Synthetic Context with LLM for Entity Linking from Scientific Tables

Anonymous ACL submission

Abstract

Tables in scientific papers contain crucial information, such as experimental results. Entity Linking (EL) is a promising technology that analyses tables and associates them with a knowledge base. EL for table cells requires identifying the referent concept of each cell while understanding the context relevant to each cell in the paper. However, extracting the relevant context from the paper is challenging because the relevant parts are scattered in the main text and captions. This study defines a 011 rule-based method for extracting broad context 012 013 from the main text, including table captions and sentences that mention the table. Furthermore, we propose synthetic context as a more refined context generated by large language models (LLMs). In a synthetic context, contexts from the entire paper are refined by summarizing, 018 injecting supplemental knowledge, and clari-019 fying the referent concept. We observe this approach improves accuracy for EL by more than 10 points on the S2abEL dataset, and our qualitative analysis suggests potential future works.

1 Introduction

027

Information analysis of scientific papers has numerous applications in accelerating science, such as paper retrieval, reading assistance, and automatic knowledge base construction. Particularly in information science, crucial information, such as experimental results, evaluation datasets, tasks, and evaluation metrics, is often recorded in tables within papers. Thus, the analysis of table information is an important research field.

For the table analysis, entity linking (EL) that associates table cells in scientific papers to a knowledge base (KB) is a promising technology, and various methods and datasets have been proposed for this purpose (Kardas et al., 2020; Yang et al., 2022; Lou et al., 2023). S2abEL (Lou et al., 2023) is a large-scale evaluation dataset for EL targeting tables in papers for the machine learning field. In



Figure 1: Example of entity linking for table cells in Devlin et al. (2019). Given a target table cell (e.g., "CoLA"), a model seeks to link it to the corresponding entity (e.g., "/dataset/cola") in *Papers with Code* by considering the contexts in the paper related to the cell.

the dataset, each table cell is linked to an entity defined in Papers With Code $(PwC)^1$, a free and open KB in the scientific domain, as illustrated in Figure 1. To correctly link the target table cell "CoLA" to the corresponding entity "/dataset/cola" in PwC, a model needs to understand the concept of CoLA from the contexts scattered in the main text, captions, and references.

043

045

047

049

051

057

060

061

062

063

064

065

066

However, extracting such contexts relevant to each cell from a paper has three technical challenges. (i) Relevant contexts for a cell text are scattered in an entire paper, and mentions are often abbreviated or paraphrased, (ii) The context or explanation for a referent concept of a cell can be insufficient, and (iii) General words such as "Ours," "Baseline," and "All" are often used in cell texts, and the referent is ambiguous. An example for the first, in the paper of The Evolved Transformer (So et al., 2019), a cell text "ET PERP" is interpreted as "the perplexity achieved by the Evolved Transformer", although the term "ET PERP" does not appear in the main text of the paper. Second, explanations for well-known methods such as LSTM are often omitted, and thus, sufficient contexts are

¹https://paperswithcode.com/

105

106

107

108

109

110

111

112

113

114

067

unavailable in the paper itself. Third, the cell text "All" stands for the entire dataset. Identifying the dataset requires understanding the context of the main text. However, the word "All" is general and frequently used in irrelevant contexts in the paper.

To address these challenges, we propose a data synthesis method for providing an EL model with supplemental contexts for table cells by using large language models (LLMs) such as Chat-GPT (OpenAI, 2023) and LLaMa2 (Touvron et al., 2023). LLMs, acquiring specialized knowledge through pre-training, can be utilized as knowledge bases (Taylor et al., 2022). They also demonstrate high performance in a zero-shot setting for abbreviation expansion (Gorman et al., 2021) and coreference resolution (Wei et al., 2022). Therefore, we expect that they can elicit supplemental information for table cells, which is not mined by previous methods.

In experiments, we use context data for EL obtained by our refined rule-based method and LLMbased method and confirm consistent improvements by both methods over the baseline proposed in S2abEL, including a 10-point improvement in accuracy for EL. Our qualitative analysis also reveals that synthetic context data captures better supplemental information through context completion and knowledge completion by LLMs².

2 Related Work

2.1 Entity Linking in Scientific Table

Table analysis is crucial for extracting experimental information and results in information extraction from scientific papers. For instance, Axcell (Kardas et al., 2020) extracts tables from the LATEX source of papers and performs linking of table cells to entities in a knowledge base. Similarly, S2abEL (Lou et al., 2023) constructs a dataset annotated with entities linked to cells, along with the type of information and the source references for that information, for comparable tasks. Axcell and S2abEL use features representing table cells, such as the cell's positional information and text from the main body that matches the cell's text. SciREX (Jain et al., 2020) and CitationIE (Viswanathan et al., 2021) aim to extract information from the entire paper, not only tables. In these works, the entire document is converted into a feature. However, relevant descriptions of specific table cells are

scattered throughout the document. Therefore, it is necessary to efficiently extract the contexts of the cells from the document. 115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

2.2 Data Augmentation/Synthetic Data

Data augmentation and synthesis using LLMs are employed in various tasks. Wang et al. (2021) employed few-shot learning to create training data from labels, and Josifoski et al. (2023) generated synthetic training data by solving inverse problems, taking advantage of the imbalance in task difficulty. Although these are effective methods for problems that are difficult to solve directly with LLM's zeroshot or few-shot capabilities, synthesizing tables or papers from entities to be linked is challenging. Therefore, this study aims to enhance the learning efficiency of current human-annotated data by utilizing synthetic data generated with LLMs.

3 Entity Linking in Scientific Tables

EL for scientific tables aims to map each table cell within a paper to an entity in a KB (PwC in our experiments) or "OutKB" if no corresponding entity is found in the KB. S2abEL divides this task into the following four subtasks:

- 1. Cell Type Classification (CTC): Classifying table cells into five types: *method*, *dataset*, *metric*, *dataset&metric*, and *other*. e.g., The cell "CoLA" is classified as *dataset*.
- 2. Attributed Source Matching (ASM): Identifying the *attributed source(s)* for a table cell within a paper. The attributed source(s) is the reference paper that originally proposed or introduced the concept that a target cell refers to. This step aims to distinguish similar surface forms to find the correct referent entities in the subsequent subtasks. e.g., The paper by Warstadt et al. (2019) is identified as the attributed source of "CoLA".
- 3. Candidate Entity Retrieval (CER): Retrieving candidate entities from the KB that are likely to be linked to a target cell. The purpose of this step is to exclude unlikely candidates to reduce computational costs. e.g., The entity "/dataset/cola" associated with the paper by Warstadt et al. (2019) in PwC is added to the candidates for "CoLA."
- 4. Entity Disambiguation (ED): Selecting the referent entity from the candidates for a given cell (or assign *OutKB* if none of them is appro-

²We will release our code and data to reproduce the experiments.



Figure 2: Generation of synthetic context: Sentences related to a particular cell (raw context) are provided to a Large Language Model (LLM). By having the LLM explain the content of the target cell, contextual information related to the target cell is extracted from the raw context.

priate). e.g., "/dataset/cola" is selected from the entity candidates.

4 Method

163

164

165

166

167

169

170

171

172

173

174

175

176

178

179

187

190

191

192

195

196

197

198

In this paper, we approach the EL task by performing these subtasks. However, CTC has already surpassed 90% accuracy in prior research, and replacing CTC predictions with the correct types contributes less than a 1% improvement in the final EL accuracy. Therefore, this study does not focus on CTC; instead, we use the correct cell types to proceed to the subsequent subtasks.

In this section, we first improve the rule-based method for context extraction to collect broad contexts and introduce data synthesis to the context. Then, we show how to apply them to the EL subtasks.

4.1 Supplementing Context Information

Meticulous cell context extraction: For EL, a model needs to interpret the concept that a cell text represents and extract appropriate contexts for it from the paper. As context information of a target cell, prior research has utilized various features, including sentences retrieved by BM25 (Robertson and Zaragoza, 2009), the cell's position in the table, and the surrounding cells. However, the retrieval method can miss relevant sentences or extract irrelevant sentences due to text fluctuation (e.g., abbreviation, paraphrasing) or the use of general words, resulting in insufficient and erroneous information sourced from the main text. To alleviate this, we first collect text fragments covering broader contexts. Specifically, we use the following features as the contexts for a target cell: (i) The cell's text. (ii) The cell type. (iii) The table caption. (iv) Sentences referring to the table: Sentences that explicitly contain references to the table, such as "Table 1." (v)

Sentences containing the cell's text. We refer to a set of the features as the *raw context* for a target cell.

200

201

202

203

204

205

206

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

227

229

230

231

232

233

For example, the raw context of a cell in the paper by Devlin et al. (2019) illustrated in Figure 1 is as follows: (i) *CoLA*. (ii) *dataset*. (iii) *"Table 1: GLUE Test results, scored by the evaluation server* [...]". (iv) *"Results are presented in Table 1"*. (v) *"CoLA The Corpus of Linguistic Acceptability is a binary single-sentence* [...]".

Synthetic context generation: To focus on essential information in a raw context and supplement it by injecting external knowledge, we employ an LLM to generate a description for a cell based on the corresponding raw context. We refer to a description generated by an LLM as a *synthetic context*. The process of synthetic context generation in this study is illustrated in Figure 2. In this research, we employ OpenAI's GPT-4 Turbo (1106) as the LLM.

For example, the synthetic context generated by the LLM for the example shown in Figure 1 is "CoLA stands for Corpus of Linguistic Acceptability. It is a dataset used for a binary single-sentence classification task in natural language processing. [...]". This exhibits the LLM's capabilities of providing a synthetic context that summarizes adequate information for EL.

4.2 Subtasks of Entity Linking

The methodology for EL used in this study is illustrated in Figure 3, and the details for each component are explained as follows:

Attributed Source Matching: We follow the approach of S2abEL (Lou et al., 2023) for ASM. The potential attributed sources for a cell are all cited papers and the current document itself. Including



Figure 3: Entity Linking: involves searching a knowledge base for entities associated with a specific table cell, taking the main text of the paper, references, and contextual information as input. It consists of three subtasks. ASM: Extract candidates of Attributed source for the target cell from the target paper and reference papers. CER: Extract candidates of Entiry for the target cell from the attributed source candidates. The number of candidates is k. ED: Select the entity to be linked to the target cell.

the document itself is necessary for the case where the cell's referent concept is newly proposed in the document. To find the attributed source from the potential source, we calculate the relevance scores between the cell and each potential source. As the features for a scoring model, we concatenate the title and abstract of a potential source and the cell's context. As the scoring model, S2abEL adopts SciBERT (Reimers and Gurevych, 2019). In this research, we employ GPT-2 (Radford et al., 2019) with synthetic context and SciBERT with raw context. We investigate the combination effect of the scoring models and contexts in §5.2.1³. A scoring model is trained with binary cross-entropy loss.

236

237

241

242

243

245

246

249

250

251

Candidate Entity Retrieval: The candidate entity set for a target cell is constructed by using its attributed sources. We sort the attributed sources obtained by the ASM model by their ASM scores and then retrieve entities for each attributed source from 253 the KB until we obtain k entity candidates for the 254 cell. In S2abEL, dense retrieval (DR) (Karpukhin 255 et al., 2020) is also employed to add more entities 256 to the candidate set. However, the experimental 257 results in S2abEL reported that the ASM-based retrieval method without using DR achieves over 259 90% recall when $k \ge 30$. Therefore, we retrieve candidate entities using ASM results only to see 261 the effect solely of ASM on EL performance. 262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

284

285

287

289

290

291

292

293

294

295

296

297

Entity Disambiguation: Following the prior work, we employ a model to calculate the score for each entity candidate by feeding the concatenation of the entity name, its description, and the cell context into the model. We adopt SciBERT for the scoring model and train it with binary crossentropy loss. The highest-scoring entity is linked to the cell when the score is greater than a predefined threshold. Otherwise, OutKB is assigned to the cell, representing being out of KB. We set the threshold to 0.5, the same as prior research.

5 Experiments

The experiments follow the setup of S2abEL, where training, validation, and test data are created from different topics. The results below represent the cross-validation average performed on the topics included in S2abEL. This allows us to compare the generalization performance of the models without overfitting to specific topics.

5.1 Entity Linking

In EL experiments, as explained in §3, the cell types are determined using the ground truth data.

5.1.1 End-to-end Entity Linking

In end-to-end entity linking experiments, we compare raw context and synthetic context against S2abEL (Lou et al., 2023) as the baseline. Dense retrieval (DR) is leveraged in CER in the baseline, but we don't use it in raw context and synthetic context conditions. Hence, we add S2abEL without the DR condition. We apply the same context to both ASM and ED. We report three metrics: InKB accuracy, OutKB F1, and Overall accuracy. InKB hit@1 accuracy shows the hit rate at the top when an entity to be linked is present. OutKB is binary and hence evaluated by the F1 score. Overall accuracy is calculated as correct if a cell is an OutKB mention and is predicted as such, or if a cell is an

³As the token length of synthetic contexts often exceeds the input token limit of SciBERT (512), we did not employ SciBERT with synthetic context.

Method	Overall acc.	OutKB F1	InKB hit@1
	k = 50		
S2abEL (Lou et al., 2023)	58.2	71.4	33.4
S2abEL w/o DR	60.8	71.7	27.1
Raw Context	69.9	76.5	47.1
Synthetic Context	70.5	76.4	53.1
	k = 20		
S2abEL w/o DR	60.2	70.3	25.3
Raw Context	68.9	75.5	44.7
Synthetic Context	70.8	76.6	52.2

Table 1: Result of End-to-end Entity Linking

InKB mention, is indicated as InKB, and hits at top1. The number of entity candidates k retrieved in CER is set to k = 50 (the same as prior work) and k = 20.



Figure 4: Evaluation of different Contexts on variation number of candidate entity.

Result Table 1 shows that the raw context and synthetic context condition have improved overall accuracy by over 10 points compared to the S2abEL baseline in k = 50. This indicates that the necessary context information for the EL task has been successfully extracted from the main text in the raw context. When comparing synthetic context to raw context, there's a slight improvement in overall accuracy and a 6-point increase in the InKB hit@1. This demonstrates that synthetic context captures the appropriate information effectively from raw context.

314

315

316

317

318

319

320

321

323

324

325

326

327

328

329

330

331

332

333

334

335

337

338

340

341

342

343

344

345

346

347

348

349

350

352

353

354

355

356

In the results with smaller entity candidates k = 20, the difference between raw context and synthetic context has become more pronounced. This suggests that in the case of k = 50, there is almost no difference in the entity candidates extracted by raw context and synthetic context. We observe the synthetic context's Overall accuracy doesn't drop in small k compared with others in Figure 4. This indicates that even at k = 20, synthetic context succeeded in extracting entities likely to be linked.

5.1.2 Evaluating Method Combinations for ASM and ED in Entity Linking Tasks

To observe the effects of context in ASM and ED, the subtasks of EL, we compare the accuracy of EL with exhaustive combinations of contexts and subtasks. The number of entity candidates is k =20 since the difference of entity candidates is small in larger k. And Dense Retrieval is not conducted in ED.

Result Table 2 shows that using either Raw or Synthetic for only ASM or ED improves accuracy compared to using the Context defined in S2abEL. When used only for ASM, the improvement in accuracy is about 1 or 2 points, while for ED, the improvement is significant, around 6 to 7 points. When comparing conditions where only the two proposed methods are applied to ASM, it is observed that synthetic context improves the overall accuracy by 1 point, indicating that synthetic context is capable of enhancing ASM. Comparing conditions where the proposed method is applied only to ED with those where it is applied only to ASM, the conditions applied to ED show significant improvements in InKB hit@1. As a result, applying the proposed method to ED also significantly improves overall accuracy. In that case, synthetic context also achieves higher accuracy. These results suggest that both raw context and synthetic context are effective for both ASM and ED tasks and that synthetic context is capable of representing more effective contexts than raw context.

5.2 Evaluating the Impact of Context in ASM

To observe the effect of Context directly, we eval-
uate using the precision of the ASM. This experi-
ment's evaluation is conducted only on cells with
attributed source papers. In the S2abEL dataset,359
360

5

303

304

305

307

310

311

312

ASM Method	ED Method	Overall acc.	OutKB F1	InKB hit@1
S2abEL	S2abEL	60.2	70.3	25.3
Raw Context	S2abEL	61.0	70.8	27.2
Synthetic Context	S2abEL	62.2	72.0	28.8
S2abEL	Raw Context	66.8	73.8	41.4
S2abEL	Synthetic Context	67.4	73.3	46.7
Raw Context	Raw Context	68.9	75.5	44.7
Synthetic Context	Synthetic Context	70.8	76.6	52.2

Table 2: Performance Comparison of ASM and ED Method Combinations in Entity Linking Tasks

some cells don't have an attributed source, and we
filtered out these cells in this experiment. Furthermore, unlike when used as a subtask of EL, it is
evaluated based on the accuracy of selecting the
Attributed source paper rather than the Entity. The
results are evaluated based on the accuracy of the
top 1 and top 5 ranked by the score.

370

5.2.1 Variation of Context and Scoring model

In the ASM task, we calculate scores for all cited references and the target paper to select the attributed source for the cell. The model calculating the score is given the cell's context and the title and abstract of each attributed source candidate. We compare variations of this scoring model and the cell context.

GPT4 Zeroshot We leverage GPT4-Turbo (1106) to ASM in a zero-shot setting. The raw contexts, along with the titles and abstracts of all cited references, are given to GPT4-Turbo, and it infers the attributed source paper directly. For GPT4 Zeroshot, the evaluation is based solely on the top 1 accuracy due to its direct selection of the cited reference without scoring.

Cell Context Three types of context will be compared. S2abEL is the context defined in S2abEL paper (Lou et al., 2023), raw context, and synthetic context.

Scoring Model SciBERT and GPT2 will be
trained with each Context. The cell context can
be longer than BERT's max input token length
(512). Hence, GPT2, which allows more input tokens (1024), is used to capture all context. The
detailed statistics of the token number of contexts
are provided in Appendix C.

Result Table 3 shows that FineTuned with synthetic context demonstrates the highest performance in both @top1 and @top5 metrics compared to other conditions. The performance of GPT4 Turbo zero-shot learning is significantly lower than that of other models undergoing FineTuning. When using the raw context with SciBERT, the @top1 accuracy is lower than the S2abEL, but the @top5 accuracy outperforms it. To compare SciBert and GPT2 with raw context, GPT2 performs less than SciBert when using raw context despite having larger parameters. 397

398

399

400

401

402

403

404

405

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

Error Analysis and Discussion

GPT4 zero-shot A drastic tendency was observed for the GPT4 Turbo zero-shot when prompted to choose attributed sources from the cited references and the target papers simultaneously. Depending on the prompt text, the GPT4 Turbo chooses them only from the target papers or only from cited references. To mitigate this issue, the prompt text was adjusted to choose the attributed source from the cited references first and then estimate whether the concept was newly proposed. If the model estimated it as a new concept, the attributed source from cited references is discarded, and the target paper is elected. Concrete examples of each are provided in Appendix B.

SciBERT vs GPT2 Raw context aims to extract sufficient information from the main text. And sentences with little relevance to the cell can be included in the context. Therefore, it is believed that GPT2, which can use all the context, did not contribute to accuracy.

Synthetic Context Analysis When comparing synthetic context and raw context, it was found that employing synthetic context improved accuracy,

		accuracy@top1			accuracy@top5		
Cell Context	Scoring Model	all	method	dataset	all	method	dataset
GPT4 Turbo zero-shot	None	22.3	30.0	1.0	-	-	-
S2abEL (Lou et al., 2023)	SciBERT	49.3	55.3	29.8	63.2	64.9	53.7
Raw Context	SciBERT	45.0	45.3	44.7	67.3	65.6	69.1
Raw Context	GPT2	41.1	43.8	35.4	62.8	63.3	60.7
Synthetic Context	GPT2	55.6	56.5	51.8	75.7	75.7	70.8

Table 3: ASM Result of varying Cell Contexts and Scoring Model



Figure 5: Completions in the synthetic contexts: There are two types of completions in synthetic contexts. **Context Completion**: LLM understands that the cell context is the abbreviated name and explains it in the synthetic context. Clarify the referent of the cell when the referent is ambiguous. **Knowledge Completion**: If the cell's content is not explained in the paper, LLM adds supplemental information.

with two types of improvements observed: context completion and knowledge completion.

Context Completion Only the abbreviated name of a method or dataset is mentioned in the cell, but the full name and description are provided in the main text. For instance, a cell in a table is "ET Perp," and it refers to the "Perplexity achieved by the Evolved Transformer." However, the expression does not appear in the main text. Thus, baseline methods failed to link the cell to the Transformer correctly. In the synthetic context, the abbreviation context is explained, and the cell is successfully linked to the correct entity.

We found that when the referent of a cell is ambiguous on its own, such as "baseline" and "All", the context supplements this by clarifying the referent in the synthetic context. Table 3 in Clark and Gardner (2018) contains a cell named "All." The synthetic context for this cell is "[...]the term 'All' likely refers to the entire dataset of TriviaQA, which [...]", making the referent of the word "All" clear. **Knowledge Completion** There may not be sufficient descriptions for well-known methods or datasets in the main text. For example, a cell content might be LSTM, and while there are mentions of LSTM in the main text, a description of the concept might be missing. Hence, previous methods misinterpret LSTM as a new concept. In the synthetic context data, the LLM supplemented the fact that LSTM stands for Long Short-Term Memory, a type of model that allows it to be correctly associated with the correct reference.

Errors in Synthetic Context Conversely, errors were also observed due to the failure of synthetic context to understand the context accurately. Specifically, in the Paper "Commonsense for generative multi-hop question answering tasks" Bauer et al. (2018), the cell labeled "*Dev*" in Table 3 refers to the evaluation dataset's Development set. However, both the GPT4 Turbo and GPT3.5 16k models mistakenly identified it as a person name "*Devi Parikh*," who is an author of one of the cited references in the paper and the word was included in raw context but unrelated to the experiments in

477 478

- 479
- 480

481

482

483

485

486

487

488

489

491

492

493

494

495

496

497

498

499

503

504

507

508

510

511

Table 3. This implies the raw context is not sufficient to identify essential information accurately.

5.2.2 Experimental Investigation of LLM for Synthetic Context

In this experiment, we measure the impact of task scores by variations of LLMs that generate synthetic contexts. In previous experiments, GPT4 Turbo was used as the LLM. The quality of the synthetic context is believed to be influenced by the language comprehension ability and specialized knowledge of the domain of the LLMs. Therefore, in this experiment, we evaluate synthetic contexts generated by the following two models in addition to GPT4 Turbo: GPT3.5 Turbo: Compared to GPT4 Turbo, it has lower language comprehension and language refinement abilities, and the learning data is assumed to be similar. TULU2 70B+DPO (Ivison et al., 2023): An Open-Sourced model that continues learning from LLAMA2 70B using the Instruction dataset and the Direct Preference Optimization (DPO) algorithm (Rafailov et al., 2023). It demonstrates performance equivalent to GPT3.5turbo in MT-Bench and AlpacaEval. Scientific documents in the Machine Learning field are included in the training dataset, such as SciERC (Luan et al., 2018) and Qasper (Dasigi et al., 2021). Hence, we expect it to have enough knowledge to perform knowledge completion.

Result Table 4 shows that TULU v2 70B+DPO demonstrates higher accuracy than GPT3.5-16k and exhibits competitive performance to GPT4-Turbo. Regardless of the type of LLM, there is a consistent trend that cells of the method type have higher accuracy than those of the dataset type. However, this trend is more pronounced in GPT3.5-16k.

Error Analysis and Discussion We confirmed 513 that TULU2 70B+DPO has knowledge about fa-514 mous methods, datasets, and evaluation metrics 515 and performs appropriate knowledge completion. 516 This indicates that when the pre-training dataset 517 for LLM includes data related to the domain, 518 Knowledge Completion can be expected. In the 519 TULU2 70B+DPO model, outputs not based on 520 facts, known as hallucinations (Ji et al., 2023), were observed. For example, a table cell in Zhong et al. 522 (2019) is "CFC (ours)", and its attributed source 523 paper is the paper itself. However, part of the syn-524 thetic context generated by TULU2 70B+DPO is The term "CFC (ours)" in the context of the scien-526

tific paper titled "modelname for Multi-evidence Question Answering" refers to a new question answering model[...], which refers to a non-existent paper title. Suppressing such hallucinations while leveraging the knowledge of LLMs is a challenge for the future.

527

528

529

530

531

532

533

534

535

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

561

562

563

564

565

566

567

568

569

570

573

Conlusion 6

In this study, we proposed new context extraction methods from the main text for Entity Linking from table cells of scientific papers. First, we propose a rule-based context extraction method (raw context) to collect broad context from a paper. Then, we introduce the synthesis date by LLM to refine the raw context (synthetic context). By employing raw context and synthetic context, we improved the accuracy of Entity Linking by more than 10 points. By qualitative analysis, we observe LLM refines raw context by supplementing context and injecting information.

Limitations 7

Model bias. Synthetic context depends on LLMs' generative capabilities and knowledge, making it susceptible to the model's bias. This study targets only English-language papers in the machine learning domain, which may limit generalization to other languages and fields.

Model Availability. The experiments in this study were conducted using OpenAI's GPT4 Turbo 1106, GPT3.5 16k, and TULU2 70B+DPO. GPT4 Turbo and GPT3.5 are accessible via the OpenAI API, but access may be lost in the future due to model version updates. Currently, these models are supported by the Azure OpenAI API.

References

- Lisa Bauer, Yicheng Wang, and Mohit Bansal. 2018. Commonsense for generative multi-hop question answering tasks. arXiv preprint arXiv:1809.06309.
- Christopher Clark and Matt Gardner. 2018. Simple and effective multi-paragraph reading comprehension. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 845–855, Melbourne, Australia. Association for Computational Linguistics.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. ArXiv, abs/2105.03011.

	a	ccuracy@t	op1	a	ccuracy@t	op5
Generation Model	all method dataset		dataset	all	method	dataset
GPT4 Turbo	55.6	56.5	51.8	75.7	75.7	70.8
GPT3.5-16k	51.4	54.3	44.8	71.3	73.4	66.7
TULU v2 70B+DPO	54.0	55.8	49.2	74.5	74.6	72.9

Table 4: ASM Result of varying LLM models for synthetic context generation

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In North American Chapter of the Association for Computational Linguistics.

574

575

576

577

578

579

581

582

583

585

586

587

589

591

592

593

594

595

599

601

611

613

614

615

616

618

619

- Kyle Gorman, Christo Kirov, Brian Roark, and Richard Sproat. 2021. Structured abbreviation expansion in context. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 995–1005, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. 2023. Camels in a changing climate: Enhancing Im adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*.
- Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. SciREX: A challenge dataset for document-level information extraction. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7506– 7516, Online. Association for Computational Linguistics.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12).
- Martin Josifoski, Marija Sakota, Maxime Peyrard, and Robert West. 2023. Exploiting asymmetry for synthetic training data generation: SynthIE and the case of information extraction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1555–1574, Singapore. Association for Computational Linguistics.
- Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. AxCell: Automatic extraction of results from machine learning papers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8580–8594, Online. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and

Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769–6781, Online. Association for Computational Linguistics. 620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

663

- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization.
- Yuze Lou, Bailey Kuehl, Erin Bransom, Sergey Feldman, Aakanksha Naik, and Doug Downey. 2023. S2abEL: A dataset for entity linking from scientific tables. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3089–3101, Singapore. Association for Computational Linguistics.
- Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3219–3232, Brussels, Belgium. Association for Computational Linguistics.
- OpenAI. 2023. Chatgpt. https://chat.openai. com.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. 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 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.

- 709 710 711
- 712 713 714 715 716 717 718 719 721 722 723

724

725

- 707 708

- David So, Quoc Le, and Chen Liang. 2019. The evolved transformer. In International conference on machine learning, pages 5877-5886. PMLR.
- Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science.
- 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 arXiv:2307.09288.

672

679

681

685

690

701

703

704 705

706

- Vijay Viswanathan, Graham Neubig, and Pengfei Liu. 2021. CitationIE: Leveraging the citation graph for scientific information extraction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 719-731, Online. Association for Computational Linguistics.
- Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. 2021. Towards zero-label language learning. arXiv preprint arXiv:2109.09193.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments. Transactions of the Association for Computational Linguistics, 7:625–641.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners.
- Sean Yang, Chris Tensmeyer, and Curtis Wigington. 2022. TELIN: Table entity LINker for extracting leaderboards from machine learning publications. In Proceedings of the first Workshop on Information Extraction from Scientific Publications, pages 20–25, Online. Association for Computational Linguistics.
- Victor Zhong, Caiming Xiong, Nitish Shirish Keskar, and Richard Socher. 2019. Coarse-grain fine-grain coattention network for multi-evidence question answering. ArXiv, abs/1901.00603.

A **Prompt for Synthetic Context** generation

Generate Synthetic Context as a summarization task in a zero-shot setting with an LLM. Synthetic context data generation was performed using GPT4-Turbo-1106-preview, GPT3.5-turbo-16k, and TULU2-70B DPO. The parameters during generation, temperature, and top-k were set to 0 to stabilize the generation. The input to the GPT4-Turbo-1106-preview and GPT3.5turbo-16k models follows a format that embeds the text of the cell (CELL_CONTENT), the paper title (PAPER_TITLE), the abstract (PA-PER_ABSTRACT), and the paper context information (PAPER_CONTEXT). The total number of prompt tokens to generate all synthetic context is 1182k, and for completion tokens, it is 1695k. Hence, generating all synthetic context with GPT4-Turbo 1106 costs about \$170.

prompts -

system_prompt : You are a researcher in the field of machine learning. You are provided with a word that appears in a certain paper and information in the paper related to that word. Please explain the word based on the information provided.

user_prompt : *Please explain the word {CELL_CONTENT}.* The title of the paper in which this word appears is "{PA-PER_TITLE}", and the abstract is "{PA-*PER_ABSTRUCT*}". The category of this word is {CELL_TYPE}. The relevant descriptions in the text are written below. {PA-PER_CONTEXT} Please provide your answer as concisely as possible.

B **GPT4** Turbo zero shot-learning prompt

In §5.2.1 of the experiment, GPT4 Turbo was evaluated directly in a zero-shot setting for Attributed Source Matching. Specifically, the main text information of the paper and the titles and abstracts of all cited references were embedded into the following template as input. As output, the ID of the cited reference that serves as the source and a flag indicating whether it represents a novel concept proposed in the paper were obtained. If it is determined that the flag represents a concept proposed in the paper, the source is cited as SourcePaper without using the cited reference ID.

726

727

732

733

735

736

737

738

prompts -

system prompt : You are tasked with identifying the source reference of the concept indicated by the cell text in a table within a machine learning academic paper. This paper is referred to as the "Source Paper" and its cited literature as "Reference Papers". The concept indicated by the cell text in the table is either a dataset or a method, which was proposed either in the cited literature. Your task is to estimate the paper in which this concept was proposed. For making your estimation, you will be provided with the cell text of the table, the type of concept that the cell text of the table is indicating, the caption of the respective table, and descriptions in the "SourcePaper" that are relevant to the respective table. You will also be presented with potential choices which include the title and abstract each of the cited literature. Please make a selection from these options. Your response should be in the following JSON format: { "estimate_result": "ID of a ReferencePaper", "is_source": "True or False" } Please input that ReferencePaper's ID into the estimate result field. Also, if you believe that the content indicated by the cell text in the table is something newly proposed in the SourcePaper, please enter True in the is_source field.

741

742

743

744

745

747 748

752

753

754

756

757

758

C Cell Contexts Statistics

We compare the statistics of the number of tokens for the input to the model used in S2abEL and the raw context and synthetic context used in this study. Table 5 shows the mean, standard deviation, maximum, and minimum number of tokens for the entire data and the mean and standard deviation of the number of tokens when the Cell type is method or dataset. Comparing the features of S2abEL and raw context, the raw context tends to have a smaller average number of tokens and a larger standard deviation. This is because, in S2abEL, information about the position of the table and surrounding cells was used as input. In contrast, in the raw context of this study, sentences that mention the table or cells in the captions or main text are added. As a result, the number of tokens varies significantly depending on the mention in the main text, leading to a larger standard deviation. The synthetic context summarizes and complements the raw context and

consistently has fewer tokens. Furthermore, when there are many mentions in the main text, only necessary information is extracted, and when there are no mentions, information is supplemented, leading to a significantly smaller standard deviation.

761

762

763

764

765

766

767

768

769

771

772

773

774

775

776

777

778

779

780

781

D Training Details

We trained all models using the AdamW optimizer (Loshchilov and Hutter, 2019) with linear decay warm-up. All models were trained using a single 48Gb NVIDIA A6000 GPU. To train GPT-2 for ASM, We added a short prompt before cell contexts. The prompt is "Given a table cell text from an academic paper in the field of Machine Learning, classify whether the information in the cell originates from provided cited literature or other. The reference information is"

E Detailed Entity Linking Result

The Entity Linking experiments were conducted using cross-validation on 10 paper categories within the S2abEL dataset. The results across all folds are presented.

	all			method		dataset		
Name	average	std	max	min	average	std	average	std
S2abEL	522.6	221.1	1292	49	569.2	220.9	434.1	192.6
Raw Context	510.2	282.7	5044	29	499.4	268.8	530.9	306.5
Synthetic Context	385.6	98.6	2121	65	394.1	98.2	369.3	97.2

Table 5: Ce	ell Contexts	Token	number	Statistics
-------------	--------------	-------	--------	------------

parameter	GPT-2	SciBERT
learning rate	2e-5	2e-5
batch size	16	32
max token length	1024	512
warm-up ratio	10%	10%

Table 6: Training Hyperparameters

Test fold	Overall acc.	OutKB F1	InKB hit@1	Overall acc.	OutKB F1	InKB hit@1	Overall acc.	OutKB F1	InKB hit@1
img_gen	48.3	55.6	26.7	53.8	57.1	37.4	48.6	46.9	43.1
misc	71.3	83.2	1.2	80.1	87.8	37.8	86.5	92.0	74.4
mt	49.2	60.6	22.5	50.0	59.9	21.2	62.2	68.8	39.4
nli	61.4	73.4	26.6	66.4	76.9	36.6	64.8	72.7	52.3
object_det	31.2	36.8	15.8	64.7	60.5	58.7	65.1	73.6	59.5
pose_estim	65.8	77.3	29.6	84.2	95.8	67.6	70.7	82.7	41.7
qa	73.7	84.4	22.6	82.3	90.2	52.2	82.5	89.6	51.7
sem_seg	63.2	73.3	49.5	67.7	66.3	55.0	76.7	73.7	68.8
speech_rec	69.0	79.3	28.1	67.2	78.0	35.7	76.9	83.0	50.2
text_class	68.8	79.0	30.6	73.0	82.6	45.1	73.8	83.0	41.4
average	60.2	70.3	25.3	68.9	75.5	44.7	70.8	76.6	52.2