: lapses of access to government websites; Agent: the Anonymous online community

**GPT-4 OUTPUT:** Answer: Agent: the Anonymous online community; Location: Australia; Patient: Australian government; Consequence: lapses of access to government websites, disruption on multiple government websites.

**GROUND-TRUTH:** Agent: Anonymous online community; Patient: Australian government; Consequence: lapses of access to government websites on 10 and 11 February 2010; Location: Australian

Table 10: Example prompt, input, output, and ground-truth for the LLMs experiments. Demonstrations are omitted.

# **D** Event Prediction Experimental Details

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1205

1208

1210

1211

1212

The future event prediction experiments (§ 5) were conducted in October and November, 2023. We use OpenAI API endpoints gpt-3.5-turbo and gpt-4 for GPT-3.5 and GPT-4 experiments, specifically. To ensure the consistency among different runs, we set temperature=0.0. Detailed instructions and example input and output are shown in Table 11.

# **E** More Experimental Results

In this section, we present more experimental results of using different proportions of training data for training (appendix E.1) and results on entity and non-entity arguments (appendix E.2).

### E.1 Analysis on Data Size

The data volume of MAVEN-ARG significantly exceeds that of commonly used datasets. To examine the benefits of increased data scale, we train models on training data of varying sizes and observe their performance on the original test set. The ex-



Figure 6: Mention-level F1 scores of investigated EAE models using different proportions of training data on MAVEN-ARG. This experiment adopts the gold trigger evaluation approach (Peng et al., 2023c) and hence Text2Event is not applicable. Due to the computation limitations, CLEVE is not evaluated in this experiment.

perimental results are shown in Figure 6, which 1213 demonstrates that more training data indeed en**PROMPT:** You need to complete the event chain prediction task. I will provide you with an event chain consisting of multiple events, where earlier events cause the occurrence of subsequent events. You need to predict what the next most likely event in this chain will be. Events will be provided in a structured JSON format, including event type, core event trigger words, and event elements. I will first give you 5 example chains: ... Now please predict the next event for the following event chain:

INPUT: Event No.1: {"event type": "Attack", "event trigger": ["attack"], "event argument": {"Agent": ["unknown assailant"], "Patient": ["Boris Nemtsov"], "Location": ["Bolshoy Moskvoretsky Bridge"]}} Event No.2: {"event type": "Use\_firearm", "event trigger": ["fired"], "event argument": {"Agent": ["unknown assailant"], "Patient": ["Nemtsov"], "Location": ["Bolshoy Moskvoretsky Bridge"]}} Event type": "Bolshoy Moskvoretsky Bridge"]} Event type": "Use\_firearm", "event trigger": ["fired"], "event argument": {"Agent": ["unknown assailant"], "Patient": ["Nemtsov"], "Location": ["Bolshoy Moskvoretsky Bridge"]}} Event type": "Bolshoy Moskvoretsky Bridge"]}}

**GPT-3.5 OUTPUT:** Based on the provided event chains, the next most likely event in the given event chain is: Event No.4: {"event type": "Criminal\_investigation", "event trigger": ["indicted"], "event argument": {"Agent": ["unknown assailant"], "Suspect": ["Boris Nemtsov"], "Incident": ["shooting on Bolshoy Moskvoretsky Bridge"], "Conclusion": ["investigation and charges"]}}

**GPT-4 OUTPUT:** Event No.4: {"event type": "Death", "event trigger": ["died"], "event argument": {"Agent": ["Boris Nemtsov"], "Cause": ["gunshot wounds"], "Location": ["Bolshoy Moskvoretsky Bridge"]}}

**GROUND-TRUTH:** Event No.4: {"event type": "Death", "event trigger": ["died"], "event argument": {"Agent": ["Boris Nemtsov"], "Location": ["Bolshoy Moskvoretsky Bridge"]}}

Table 11: Example prompt, input, output, and ground-truth for the future event prediction experiments. Demonstrations are omitted and the JSON strings are unformatted to avoid taking up to much space.

hances model performance and allows for a com-1215 prehensive comparison of different models. The 1216 extensive data of MAVEN-ARG make it feasible 1217 to train a large language model (LLM) for gen-1218 eral event understanding, which we leave as future 1219 work. Table 12 shows the detailed experimental 1220 results, i.e., mention level, entity coreference level, 1221 and event coreference level. 1222

# E.2 Entity and Non-Entity Arguments

1223

1224

1225

1226

1227

1228

Table 13 presents the overall results on entity and non-entity arguments of MAVEN-ARG. The nonentity arguments do not have coreferential relationship with each other and hence there is no entity coreference level evaluation for them.

		Mentio	n Level		E	ntity Co	oref Lev	vel	E	vent Co	oref Lev	el
Proportion	Р	R	F1	EM	P	R	F1	EM	P	R	F1	EM
					DMI	BERT						
1%	83	7.8	7.9	7 9	10.4	95	0.7	8.6	9.4	8.6	8.8	77
3%	11.1	10.7	10.8	10.1	11.9	$\frac{9.3}{11.2}$	11.3	10.3	10.2	9.6	9.7	8.8
5%	15.2	14.7	14.8	14.0	14.8	13.9	14.0	12.9	13.1	12.3	12.5	11.5
7%	17.0	16.4	16.5	15.7	17.9	16.8	17.1	15.7	16.4	15.4	15.6	14.4
10%	18.0	17.4	17.5	16.6	17.5	16.4	16.6	15.4	16.0	15.0	15.2	14.1
20%	22.6	22.0	22.1	21.2	21.0	19.8	20.0	18.6	19.3	18.2	18.4	17.1
30%	25.7	25.0	25.2	24.2	23.1	21.7	21.9	20.4	21.5	20.2	20.4	19.0
50%	26.9	26.3	26.4	25.5	23.9	22.7	22.9	21.4	22.1	21.0	21.2	19.8
70%	27.5	26.9	27.0	26.1	24.0	22.7	23.0	21.5	22.1	21.0	21.2	19.9
90%	29.2	28.6	28.7	27.8	24.7	23.5	23.8	22.2	22.9	21.8	22.0	20.6
					BERT	F+CRF						
1%	16.4	14.8	14.9	11.6	21.7	19.6	19.8	15.4	21.1	19.1	19.3	15.1
3%	25.3	24.1	23.8	19.2	31.4	29.6	29.4	23.4	30.4	28.8	28.5	22.9
5%	32.5	31.2	30.9	25.8	38.8	36.8	36.5	29.6	38.0	36.2	35.9	29.2
7%	33.9	32.8	32.4	27.1	41.6	39.7	39.3	31.9	40.7	38.9	38.5	31.3
10%	36.8	36.0	35.5	30.0	43.1	41.8	41.1	33.9	42.3	41.1	40.4	33.4
20%	40.5	39.7	39.1	33.7	46.6	45.2	44.5	37.0	45.7	44.4	43.7	36.4
30%	42.0	41.3	40.7	35.1	47.2	45.9	45.2	37.6	46.4	45.3	44.5	37.0
50%	42.1	42.2	41.5	30.3	48.3	4(.2	40.4	39.0	41.3	40.4	45.0	38.3
70%	44.3	43.9	43.2	31.8	49.5	48.0 47 E	41.0	40.0	48.0	41.1	40.8 45 5	39.4 20 E
90% 100%	43.8	43.7 73.0	42.0 /3.1	37.3	47.9	47.0 47.5	40.4	39.∠ 30.3	40.9	40.0 46.6	45.0 45.8	38.7
	11.0	10.0	10.1	01.1	10.1 FF	11.0 M	10.0	00.0	11.1	10.0	10.0	
	110	10.0	10.0					10.0			14.0	
1%	14.0	13.2	12.9	9.9	16.6	16.3	15.7	12.2	15.7	15.4	14.8	11.5
3% F07	20.5	18.8	18.7	14.4	23.3	22.1	21.7	17.0	22.4	21.3	20.8	10.4
370 70%	24.0 24.2	21.0 21.8	21.7 21.0	17.0 17.1	27.1	25.0	25.0	20.0	20.1	24.2 24.5	24.0 24.4	19.2 10.5
1/0	24.2	21.0 23.4	21.9 23.5	185	27.1	25.2 27.0	25.0	20.0 21.7	20.4	24.0	24.4 26.2	19.0
20%	25.3 26.7	$\frac{20.4}{24.5}$	$\frac{20.0}{24.5}$	10.5	29.1	$\frac{21.0}{28.2}$	$\frac{20.9}{28.0}$	$\frac{21.7}{22.0}$	20.5	20.4 27.4	20.2 27.2	$\frac{21.1}{22.2}$
30%	26.1	$\frac{24.0}{24.1}$	24.0 24.1	19.0	29.6	$\frac{20.2}{27.8}$	$\frac{20.0}{27.7}$	$\frac{22.5}{22.7}$	28.6	27.4 27.0	21.2 26.7	21.2
50%	27.8	25.5	25.6	20.7	31.3	29.4	$\frac{29.2}{29.2}$	24.0	30.3	28.5	28.3	23.2
70%	28.1	25.7	25.8	20.9	31.4	29.5	29.3	24.2	30.5	$\frac{-0.0}{28.7}$	$\frac{-0.0}{28.5}$	23.5
90%	28.2	26.0	26.0	21.0	31.5	29.8	29.6	24.4	30.7	29.1	28.8	23.7
100%	28.3	26.1	26.1	21.1	31.6	29.9	29.6	24.5	30.8	29.2	28.9	23.7
					PA	IE						
1%	33.6	32.7	31.8	25.3	39.3	38.4	37.4	30.1	39.1	38.3	37.3	30.0
3%	37.0	36.0	35.2	28.6	43.2	42.0	41.1	33.7	43.1	42.0	41.1	33.7
5%	38.6	37.4	36.7	30.0	45.9	44.4	43.6	35.8	46.2	44.8	43.9	36.1
7%	39.8	39.0	38.1	31.5	45.8	45.0	43.9	36.5	46.1	45.3	44.2	36.7
10%	40.6	40.0	39.0	32.4	46.9	46.2	45.1	37.6	47.2	46.6	45.4	37.9
20%	43.2	42.1	41.3	34.7	49.5	48.3	47.4	39.9	49.8	48.6	47.7	40.1
30%	43.2	42.1	41.3	34.7	49.5	48.3	47.4	39.9	49.8	48.6	47.7	40.1
50%	43.4	42.6	41.8	35.3	49.9	49.0	48.0	40.6	50.3	49.4	48.4	40.9
70%	44.0	43.0	42.3	35.8	50.7	49.5	48.7	41.3	51.2	50.1	49.1	41.7
90%	44.4	43.4	42.7	36.2	51.3	49.9	49.1	41.8	51.5	50.3	49.4	42.0
100%	44.5	43.4	42.7	36.3	50.8	49.4	48.7	41.4	51.1	49.8	49.0	41.7

Table 12: Experimental results (%) of the EAE models using different proportions of training data of MAVEN-ARG. In this experiment, we adopt the gold trigger evaluation approach (Peng et al., 2023c).

Madal	#Donoma		Mentio	n Level		E	ntity Co	oref Lev	el	Event Coref Level			
WIGUEI	#Parallis	P	R	F1	EM	Р	Ŕ	F1	EM	Р	R	F1	EM
Entity Argument													
DMBERT	110M	19.7	19.7	19.7	19.5	12.5	12.4	12.4	12.3	11.8	11.8	11.8	11.6
CLEVE	355M	22.1	22.1	22.1	22.0	13.2	13.2	13.2	13.0	12.3	12.2	12.2	12.1
BERT+CRF	110M	18.6	18.5	18.5	17.8	12.7	12.6	12.6	12.0	12.0	11.8	11.8	11.3
EEQA	110M	6.3	6.2	6.2	5.6	9.1	9.1	9.0	8.3	8.5	8.5	8.4	7.7
Text2Event	770M	5.5	5.6	5.5	5.2	4.0	4.0	4.0	3.7	3.2	3.2	3.1	2.9
PAIE	406M	20.4	20.5	20.3	19.2	21.0	21.1	20.9	19.8	20.0	20.1	19.9	18.9
				N	on-Enti	ty Argu	ment						
BERT+CRF	110M	24.8	24.8	24.0	19.4	-	_	_	_	25.3	25.3	24.5	19.6
EEQA	110M	18.9	17.6	17.5	13.9	—	_	_	_	18.6	17.4	17.2	13.7
Text2Event	770M	1.7	1.7	1.6	1.1	_	_	_	_	1.5	1.6	1.5	1.1
PAIE	406M	39.4	38.4	37.6	30.4	-	_	_	_	38.7	37.8	36.9	29.7

Table 13: Experimental results (%) of existing state-of-the-art fine-tuned EAE models on entity and non-entity arguments of MAVEN-ARG. Classification-based models, e.g., DMBERT and CLEVE, are inapplicable to non-entity arguments.