CycleOIE: An Unsupervised Low-Resource Training Framework For Open Information Extraction

Anonymous ACL submission

Abstract

In Open Information Extraction (OpenIE), the acquisition of manually annotated sentenceextraction pairs is expensive, while automatically labeled datasets may struggle to accurately reflect real-world requirements for OpenIE systems. Existing neural models often demonstrate impressive performance on largescale training sets but falter when tested on smaller-scale datasets due to the discrepancy in attributes between the training and test sets. In real-world scenarios, it is crucial for OpenIE systems to align closely with test sets, even when faced with limited annotated data for training.

> This paper introduces CycleOIE, a novel training framework applied to a pair of inverse text-to-text models. Through CycleOIE, we train a pair of T5 models on our curated dataset, LSOIE-g, achieving performance levels that surpass baselines trained on significantly larger fully supervised training sets. The ablation study offers a detailed comparison between fully supervised training and CycleOIE, highlighting the effectiveness of CycleOIE on LSOIE-g as the primary factor in enhancing T5's OpenIE performance.

1 Introduction

011

014

016

017

021

022

027

037

041

Open Information Extraction (OpenIE) is a type of Information Extraction task that extracts structured information, such as triples (subject; relation; object), from a given sentence (Yates et al., 2007; Angeli et al., 2015). This task is crucial for various downstream NLP applications, including summarization, knowledge graph construction, and knowledge base question answering.

Neural models for OpenIE demand substantial training data, but manually labeling datasets meeting requirements of real-world applications is costly. Publicly available large-scale OpenIE datasets like OIE2016(Stanovsky and Dagan, 2016), IMoJIE(Kolluru et al., 2020b) and LSOIE(Solawetz and Larson, 2021) are often weakly labeled. Annotations of these datasets are either high confidence extractions filtered out of predictions of existing SOTA OpenIE systems or automatically converted from datasets of other NLP tasks. Though these weakly labeled datasets are employed for training, they are seldom used for evaluation for limited quality. Some OpenIE benchmarks like WiRe57(Lechelle et al., 2019), CaRB(Bhardwaj et al., 2019) and BenchIE(Gashteovski et al., 2022) are presented. The data scale of these benchmarks are so small that they can only be used for evaluation. They have no split for training set for their small scales which are usually no more than 1k sentences. Meanwhile, these small-scale test sets are either labeled by experts or crowdsourcing under a explicitly stated annotation guidelines. Containing richer and inferred relations(Pei et al., 2023), these smallscale test sets play better roles in simulating the scenarios of real world applications and evaluating OpenIE systems.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

Researchers design novel neural architectures to fit large-scale weakly labeled training sets, aiming to enhance performance on small-scale test sets labeled by experts or crowdsourcing. However, disparities in attributes between training and test sets, such as the source and number of inferred relations persist(Pei et al., 2023). Consequently, model architectures, regardless of novelty, only impact how well the model fits the training set. Fine-tuning is then required for neural models to align with test data, but obtaining annotated data for fine-tuning is challenging in both academic and real application scenarios.

This paper introduces CycleOIE, an unsupervised training framework tailored for low-resource scenarios. We first leverage gpt-3.5 to generate extractions, constructing a training set for CycleOIE. Specifically, we instruct gpt-3.5 with the prompt words consisting the sentence from LSOIE

wiki¹ validation set and annotation guidelines of BenchIE which adapts to both BenchIE and CaRB 084 dataset. The combination of LSOIE wiki validation sentences and extractions generated by gpt-3.5 is named as LSOIE-g. Subsequently, we train our OpenIE model with CycleOIE. Though LSOIEg is an parallel dataset, during the cycle training process, its sentences and extractions are shuffled and loaded separately, which means only one side 091 of the sentence-extraction pair is utilized in training one model. This characteristic defines CycleOIE as unsupervised, aligning with the practical consideration that acquiring one side of the data is often more accessible in real-world scenarios. In each epoch of CycleOIE, there are two cycles, which are Sentence-Extraction-Sentence (SES) cycle and Extraction-Sentence-Extraction (ESE) cycle. Each cycle needs two models which are 100 Sentence-to-Extraction(S2E) model and Extraction-101 to-Sentence(E2S) model. The data is input to 102 model- 1^2 in eval mode. Then the prediction of the 103 model-1 is input to model-2 in train mode. So, only model-2's parameter is updated in a cycle by computing the loss between the prediction of model-2 106 107 and the input data to model-1 which is from one side of LSOIE-g. In SES cycle, S2E model is model-1 and E2S model is model-2 while in ESE 109 cycle, E2S model is model-1 and S2E model is 110 model-2. In CycleOIE, we run SES cycle and ESE 111 cycle by turns so that E2S model and S2E model 112 can be trained as model-2 by turns. In evaluation, 113 only S2E model is needed to predict extractions. 114

> The central idea behind CycleOIE is that a good extraction should contain all facts appeared in the sentence. Therefore, an extraction can be predicted by an S2E model, and conversely, a sentence can be predicted by an E2S model. No parallel annotation data (i.e., sentence-extraction pair) is needed for training, assuming a well-performing pretrained model capable of predicting the other side of the data. In a cycle, we freeze the parameters of model-1 and train model-2. We use model-1 to predict the other side of data so the input of model-2 is obtained. Then compute the loss of model-2 with the data pair composed up of the output of model-2 and the input of model-1 (i.e., the original input data). We assume model-1 is a well performing pretrained model, capable to predict data of the other side while it is not. So, in the next cycle,

115

116

117

118

119

121

122

123

124

125

126

127

128

129

130

131

we reverse the position of two models and let it get 132 trained as model-2. In CycleOIE, we run SES cycle 133 and ESE cycle by turns so that E2S model and S2E 134 model can be trained as model-2 by turns. After 135 each epoch, the parameters of model-2 is updated 136 and various intermediate products will be output 137 in the next cycle. So CycleOIE can be viewed as a 138 type of data augmentation along with the training 139 process. This approach eliminates the need for a 140 highly performing pretrained S2E or E2S model 141 at the outset and avoids the requirement for large-142 scale, high-quality parallel data. 143

sentence	Earlier today, Thailand 's Prime Min-
	ister Yingluck Shinawatra formally
	dissolved the country 's parliament
	and called for new elections.
extraction	subject <is> the country 's parliament</is>
	<and> relation <is> dissolved <and></and></is></and>
	object <is> Earlier today Thailand</is>
	's Prime Minister Yingluck Shinawa-
	tra <then> subject <is> Thailand 's</is></then>
	Prime Minister Yingluck Shinawatra
	<and> relation <is> called <and> ob-</and></is></and>
	ject <is> Earlier today new elections</is>

Table 1:	<is>,</is>	<and>,</and>	<then></then>	is	added	to	format
dataset to	comply	with a to	ext-to-tex	xt r	nodel.		

Our S2E model and E2S model are both finetuned from flan-t5-base³(Chung et al., 2022), a text-to-text transformer that utilizes textual input and output. To adapt T5 for the S2E and E2S tasks, we design a template for extractions with three additional tokens: $\langle is \rangle$, $\langle and \rangle$, $\langle then \rangle$ respectively, as illustrated in Table 1. 144

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

162

164

165

We conduct experiments on various OpenIE benchmarks, revealing performance that surpasses baselines trained on much larger-scale training set in full supervision. Our contributions are summarized as follows:

(i) We introduce CycleOIE, a training framework that employs cycle training for the OpenIE task, implemented on a text-to-text pretrained language model, T5.

(ii) We design a format for the extraction sequence to align with both the input and output format of a text-to-text model.

(iii) We construct LSOIE-g, a training set for CycleOIE containing approximately 2k sentences, significantly fewer than typical OpenIE training

¹https://huggingface.co/datasets/wardenga/lsoie/viewer/wiki ²In this ariticle, we use **model-1** to refer to <u>the first model</u> and **model-2** to refer to <u>the second model</u> of current cycle.

³https://huggingface.co/google/flan-t5-base

166 167

168

170

171

172

Related Work

from a given sentence. Early OpenIE methods

(Yates et al., 2007; Angeli et al., 2015; Del Corro

and Gemulla, 2013; Gashteovski et al., 2017) rely

on linguistic expertise for sentence parsing, while

recent neural methods harness deep neural net-

works to represent the semantics with hidden states.

gorized into sequence labeling and sequence gen-

eration methods. Sequence Labeling methods (Ro

et al., 2020; Kolluru et al., 2020a) usually use an

encoder (e.g., BERT) to encode the given sentence

into its embedding. A sequence labeling head con-

nected to the encoder labels each token in the se-

quence. On the other hand, sequence generation

methods (Cui et al., 2018; Kolluru et al., 2020b)

design prompts or generation templates to trans-

form information extraction task into text genera-

tion task. Generally, sequence labeling methods

offer faster predictions since they can be performed

in parallel while their extractions are limited to

tokens from the given sentence, potentially lead-

ing to incoherence or grammatical mistakes. In

contrast, sequence generation methods are slower

as the decoder generates from left to right sequen-

tially, while their generated outputs tend to be more

Cycle training, also known as cycle-consistency

training, is a training framework that use unpaired

data to train a pair of inverse models (i.e., out-

put and input of model-1 could become the input

and output of model-2 in reverse.). It is initially

introduced in the machine translation task as the

term iterative back-translation(Hoang et al., 2018)

to solve the challenge of lack of sentence pairs

composed of source language sentences and target

coherent and adaptive.

2.2 Cycle training

language sentences.

Neural methods in OpenIE can be broadly cate-

2

2.1

- 173
- 174 175

176

178 179

183

184

187

190 191

192

195

196 197

198

199

206

210 211

214

Cycle training is widely used in text generation tasks. Hoang et al. (2018); Wei et al. (2020); Dou

sets. LSOIE-g can be used to mirror the highly et al. (2020) manage to use this method to overlimited annotation data in real world scenarios. come the scarcity of paired sentence datasets in (iv) We conduct extensive experiments demonmachine translation area. Iovine et al. (2022b) instrating that unsupervised CycleOIE is able to outtroduce CycleKQR, leveraging cycle training to perform some fully supervised OpenIE systems, enhance Question Answering (QA) performance even with a very limited amount of training data. by rewriting queries into appropriate forms while retaining semantics. Wang et al. (2023) evaluate the effectiveness of cycle training in ensuring consistency between structured data and text, achiev-**Open Information Extraction** ing performance comparable to fully supervised Open Information Extraction is a fundamental NLP approaches for data-to-text generation tasks. Retask aiming at extracting structured information

cently, cycle training has been applied to train information extraction models. Iovine et al. (2022a) apply cycle training on T5 to address the lack of indomain annotation data, achieving competitive performance with fully supervised models on Named Entity Recognition (NER) tasks. With NER and many other information extraction tasks, OpenIE share a similar challenges of insufficient in-domain annotation data. this highlights the potential of cycle training to handle OpenIE tasks. The key lies in designing templates that induce extraction

generations consistent with their source sentences.

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

240

241

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

3 Method

3.1 CycleOIE

We present CycleOIE to address the challenge of lack of paralleled training set and train our model for OpenIE. CycleOIE is implemented by two textto-text models (with T5 serving as our backbone model), referred to as S2E(Sentence to Entity) model and E2S(Entity to Sentence) model as shown in Figure 1. S2E model (at the top-right corner of Figure 1) is expected to predict the extraction given a sentence as input. Each extraction is made up of triples consisting a head entity, a relation entity and a tail entity. (subject, relation and object syntactically.) E2S model is expected to predict the corresponding sentence given an extraction. CycleOIE is trained through two cycles. In each cycle, the first model is freezed and only the second model's weights will be updated.

Specifically, in ESE cycle (the left part of Figure 1, the e is input to the E2S model in eval mode to predict \hat{s} . Then \hat{s} is input to the S2E model (in train mode) to predict \hat{e} . Loss between \hat{e} and e is computed to update the parameters of the S2E model through back propagation. After training the S2E model for 1 epoch, the training process comes to SES cycle. In SES cycle (the right part of Figure 1, the s is input to the S2E model (in



Figure 1: CycleOIE training framework

eval mode) to predict \hat{e} . Then \hat{e} is input to the E2S model (in train mode) to predict \hat{s} . Loss between \hat{s} and s is computed to update the parameter of the E2S model through back propagation. Once both models have completed one epoch of training, CycleOIE is one epoch updated. Table 2 demonstrates CycleOIE conducted on batched training data.

270

271

272

273

275

276

281

Algorithm: CycleOIE
Input: Dataset of sentences D_S extractions D_E
Output: S2E model M_{S2E} and E2S model M_{E2S}
while M_{S2E} and M_{E2S} have not converged do
for every batch S in D_S :
Transform S into E' using M_{S2E}
Train M_{E2S} with (E', S)
end for
for every batch E in D_E :
Transform E into S' using M_{E2S}
Train M_{S2E} with (S', E)
end for
end while

Table 2: Algorithm: CycleOIE

3.2 OpenIE as sequence generation

For taking better advantages of pretrained text-totext model, we design a generation template that guides the language model to generate textual extractions. The model's tokenizer is extended with three additional tokens: *<is>*, *<and>*, *<then>*. As illustrated in Table 1 , *<is>* is utilized, in combination with its previous token, to indicate which part of a triple should be generated next; *<and>* is used to separate head entity, relation entity and tail entity; *<then>* is used to delimit triples in a extraction.

For training the S2E model, the extraction se-

quence, denoted as $e = \{e_1, e_2, e_3, ..., e_n\}$, is input to the freezed E2S model, letting E2S model to predict a sentence sequence $\hat{s} = \{\hat{s}_1, \hat{s}_2, \hat{s}_3, ..., \hat{s}_n\}$. The generated sentence sequence \hat{s} is then fed into the S2E model, aiming to predict $\hat{e} =$ $\{\hat{e}_1, \hat{e}_2, \hat{e}_3, ..., \hat{e}_n\}$ that aligns with the original extraction sequence e. The loss is computed between \hat{e} and e.

289

290

291

292

293

294

297

298

299

300

301

302

303

304

305

308

309

310

311

312

$$\mathcal{L}_{\theta}(E, \hat{E}) = -\frac{1}{|E|} \sum_{e \in E} \frac{\sum_{i < |e|} p(e_i) \log p(\hat{e}_i)}{|e|}$$
(1)

Here, E represents a batch of e and |E| is the batch size, θ denotes the parameter of S2E model, p(.)signifies probability of token, e_i is the i-th token in e and \hat{e}_i is the i-th token in \hat{e} , and |e| represents the length of the sequence of extractions.

For training E2S model, the sentence sequence $s = \{s_1, s_2, s_3, ..., s_n\}$ is input to the freezed S2E model, letting S2E model to predict the extraction sequence $\hat{e} = \{\hat{e}_1, \hat{e}_2, \hat{e}_3, ..., \hat{e}_n\}$. The generated extraction sequence \hat{e} is then fed into the E2S model, expecting its prediction $\hat{s} = \{\hat{s}_1, \hat{s}_2, \hat{s}_3, ..., \hat{s}_n\}$ to fit s. The loss is then calculated between \hat{s} and the gold s.

$$\mathcal{L}_{\phi}(S, \hat{S}) = -\frac{1}{|S|} \sum_{s \in S} \frac{\sum_{i < |s|} p(s_i) \log p(\hat{s}_i)}{|s|} \quad (2)$$

Here S represents a batch of s and |S| is the batch size, ϕ denotes the parameters of the E2S model, s_i is the i-th token in s, \hat{s}_i is the i-th token in \hat{s} , and |s|represents the length of text sequence of the given sentence.

3.3 Build CycleOIE Dataset

313

314

315

316

317

321

327

334

335

337

339

341

342

347

349

353

As pointed out by Pei et al. (2023), attribute discrepancies generally exists between weakly labeled large-scale training sets and manually labeled small-scale test sets. These discrepancies manifest in factors such as the source of sentences and the proportion of inferred relations and N-ary relations they retain. Table 3 illustrates these differences; LSOIE lacks inferred relations presented in WiRe57, CaRB, while IMoJIE lacks N-ary relations found in WiRe57 and CaRB. Both IMoJIE and LSOIE sentences originate from sources different from CaRB. Achieving high performance on these human-labeled test sets through training on weakly labeled training sets poses a challenge. To address this, we utilize gpt-3.5 to generate extractions adhere to the annotation guidelines of BenchIE⁴ given LSOIE wiki validation set sentences as input. Thus we obtain extractions for those sentences. We name the dataset LSOIE-g which comprises LSOIE wiki validation set sentences and gpt-3.5 generated extractions. Although LSOIE-g pairs these extractions with their input sentences, during the cycle training, we use sentences in the SES cycle and extractions in the ESE cycle separately. Through experiments, we observe that our model cycle trained on LSOIE-g outperforms some neural baselines fine-tuned in fully supervised settings on much larger training sets.

4 Datasets

As mentioned in the introduction, neural OpenIE systems are typically trained and tested on different datasets. Datasets used in our experiments are listed in Table 3.

LSOIE(Solawetz and Larson, 2021) is a large-scale OpenIE dataset converted from QA-SRL 2.0 with similar conversion method to OIE2016(Stanovsky and Dagan, 2016). IMo-JIE(Kolluru et al., 2020b) aims to construct a high quality OpenIE dataset with Wikipedia sentences and high confidence extractions predicted by former OpenIE systems including OpenIE4, ClausIE, and RNNOIE. Wire57(Lechelle et al., 2019) releases a tiny dataset comprising 57 sentences from 5 documents and extractions annotated by 2 experts. CaRB(Bhardwaj et al., 2019) is an OpenIE benchmark with 641 sentences in its test set. Obtained through crowdsourcing, CaRB's annotations are generally considered as an noise reduction from OIE2016(Stanovsky and Dagan, 2016). BenchIE(Gashteovski et al., 2022) releases a dataset of 300 sentences. Two expert annotators try to exhaust every possible extraction of facts in a sentence. Annotations of the same fact are classified into one cluster. Unlike WiRE57 and CaRB utilizing token-level scoring functions, BenchIE judges an predicted triple to be true only when it exactly matches a gold triple in the cluster of the fact. we find this scoring function to be too strict, resulting in significant differences in results on BenchIE compared to other benchmarks. To address this, We introduce a new benchmark, CaRB-B, which combines the scoring function of CaRB and the dataset of BenchIE, replacing BenchIE's scoring function and retaining its dataset.

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

388

389

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

CaRB, due to its quantity advantage among those human-labeled datasets, is selected as our main target for evaluation. Considering the similarity between LSOIE and CaRB⁵, we use sentences of LSOIE wiki validation set to compose our training set for CycleOIE to simulate a real scenario where the number of extraction annotations is very limited.

5 Experiment

5.1 Baseline Systems

We compare our methods with 4 previous OpenIE systems: ClausIE(Del Corro and Gemulla, 2013), MINIE(Gashteovski et al., 2017), Multi²OIE(Ro et al., 2020), and OpenIE6(Kolluru et al., 2020a). The former two are classical rule-based methods while the latter two are neural-based methods.

ClausIE(Del Corro and Gemulla, 2013) takes the advantage of English grammatical knowledge to extract information from a sentence without any training. By its core step, Dependency Parsing, the syntactic relation between words in a sentence is analysed. Then matching the dependency to the clauses, combining the properties of predicates and the structure of clauses, extractions from clauses are generated. Some worries are presented that the extraction of ClausIE may consist multiple propositions merging into a over specific value(subject,

⁴BenchIE samples 300 sentences and re-annotated by experts under their guideline. Through our research on CaRB and BenchIE, we think the annotation guidelines of BenchIE has many common requirements with CaRB and can instruct gpt-3.5 to generate extraction annotations for CaRB sentences.

⁵OIE2016 is automatically converted from QA-SRL(He et al., 2015) while LSOIE is automatically converted from QA-SRL 2.0(FitzGerald et al., 2018). CaRB is a re-annotation of OIE2016

	Dataset	Source	Sentences	Extractions	Inferred Relations	N-ary relations
training sate	IMoJIE	wikipedia	71209	215K	3K	0
training sets	LSOIE(wiki)	QA-SRL, wikipedia	12832	101K	0	32K
	LSOIE-g	LSOIE, gpt-3.5	2147	6820	383	0
	WiRe57	Wikipedia, Newswire	57	343	173	79
test sets	CaRB	OIE2016	641	2715	736	683
	BenchIE	CaRB	300	783	0	0

Table 3: Statistics of datasets used in our experiments

relation, or object) of a tuple representing an extraction result. These extraction results of less conciseness probably don't meet the demand of its downstream tasks like knowledge graph constructions. Minie(Gashteovski et al., 2017) is an OpenIE system built on the top of ClausIE, aiming to provide compact extractions, i.e., minimizing each value in the tuple representing extraction result.

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

449

443

444

445

Multi²OIE(Ro et al., 2020) first leverages a pretrained language model to encode the context of the sentence. Multi²OIE regards OpenIE as a kind of sequence labeling task. Predicates are first labeled by a predicate classifier head concatenated to the BERT model. Then arguments are predicted by an argument classifier given the output of BERT model including the hidden state vector and predicate average vector. Among neural methods, While methods of sequence labeling outputs predictions with faster speed, methods of generation produce extractions of better quality. OpenIE6(Kolluru et al., 2020a), using Iterative Grid Labeling architecture to label a 2-D grid which represents an extraction result with a row, surpassing generative methods(Cui et al., 2018; Kolluru et al., 2020b) achieves SOTA at that time.

5.2 Main results

We compare our CycleOIE with baseline systems. As shown in Table 4, our unsupervised method, CycleOIE, achieves the highest precision and F1 on both CaRB and CaRB-B benchmark, comparing with tradition rule-based methods and two fully supervised neural methods. When evaluating on WiRe57 CycleOIE's F1 only defeats Multi²OIE. Neural methods' performance can be influenced by various factors beyond model architecture, such as the scoring function and attribute discrepancies between the training set and the test set. In Table 5, we conduct additional experiments controlling variables to evaluate the performances of neural methods when they are trained on LSOIE whose sentences share a same source with our LSOIE-g. We compare CycleOIE's performance with neural methods trained on LSOIE wiki training set. The results indicate CycleOIE outperforms these two OpenIE6 and Multi²OIE on WiRe57⁷, CaRB⁸ and WiRe57-C⁹ benchmark even though CycleOIE is training on a dataset with a significantly smaller scale than the LSOIE training set. Notably, when evaluating on CaRB, the performance gap between CycleOIE and other baseline methods becomes more pronounced than the gap illustrated in Table4. 446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

5.3 Ablation study

We conduct an ablation study and results are presented in Table 6. Here are the details of our experimental settings.

Setting (1) **sft on LSOIE** stands for training flant5-base on LSOIE wiki training set in fully supervised fine-tuning and (2) **sft on LSOIE+IMoJIE** stands for a similar setting except for training on the union of LSOIE wiki and IMoJIE training set.

In (3) sft on LSOIE \rightarrow sft on LSOIE-g setting, we train the model by 2 steps. At step 1, we do the same training as (1). At step 2, we train the weights output by step 1 on LSOIE-g data with full supervision. In the step1 of (4) sft on LSOIE+IMoJIE \rightarrow sft on LSOIE-g, we do the same training as (2), and the other operations keep the same with the previous one.

In (5) sft on LSOIE \rightarrow CycleOIE setting, we train flan-t5-base in fully supervised fine-tuning on LSOIE wiki training set at step 1, then cycle trained the model output by step 1 on LSOIE-g data. While in setting (6) sft on LSOIE+IMoJIE \rightarrow CycleOIE the training set at step 1 is augmented.

Setting (10) CycleOIE stands for we cycle train flan-t5-base on LSOIE-g without any supervised fine-tuning to bootstrap. Setting (7), (8), (9) stands for 1/4, 1/2, 3/4 of sentences and extractions in

⁷WiRe57 scoring function + Wire57 test set

⁸CaRB scoring function + CaRB test set

⁹WiRe57 scoring function + CaRB test set

	WiRe57			CaRB			CaRB-B		
	Р	R	F1	Р	R	F1	Р	R	F1
ClausIE	0.401	0.298	0.342	0.411	0.496	0.450	0.580	0.534	0.556
MINIE	0.400	<u>0.323</u>	<u>0.358</u>	0.429	0.382	0.404	0.466	0.436	0.441
OpenIE6	0.465	0.326	0.383	0.589	<u>0.476</u>	<u>0.527</u>	0.478	0.671	0.559
$Multi^2OIE$	<u>0.457</u>	0.182	0.261	<u>0.609</u>	0.458	0.523	<u>0.598</u>	<u>0.613</u>	<u>0.605</u>
CycleOIE	0.388	0.205	0.268	0.691	0.427	0.528	0.603	0.613	0.608

Table 4: Overall evaluation. CaRB-B stands for combining the scoring function of CaRB and the dataset of BenchIE, while WiRe57 and CaRB use the scoring function and dataset of their own. In this table, the author-published baselines are evaluated, which means the training set is defined by their author and probably not unified. OpenIE6 was trained on the OpenIE4 training dataset used to train IMoJIE. Mult2OIE is trained on the training dataset of SpanOIE⁶(Zhan and Zhao, 2020).

scoring function			CaRB			WiRe57		
model	training set	test set	Р	R	F1	Р	R	F1
OpenIE6	LSOIE	WiRe57	0.311	0.247	0.275	0.311	<u>0.194</u>	<u>0.239</u>
$Multi^2OIE$	LSOIE	WiRe57	0.440	0.202	<u>0.276</u>	0.440	0.128	0.198
CycleOIE	LSOIE-g	WiRe57	<u>0.385</u>	0.271	0.318	<u>0.388</u>	0.205	0.268
OpenIE6	LSOIE	CaRB	0.403	0.389	0.396			-
$Multi^2OIE$	LSOIE	CaRB	<u>0.611</u>	0.369	<u>0.461</u>	-	-	-
CycleOIE	LSOIE-g	CaRB	0.691	0.427	0.528	-	-	-

Table 5: Above the dashline, we compare our CycleOIE, which is trained on LSOIE-g to those models trained on LSOIE wiki training set, on the test data of WiRe57 with CaRB's scoring function and WiRe57's scoring function respectively. Below the dash line, we compare our CycleOIE, which is trained on LSOIE-g, to those models trained on LSOIE wiki training set, on the test data of CaRB, with scoring function of CaRB.

LSOIE-g is sampled as training set respectively.

In settings above, when models are trained on LSOIE-g in supervised fine-tuning, the LSOIE-g is paralleled, which means the sentences and the corresponding extractions are paired up. When models are trained on LSOIE-g in CycleOIE, LSOIE-g is unpaired to simulate the real world scenarios, which means sentences and extractions are shuffled and loaded separately.

Table 6 indicates that the (10) CycleOIE setting, without any supervised fine-tuning bootstrap, achieves similar performance with (6) sft on $LSOIE+IMoJIE \rightarrow CycleOIE$, demonstrating the promise of CycleOIE in real-world OpenIE scenerios, where annotation data is often insufficient.

Comparing setting (1) and (2), we observe that augmenting the training set with IMoJIE significantly improves recall, albeit at the cost of precision. Comparing (3)(4) or (5)(6) to (1)(2), it is observed that once the model is trained on LSOIEg (either in sft or CycleOIE), precision and f1 show significant improvements. This suggests that our built LSOIE-g is effective as a training dataset for CaRB and CaRB-B. Our designed prompt words, which are derived from BenchIE annotation guidelines, successfully instruct gpt-3.5 to be a qualified annotator. 506

507

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

Comparing (5) to (3), the former exhibits higher performances in precision, recall and f1. Comparing (6) to (4), the former demonstrates advantages in recall and F1, but the gap, taking F1 as the main metric for comparison, is smaller than the gap between (5) and (3). This aligns with our understanding that cycle training can be considered a form of continuous data augmentation along with the training process, significantly improving the performance when the model has not yet acquired enough domain knowledge from existing training data. Conversely, When the existing training data contains sufficient domain knowledge for this task, cycle training doesn't have a prominent advantage over sft.

We observe that (5), (6) and (10) have similar performances, indicating the model cycle trained

505

		CaRB			5	
	Р	R	F1	Р	R	F1
(1) sft on LSOIE	0.573	0.339	0.426	0.488	0.459	0.473
(2) sft on LSOIE+IMoJIE	0.402	0.453	0.426	0.332	0.598	0.427
(3) sft on LSOIE \rightarrow sft on LSOIE-g	0.654	0.418	0.51	0.592	0.576	0.584
(4) sft on LSOIE+IMoJIE \rightarrow sft on LSOIE-g	0.704	0.415	0.523	0.636	0.576	0.604
(5) sft on LSOIE \rightarrow CycleOIE	0.694	0.424	0.526	0.600	0.589	0.594
(6) sft on LSOIE+IMoJIE \rightarrow CycleOIE	<u>0.694</u>	0.428	0.530	<u>0.618</u>	0.602	0.610
(7) CycleOIE(1/4 LSOIE-g)	0.557	0.386	0.456	0.521	0.56	0.54
(8) CycleOIE(1/2 LSOIE-g)	0.616	0.412	0.494	0.557	0.577	0.567
(9) CycleOIE(3/4 LSOIE-g)	0.651	0.419	0.510	0.594	0.598	0.596
(10) CycleOIE	0.691	<u>0.427</u>	<u>0.528</u>	0.603	0.613	<u>0.608</u>

Table 6: Ablation study



Figure 2: Loss and F1 score of each epoch.

on LSOIE-g can acquire sufficient knowledge to solve the OpenIE task and hit the test data of CaRB. Supervised fine-tuning on wealky labeled largescale training sets helps to bootstrap but provide almost no additional domain knowledge compared to CycleOIE on LSOIE-g.

528

529

530

531

532

534

535

538

The bottom four lines demonstrates the performance of CycleOIE given different scale of LSOIEg training data. We observe that the model performs better on every metric with an increasing scale of training data, indicating our built LSOIE-g is suitable as the training set of CycleOIE, especially when target benchmark is CaRB and CaRB-B. More detailed influences caused by the training data size is illustrated in Figure 2 where we record the loss θ and F1 score of each epoch of CycleOIE. We observe that training on 1.00LSOIE-g, the loss descends in the fastest rate and then remains at the lowest position and the F1 score converges in the first place at after about 20 epochs. The fluctuation of the F1 score also gets smaller when CycleOIE is conducted on a larger training set. 539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

565

566

567

569

570

6 Conclusion

In this paper, we present CycleOIE, a novel framework training a pair of inverse text-to-text model to address OpenIE. We design a template to ensure extractions align with the input and output format of T5. Leveraging the capabilities of gpt-3.5, we construct LSOIE-g as the training set for CycleOIE. Experimental results demonstrate that our approach outperforms baselines trained under full supervision on much larger-scale OpenIE training sets.

Limitations

We implemented CycleOIE based on T5 without introducing novel model structures or task pipelines for further performance enhancement. Due to current constraints related to available GPUs and time limitations, we refrained from applying CycleOIE to fine-tune large language models such as LLaMA(Touvron et al., 2023) and ChatGLM¹⁰(Du et al., 2022), which are more powerful generative models. A speculative avenue for future exploration is the potential of large language models to

¹⁰https://github.com/THUDM/ChatGLM2-6B

seamlessly transition between sentence-generating
and extraction-generating modes by altering instructions after instruction tuning. This implies
the possibility of training a single model instead of
a pair. We aspire to delve into this aspect of large
language models in future research.

577 References

578

579

581

582

586

589

590

591

592

594

595

597

600

608

609

610

611

612

613

614 615

616

617

618

619

623

- Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D. Manning. 2015. Leveraging linguistic structure for open domain information extraction. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 344–354, Beijing, China. Association for Computational Linguistics.
 - Sangnie Bhardwaj, Samarth Aggarwal, and Mausam Mausam. 2019. CaRB: A crowdsourced benchmark for open IE. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 6262–6267, Hong Kong, China. Association for Computational Linguistics.
 - Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. arXiv preprint arXiv:2210.11416.
 - Lei Cui, Furu Wei, and Ming Zhou. 2018. Neural open information extraction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 407–413, Melbourne, Australia. Association for Computational Linguistics.
 - Luciano Del Corro and Rainer Gemulla. 2013. Clausie: clause-based open information extraction. In *Proceedings of the 22nd international conference on World Wide Web*, pages 355–366.
- Zi-Yi Dou, Antonios Anastasopoulos, and Graham Neubig. 2020. Dynamic data selection and weighting for iterative back-translation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5894–5904, Online. Association for Computational Linguistics.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, Dublin, Ireland. Association for Computational Linguistics.

Nicholas FitzGerald, Julian Michael, Luheng He, and Luke Zettlemoyer. 2018. Large-scale QA-SRL parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2051–2060, Melbourne, Australia. Association for Computational Linguistics. 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

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

669

670

671

672

673

674

675

676

677

678

- Kiril Gashteovski, Rainer Gemulla, and Luciano del Corro. 2017. MinIE: Minimizing facts in open information extraction. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 2630–2640, Copenhagen, Denmark. Association for Computational Linguistics.
- Kiril Gashteovski, Mingying Yu, Bhushan Kotnis, Carolin Lawrence, Mathias Niepert, and Goran Glavaš. 2022. BenchIE: A framework for multi-faceted factbased open information extraction evaluation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4472–4490, Dublin, Ireland. Association for Computational Linguistics.
- Luheng He, Mike Lewis, and Luke Zettlemoyer. 2015. Question-answer driven semantic role labeling: Using natural language to annotate natural language. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 643–653, Lisbon, Portugal. Association for Computational Linguistics.
- Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. Iterative backtranslation for neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24, Melbourne, Australia. Association for Computational Linguistics.
- Andrea Iovine, Anjie Fang, Besnik Fetahu, Oleg Rokhlenko, and Shervin Malmasi. 2022a. Cyclener: An unsupervised training approach for named entity recognition. In *Proceedings of the ACM Web Conference 2022*, WWW '22, page 2916–2924, New York, NY, USA. Association for Computing Machinery.
- Andrea Iovine, Anjie Fang, Besnik Fetahu, Jie Zhao, Oleg Rokhlenko, and Shervin Malmasi. 2022b. CycleKQR: Unsupervised bidirectional keywordquestion rewriting. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11875–11886, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Keshav Kolluru, Vaibhav Adlakha, Samarth Aggarwal, Mausam, and Soumen Chakrabarti. 2020a. OpenIE6: Iterative Grid Labeling and Coordination Analysis for Open Information Extraction. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 3748–3761, Online. Association for Computational Linguistics.

Keshav Kolluru, Samarth Aggarwal, Vipul Rathore,

Mausam, and Soumen Chakrabarti. 2020b. IMo-

JIE: Iterative memory-based joint open information

extraction. In Proceedings of the 58th Annual Meet-

ing of the Association for Computational Linguistics,

pages 5871-5886, Online. Association for Computa-

William Lechelle, Fabrizio Gotti, and Phillippe Langlais.

2019. WiRe57 : A fine-grained benchmark for open

information extraction. In Proceedings of the 13th Linguistic Annotation Workshop, pages 6–15, Flo-

rence, Italy. Association for Computational Linguis-

Kevin Pei, Ishan Jindal, Kevin Chen-Chuan Chang,

Association for Computational Linguistics.

Association for Computational Linguistics.

Youngbin Ro, Yukyung Lee, and Pilsung Kang. 2020. Multi²OIE: Multilingual open information extraction based on multi-head attention with BERT. In

Findings of the Association for Computational Lin-

guistics: EMNLP 2020, pages 1107–1117, Online.

Jacob Solawetz and Stefan Larson. 2021. LSOIE: A

large-scale dataset for supervised open information extraction. In *Proceedings of the 16th Conference of*

the European Chapter of the Association for Compu-

tational Linguistics: Main Volume, pages 2595–2600,

Online. Association for Computational Linguistics.

Gabriel Stanovsky and Ido Dagan. 2016. Creating a large benchmark for open information extraction.

In Proceedings of the 2016 Conference on Empir-

ical Methods in Natural Language Processing, pages 2300–2305, Austin, Texas. Association for Computa-

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay

Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti

Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint*

Zhuoer Wang, Marcus Collins, Nikhita Vedula, Simone

Filice, Shervin Malmasi, and Oleg Rokhlenko. 2023. Faithful low-resource data-to-text generation through

cycle training. In Proceedings of the 61st Annual

Meeting of the Association for Computational Lin-

guistics (Volume 1: Long Papers), pages 2847-2867,

Toronto, Canada. Association for Computational Lin-

Hao-Ran Wei, Zhirui Zhang, Boxing Chen, and Wei-

hua Luo. 2020. Iterative domain-repaired back-

translation. In Proceedings of the 2020 Conference

tional Linguistics.

arXiv:2307.09288.

guistics.

ChengXiang Zhai, and Yunyao Li. 2023. When to

use what: An in-depth comparative empirical analysis of OpenIE systems for downstream applications. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 929–949, Toronto, Canada.

tional Linguistics.

tics.

- 6 6 6
- 693 694 695
- 6 6 6
- 6 7
- 70
- 704 705
- 70
- 707 708
- 709
- 711 712
- 713 714

715

716 717 718

- 719 720 721
- 723
- 725

726

727 728 729

730 731

- 1
- 733
- 734 735

on Empirical Methods in Natural Language Processing (EMNLP), pages 5884–5893, Online. Association for Computational Linguistics.

736

738

739

740

741

742

743

745

746

747

748

749

750

- Alexander Yates, Michele Banko, Matthew Broadhead, Michael Cafarella, Oren Etzioni, and Stephen Soderland. 2007. TextRunner: Open information extraction on the web. In Proceedings of Human Language Technologies: The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT), pages 25–26, Rochester, New York, USA. Association for Computational Linguistics.
- Junlang Zhan and Hai Zhao. 2020. Span model for open information extraction on accurate corpus. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9523–9530.