IIR: Enhancing LLMs in Multi-Document Scientific Summarization via Iterative Introspection based Refinement

Anonymous ACL submission

Abstract

The current task setting and constructed datasets for Multi-Document Scientific Summarization (MDSS) have led to a significant gap between existing research and practical applications. However, the emergence of Large Language Models (LLMs) provides us with 007 an opportunity to address MDSS from a more practical perspective. To this end, we redefine MDSS task based on the scenario that automatically generates the entire related work section, and then construct a corresponding new dataset, ComRW. We first conduct a compre-013 hensive evaluation of the performance of different LLMs on the newly defined task, and identify three common deficiencies in their ability to address MDSS task: low coverage of reference papers, disorganized structure, and high redundancy. To alleviate these three deficiencies, we propose an Iterative Introspection based Refinement (IIR) method that utilizes LLMs to generate higher-quality summaries. The IIR method uses prompts equipped with Chain-of-Thought and fine-grained operators to treat LLMs as an evaluator and a generator to evaluate and refine the three deficiencies, respectively. We conduct thorough automatic and human evaluation to validate the effectiveness of our method. The results demonstrate that the proposed IIR method can effectively mitigate the three deficiencies and improve the quality of summaries generated by different LLMs. Moreover, our exploration provides insights for better addressing MDSS task with LLMs.

1 Introduction

034

042

Multi-Document Scientific Summarization (MDSS) aims to generate a concise and condensed summary for a group of topic-relevant scientific articles. In order to meet the training demand of data-driven abstractive summarization models, the existing MDSS studies (Chen et al., 2021, 2022; Wang et al., 2023a) mainly focus on the scenario of automatically generating related work of academic Recent studies usually present the task of relation classification in a supervised perspective, and traditional supervised approaches can be divided into feature based methods and kernel methods.

Feature based methods focus on extracting and selecting relevant feature for relation classification Kambibita (2004) leverages levical, syntactic and semantic features, and feeds them to a maximum entropy model. Hendrickx et al. (2010) show that the winner of Semfval-2010 Task & used the most types of features and resources, among all participants. Nevertheless, it is difficult to find an optima feature set, since traversing all combinations of features is time-consuming for feature based methods.

To remedy the problem of feature selection mentioned above, kernel methods represent the input data by computing the structural commonness between sentences, based on carefully designed kernels. Mooney and Bunescu (2005) split sentences into subsequences and compute the similarities using the proposed subsequence kernel. Bunescu and Mooney (2005) propose a dependency tree kernel and extract information from the Shortest Dependency Path (SDP) between marked entibles. Since kernel methods require similarity computation between input samples, they are relatively computationally expensive when facing large-scale datasets.

Figure 1: Example of related work section

043

044

045

046

049

051

054

060

061

062

063

064

065

066

067

068

069

070

071

073

074

papers. When constructing the corresponding datasets, such as Multi-Xscience (Lu et al., 2020), TAD (Chen et al., 2022) and TAS2 (Chen et al., 2022), individual paragraphs of a related work section are used as gold standard summaries, and the abstract section of the target paper and the reference papers are used as input documents. Such task setting and constructed datasets have greatly advanced research on MDSS.

However, we argue that the above task setting and constructed datasets induce three drawbacks: (1) The gold standard summary is merely a paragraph of a related work section in the current task setting. However, the content and structural styles of paragraphs in different positions of the related work section vary significantly, as shown in Figure 1. Therefore, datasets built based on this task setting are prone to problems like missing context and incomplete structure. (2) The input documents of the datasets are only the abstract section of the papers. However, the information required to generate the summary may come from other sections of the papers. Therefore, incomplete input information may make it difficult to infer parts of the gold summary from the input, known as the intrinsic hallucination issue (Maynez et al., 2020; Ji et al., 2023). (3) In existing datasets, all citation markers (such as "Kambhatla (2004)" in Figure 1) are normalized to a particular symbol "@cite", making it difficult to locate different reference papers in the generated summaries. The above three drawbacks have led to a significant gap between existing

076 084 086

075

101

research on MDSS and practical applications, resulting in the neglect of content consistency and structural rationality which should be emphasized in MDSS.

Recently, Large Language Models (LLMs), such as GPT-3.5 (Ouyang et al., 2022) and GPT-4 (Achiam et al., 2023), have demonstrated remarkable capabilities in tackling numerous reasoning and text generation tasks. These capabilities offer exciting new solutions for MDSS task, that is, leveraging the powerful text generation and in-context learning (Brown et al., 2020) ability of LLMs to solve MDSS task more flexibly from a perspective closer to practical applications.

In this regard, although previous researchers (Haman and Školník, 2023; Huang and Tan, 2023; Agarwal et al., 2024; Martin-Boyle et al., 2024) have attempted to utilize LLMs to address MDSS from the perspective of practical applications, their work has only stayed at the level of qualitative analysis of LLMs. For instance, Martin-Boyle et al. (2024) use citation graphs to analyze the difference in structural complexity between humanwritten summaries and GPT-4 generated summaries. Huang and Tan (2023) discuss the role and advantages of LLMs in assisting the literature review process. However, we argue that these studies fail to provide a systematic and comprehensive evaluation of the performance of LLMs on MDSS task by constructing reasonable datasets, rendering the shortcomings of LLMs in addressing MDSS remaining unknown.

To solve the above issue, we start from the perspective of practical applications of MDSS and redefine MDSS task as given the full text of a target paper and all the reference papers cited by it as input documents, the goal is to generate the entire related work section of the target paper. Based on the definition, we construct a new dataset called ComRW, which contains 60 instances, each including a target paper, several reference papers, and a gold summary.

Based on ComRW dataset, we conduct a comprehensive evaluation of the performance of LLMs on MDSS task. Specifically, the evaluation is con-119 ducted on different closed-source LLMs and open-120 source LLMs, and compared with fully-trained 121 122 models BART (Lewis et al., 2020) and EDITSum (Wang et al., 2023a). The results reveal that al-123 though LLMs are not yet comparable to EDIT-124 Sum in terms of ROUGE (Lin, 2004) metric, both 125 BERTScore (Zhang et al., 2019) metric and human 126

evaluation results indicate that the quality of summaries generated by LLMs is higher, showcasing their strong capability in addressing MDSS task. According to the results, we also identify three major common deficiencies of LLMs in generating summaries: (1) Low Coverage of Reference Papers: LLMs tend to omit some input reference papers in the generated summaries; (2) Disorganized Structure: the structure of summaries generated by LLMs is unclear, with disorganized sub-topics; (3) High Redundancy: the summaries generated by LLMs contain much redundant or repetitive content.

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

Regarding the above three deficiencies, we further propose an Iterative Introspection based Refinement (IIR) method that utilizes LLMs to generate higher-quality summaries. Specifically, IIR divides the summary generation process into draft generation and iterative refinement stages. While the concept of iterative refinement has been widely employed in text editing (Iso et al., 2020; Awasthi et al., 2019; Schick et al., 2022), the novelty of our work lies in leveraging the powerful natural language evaluation capability (Liu et al., 2023a; Fu et al., 2023; Chiang and Lee, 2023) and instructionfollowing ability of LLMs by designing reasonable prompts. Concretely, we design prompts equipped with Chain-of-Thought (Wei et al., 2022) and finegrained operators to treat LLMs as an evaluator and a generator to evaluate and refine the three deficiencies, respectively and iteratively.

We conduct both automatic and human evaluation to validate the effectiveness of our IIR method. The results indicate that IIR method can effectively alleviate the three deficiencies of LLMs, thereby enhancing the quality of generated summaries.

Our contributions are: (1) We redefine MDSS task from the perspective of practical applications and conduct a comprehensive evaluation of the performance of LLMs on MDSS. (2) We propose IIR^1 method to mitigate the three deficiencies of LLMs in addressing MDSS task. (3) Both automatic and human evaluations validate the effectiveness and universality of our IIR method.

Task Redefinition 2

The existing task setting of MDSS and constructed datasets lead to a significant gap between existing research on MDSS and practical applications. Hence, in this paper, we redefine MDSS task from a more practical perspective. The new definition

¹The code will be released if accepted.

is: Given the full text of a target paper that needs to generate a related work section, along with the full text of all reference papers in the related work section of the target paper as input, the goal is to generate the entire related work section of the target paper.

178

179

187

190

191

192

195

196

197

198

199

203

207

210

211

212

213

214

215

216

218

219

Our new definition differs from the previous one in the following three aspects: (1) In our setting, the gold summary is the full text of the related work section, avoiding the problems of missing context and incomplete structure caused by using only paragraphs as gold summary. (2) In our setting, the input documents consist of the full texts of the target paper and reference papers, thus avoiding the intrinsic hallucination issue caused by incomplete input information. (3) We retain all citation markers within the gold summary, which facilitates precise location of different reference papers and enables us to assess content consistency of the generated summary.

3 Basic Performance Analysis of LLMs

According to the above task definition, we first construct a new dataset ComRW. The construction process and dataset analysis of ComRW are introduced in Appendix A. Please refer to Appendix A for more details.

In this section, we conduct a comprehensive evaluation of LLMs' performance on MDSS task based on ComRW dataset.

3.1 Evaluation Setup

Model Selection We test the performance on: (a) Closed-source LLMs, represented by models like GPT-3.5² (Ouyang et al., 2022), GPT-4³ (Achiam et al., 2023), and Claude 3.5⁴ (Anthropic, 2024), (b) Open-source LLMs, represented by DeepSeekv3 (Liu et al., 2024) and Llama-3.1-8B (Dubey et al., 2024). We use one-shot prompting to interact with LLMs. The prompt design strategies for LLMs are introduced in Appendix B. To effectively demonstrate the performance of LLMs, we compare them with previous fully-trained MDSS models. For this purpose, we choose the state-of-the-art MDSS model EDITSum (Wang et al., 2023a) and the widely-used pretrained text generation model BART (Lewis et al., 2020) for comparison. The detailed settings for EDITSum and BART are introduced in Appendix C.

Table 1_{77} Automatic evaluation of LLMs and other models on ComRW dataset.

Model	R-1(%)	R-2 (%)	R-L(%)	BS (%)	G-Eval
BART	42.53	11.13	40.35	84.21	1.69
EDITSum	48.51	12.11	44.42	84.79	2.01
Llama-3.1-8B	42.50	9.72	39.70	84.97	2.14
DeepSeek-v3	47.11	12.03	43.95	86.83	3.03
Claude 3.5	46.11	11.89	43.02	86.56	2.74
GPT-3.5	43.83	11.38	40.59	86.26	2.46
GPT-4	46.43	11.96	43.31	86.7	3.49

Evaluation Metrics We use ROUGE-1/2/L (R-1/R-2/R-L) (Lin, 2004) and BERTScore (BS) (Zhang et al., 2019) as the automatic metrics. We also employ a LLM-based metric G-Eval (Liu et al., 2023a), which utilizes GPT-4 with Chain-of-Thought and a form-filling paradigm to assess summary quality, with scores ranging from 1 to 5.

Furthermore, we also conduct human evaluation to ensure a more reliable and comprehensive assessment.

3.2 Evaluation Results

The result of automatic evaluation is shown in Table 1. We conclude two observations from it.

Firstly, apart from GPT-3.5 and Llama-3.1-8B, other LLMs are able to outperform BART on most metrics such as ROUGE-1/L, BERTScore, and G-Eval. However, when compared with ED-ITSum, we can find that all LLMs variants lag behind EDITSum on ROUGE metric. The bestperforming LLM variant is DeepSeek-v3, achieving ROUGE-1/2/L scores of 47.11/12.03/43.95, which show a noticeable gap compared with ED-ITSum's performance of 48.51/12.11/44.42. However, on BERTScore and G-Eval, all LLM variants surpass EDITSum. The best-performing model on BERTScore, DeepSeek-v3, achieves a score of 86.83, which exceeds EDITSum by 2.04%. Similarly, the leading model on G-Eval, GPT-4, achieves a score of 3.49, exceeding EDITSum by 1.48. The above result demonstrates that LLMs have strong zero-shot learning ability and can achieve satisfactory results on MDSS task.

Secondly, the best performing closed-source LLM is GPT-4, while the best performing open-source LLM is DeepSeek-v3. Meanwhile, DeepSeek-v3 outperforms GPT-4 on most metrics except G-Eval. This will encourage more researchers to use open-source LLMs to solve MDSS task at a lower cost and in a more flexible way.

The result of human evaluation is introduced in Appendix D. Please refer to Appendix D for the detailed human evaluation settings and results.

 $^{^{2}}$ We use the gpt-3.5-turbo-0125 variant.

³We use the gpt-4-0125-preview variant.

⁴We use the claude-3-5-sonnet-20240620 variant.



Figure 2: The framework of our IIR method.

3.3 Deficiencies of Summaries Generated by LLMs

Human evaluation result in Appendix D shows that LLMs tend to overlook some reference papers, resulting in low coverage of references in the generated summaries. Additionally, we also identify two other common deficiencies of LLMs: **disorganized structure** and **high redundancy**. Disorganized structure refers to the structure of summaries generated by LLMs is unclear, with disorganized sub-topics, while high redundancy refers to the summaries generated by LLMs contain much redundant or repetitive content. We provide detailed analyses of the two deficiencies in Appendix E.

4 Method

266

271

272

273

274

275

278

279

281

284

285

291

In this section, we propose Iterative Introspection based **R**efinement (**IIR**) method, which utilizes LLMs with prompt engineering to mitigate the above three deficiencies in LLM-generated summaries. IIR consists of four modules: *Key Aspects Extraction, Reference Paper Supplement, Structural Rationality Enhancement*, and *Content Succinctness Enhancement*. The framework of IIR is illustrated in Figure 2.

4.1 Key Aspects Extraction

We use the LLM-generated summary from Section 3 as the draft for further refinement. Due to the context window limitation of LLMs, only the Abstract, Introduction, and Conclusion sections are used as input in Section 3, which may cause some key information missing when summarizing. To ensure the integrity of input information during refinement, we extract the key aspects of each paper as additional input, given the limited context window of LLMs.

To this end, we refer to the scientific concept classification scheme proposed by Teufel (2010) to classify aspects of scientific articles relevant to summarization tasks into the following seven categories: *Objective, Motivation, Method, Results, Conclusion, Advantage,* and *Limitation.* Then, we employ LLMs as a Key Aspects Extractor to extract or generate statements for each aspect from every input paper. The prompt used for the Key Aspects Extractor is shown in Appendix H.2. 300

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

346

347

349

350

4.2 Reference Paper Supplement

After Key Aspects Extraction, we utilize LLMs to add the missing reference papers to the summary. In the prompt setting for interacting with LLMs, the target paper, the reference papers, and the draft are provided in the form of key-value pairs in JSON format. We adopt a Chain-of-Thought (Wei et al., 2022) based prompting method, requiring LLMs to first count the number of the input reference papers, then count those included in the draft, and compare the two to judge if they are equal. If not, the draft must be revised to include the missing reference papers. This process iterates until LLMs determine that no further modifications are needed. The prompt used for Reference Paper Supplement (Ref_Supple) is shown in Appendix H.3.

4.3 Structural Rationality Enhancement

After Ref_Supple, we take the draft obtained from it, along with key aspects of the target paper and reference papers, as input for Structural Rationality Enhancement (Struc_Enhance). We employ LLMs as an evaluator and a generator, respectively. The evaluator gives feedbacks and refinement suggestions on structural rationality of the draft, while the generator refines the draft based on the feedbacks and refinement suggestions.

In the preliminary experiment, we empirically observe that, when providing general and vague revising feedback, the generator tends to make extensive revisions to the draft, which causes two problems: First, it is difficult to track the modification trajectory of LLMs and difficult to evaluate the effectiveness of the modifications; Second, LLMs are prone to omitting some reference papers again when revising the draft, rendering the Ref_Supple step ineffective.

To address the above two problems, we design a fine-grained and controllable prompt method equipped with Chain-of-Thought and fine-grained operators for the evaluator and generator. Specif-

4

ically, we refer to the operations commonly used 351 in text editing systems (Reid and Neubig, 2022; Liu et al., 2023b), and predefine five types of possible refinement operations: Modify, Delete, Insert, Move and Merge. Details about these operations are listed in Table 7 of Appendix F. The five types of operations are applied at the sentence level and 357 each draft sentence is labeled with a unique identifier "<SENTENCE_?>". This setting guarantees the generated feedbacks and suggestions are specific and easily traceable.

> When prompting LLMs as the evaluator, we require LLMs to identify all sentences of the draft into different sub-topics, and then determine whether the division of these sub-topics is appropriate or whether they can be merged. This process helps identify structural irrationalities in the current draft and provides corresponding suggestions. The suggestions should be from the predefined operations of Table 7. The prompt for the evaluator is shown in Appendix H.4.

364

367

370

371

372

374

376

378

391

When prompting LLMs as the generator, we require LLMs to revise the draft strictly in accordance with the suggestions from the evaluator. The prompt for the generator is also shown in Appendix H.4. Finally, to prevent conflicts of sentence identifiers after different operations, the evaluator is required to give only one suggestion at a time, ensuring that there are no conflicts between suggestions. The evaluation-generation process then proceeds iteratively to continuously improve the structural rationality of the draft. The complete process is shown in algorithm 1 of Appendix.

4.4 **Content Succinctness Enhancement**

After Struc Enhance, we further take the draft from it as input for Content Succinctness Enhancement (Cont_Enhance). We employ LLMs as a content succinctness evaluator and a content succinctness generator. The evaluator needs to inspect and provide feedbacks on the corresponding three aspects of high redundancy illustrated in Appendix E.2. We also predefine three types of text editing operations: Modify, Delete, and Merge. Details of these operations are listed in Table 8 of Appendix F. Since the operations in this step are simpler than those required for Cont_Enhance, no iteration is required 396 for this step. The revision of the draft is completed in only one evaluation-generation process. The prompts for the content succinctness evaluator and generator are shown in Appendix H.5. 400

5 **Experiments**

In this section, we conduct experiments to validate the effectiveness of the proposed IIR method.

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

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

5.1 Experimental Setup

Metrics We employ the same automatic and human evaluation as in Section 3.1.

Chosen LLMs We choose GPT-4 and DeepSeekv3 as representatives of the closed-source and opensource LLMs, respectively.

Compared Prompting Method To show the superiority of our IIR method, we compare it with other LLM prompting methods. Specifically, we introduce a new direct prompting method called Single-Turn Prompt (SinTurn). SinTurn also utilizes LLMs as both the evaluator and the generator. However, it differs in that, the evaluator of SinTurn directly evaluates the six aspects of related work: Critical Analysis, Structural Rationality, Grammatical Fluency, Content Succinctness, Reference Coverage, and Content Consistency, and then it directly provide feedbacks and suggestions without predefined operations. Subsequently, the generator revises the draft based on the feedbacks and suggestions from the evaluator.

More experimental details are introduced in Appendix G.

5.2 Experimental Results

5.2.1 Automatic Evaluation

In automatic evaluation, we report the progressive performance of each step of IIR: Reference Paper Supplement (Ref_Supple), Structural Rationality Enhancement (Struc_Enhance), and Content Succinctness Enhancement (Cont_Enhance). The results of GPT-4 and DeepSeek-v3 are shown in Table 2 and 3. We have the following two observations.

(1) The compared method SinTurn fails to improve the performance of the drafts, with notable decreases across various metrics. This indicates that it is challenging for LLMs to simultaneously enhance multiple aspects that affect summary quality. Additionally, without predefined operations, the evaluator can only provide general and vague suggestions, which leads to extensive revisions and causes the quality of the revised draft drop significantly. Conversely, our IIR method addresses the three main deficiencies of LLMs through iterative introspection based refinement with predefined operations, therefore bringing substantial improvements on summary performance.

Summary Type	R-1 (%)	R-2 (%)	R-L (%)	BS (%)	G-Eval
Initial Draft SinTurn	46.43 44.17	11.96 10.36	43.31 42.16	86.7 85.85	3.49 3.08
IIR					
After Ref_Supple	46.85 (<u>↑</u> 0.42)	$12.68 (\uparrow 0.72)$	43.67 (<u>† 0.36</u>)	86.81 (<u>† 0.11</u>)	3.53 (<u>↑ 0.04</u>)
After Struc_Enhance (#1)	47.13 (<u>† 0.28</u>)	12.75 (<u>† 0.07</u>)	43.85 (<u>† 0.18</u>)	86.77 (J 0.04)	3.56 (<u>† 0.03</u>)
After Struc_Enhance (#2)	47.28 († 0.15)	$12.73 (\downarrow 0.02)$	43.96 (<u>† 0.11</u>)	86.75 (J 0.02)	$3.55 (\downarrow 0.01)$
After Struc Enhance (#3)	47.32 (10.04)	$12.76 (\uparrow 0.03)$	44.11 († 0.15)	86.74 (J 0.01)	3.58 († 0.03)
After Cont_Enhance	47.58 († 0.26)	12.68 (↓ 0.08)	44.29 († 0.18)	86.76 († 0.02)	3.56 (J 0.02)

Table 2: Automatic evaluation results on GPT-4. Structural Rationality Enhancement (Struc_Enhance) includes three iterations (#1, #2, #3).

Summary Type	R-1 (%)	R-2 (%)	R-L (%)	BS (%)	G-Eval
Initial Draft SinTurn	47.11 45.83	12.03 10.67	43.95 41.94	86.83 86.42	3.03 2.77
IIR After Ref_Supple After Struc_Enhance (#1) After Struc_Enhance (#2) After Struc_Enhance (#3) After Cont_Enhance	$47.79 (\uparrow 0.69)$ $47.87 (\uparrow 0.08)$ $47.9 (\uparrow 0.03)$ $48.04 (\uparrow 0.14)$ $48.83 (\uparrow 0.79)$	12.79 (\uparrow 0.76) 12.84 (\uparrow 0.05) 12.89 (\uparrow 0.05) 12.95 (\uparrow 0.06) 12.98 (\uparrow 0.03)	44.48 (\uparrow 0.53) 44.7 (\uparrow 0.22) 44.74 (\uparrow 0.04) 44.78 (\uparrow 0.04) 45.5 (\uparrow 0.72)	86.96 (↑ 0.13) 86.92 (↓ 0.04) 86.91 (↓ 0.01) 86.9 (↓ 0.01) 86.87 (↓ 0.03)	3.08 († 0.05) 3.15 († 0.07) 3.16 († 0.01) 3.2 († 0.04) 3.17 (↓ 0.03)

Table 3: Automatic evaluation results on DeepSeek-v3.

(2) On both GPT-4 and DeepSeek-v3, each module of IIR can enhance the performance of summary on most metrics. After Ref_Supple, the summary achieves obvious improvements in ROUGE-1/2. This is because this module supplements the missing reference papers in the summary, thus increasing the informativeness of the summary. After Struc_Enhance, the summary shows improvements over the Ref Supple module in all metrics except for BERTScore metric. When looking at each step of this module (#1, #2, #3), since the generator performs only one operation each time, the performance change before and after each iteration is minimal. After Cont_Enhance, a large increase in ROUGE-1/L can be observed compared with Struc Enhance. Finally, comparing the final refined summary to the initial draft, we observe noticable improvements on the five metrics. Specifically, the performance for GPT-4 increases by 1.15%, 0.72%, 0.98%, 0.06%, and 0.07, while for DeepSeek-v3, the increases are 1.72%, 0.95%, 1.55%, 0.04% and 0.14. The result demonstrates the effectiveness and universality of our IIR method in improving the quality of summaries generated by different types of LLMs.

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

5.2.2 Human Evaluation

We further conduct human evaluation to analyze
the impact of IIR on summary quality in a more
specific and comprehensive way.

480 Overall Performance The first human evalua481 tion compares our IIR method against SinTurn and
482 the initial draft. The evaluation settings are gener-

Table 4: Human evaluation results of different promptmethods on ComRW dataset.

Summary Type	CA	SR	GF	CS	RC	CC
Initial Draft	2.067	1.967	2.533	1.633	73.53%	2.667
SinTurn	2.5	2.467	2.0	2.133	/1.02%	2.00/
IIR	2.267	2.7	2.6	2.733	88.94%	2.6

ally the same as those of Appendix D, but differ in that the ranking score is from 3 (best) to 1 (worst). We use the summaries generated by GPT-4 for human evaluation.

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

The result is shown in Table 4. We draw three conclusions from it: (1) Comparing the initial draft with IIR, we find that IIR brings obvious improvements on Reference Coverage (RC), Structural Rationality (SR), and Content Succinctness (CS), which demonstrates the effectiveness of our method in addressing the deficiencies in summaries generated by LLMs. (2) Comparing IIR with Sin-Turn, it is evident that IIR can help achieve higher human scores in multiple aspects, indicating that our iterative introspection based refinement method is more conducive to improving summary performance than the single-turn prompting method. (3) It is worth noting that although SinTurn requires LLMs to improve Reference Coverage (RC) of the draft, the RC result is only 71.02%, which is even worse than the initial draft's 73.53%. This indicates that LLMs still struggle to understand complex instructions on multi-dimensional summary evaluation. Therefore, decomposing complex instructions into simple and specific instructions is an effective strategy to harness the power of LLMs.



Figure 3: Human evaluation of Structural Rationality Enhancement and Content Succinctness Enhancement.

Module Performance We conduct another human evaluation to analyze the effectiveness of the modules of our IIR method. We set up two sets of pairwise comparisons on Structural Rationality Enhancement (Struc_Enhance) and Content Succinctness Enhancement (Cont_Enhance). We randomly select 10 summaries generated by GPT-4 and invite three assessors with expertise in natural language processing. Take Struc_Enhance as an example, the assessors are asked to compare the two drafts, before and after operation, to determine which one is better, or choose a tie. Since the effectiveness of Reference Paper Supplement module has already been demonstrated before, it will not be repeated here.

509

510

511

512

513

515

516

518

519

521

525

530

532

533

534

535

537

538

541

543

547

548

550

551

552

The result is shown in Figure 3. We observe that the assessors have clear preferences for after-operation draft on both Struc_Enhance and Cont_Enhance. Specifically, regarding Struc_Enhance, after-operation draft obtains an average of 53.5% preference, whereas the average preference of before-operation draft is 40%. Similarly, for Cont_Enhance, the average preference of after-operation draft is 83.3%, notably higher than the 13.3% preference for before-operation draft. The above results indicate the effectiveness of our IIR method in handling deficiencies in structural rationality and content succinctness.

5.3 More Analyses on IIR

5.3.1 Analysis of Reference Paper Supplement We first count the number of modification iterations and the number of reference papers added in each iteration for each instance. The result of GPT-4 is shown in Figure 4.

We can find that, the average number of iterations for Ref_Supple is 1.12. Most instances require only one iteration of revision, with the first iteration introducing an average of 3.82 reference papers. Only nine instances require a second iteration of revision, which generally occurs when the first iteration is unsatisfactory, and the second iteration introduces an average of 1.56 reference papers. Only one instance requires a third iteration, supplementing 2 reference papers.



Figure 4: Statistical results of the number of modification iterations and reference papers added in each iteration for each instance.



Figure 5: TA results of summaries generated at different steps of IIR.

5.3.2 Analysis of Structural Rationality Enhancement

Statistical Result of TA We first define the concept of Topic Aggregation Degree (TA) to quantitatively analyze structural rationality of summaries. TA is introduced detailedly in Appendix E.1. We count TA of summaries generated at different steps of IIR and the results of GPT-4 are shown in Figure 5.

We can find that after Struc_Enhance, TA increases from 1.88 of the initial draft to 3.62. Each iteration of Struc_Enhance contributes to this improvement, with scores rising from 2.66 to 2.86, and finally to 3.62. These results indicate that our Struc_Enhance module can effectively enhance the structural rationality of summaries.

Predefined Operation Analysis We predefine five types of operations: *Modify, Delete, Insert, Move* and *Merge*, in Struc_Enhance. We now count the proportions of the five operations to clarify the modification strategy used by LLMs.

The result of GPT-4 is shown in Figure 6 (a). We can find that *Merge* operation accounts for the highest proportion at 59.62%, indicating that the primary operation taken by LLMs to improve structural rationality is merging dispersed sub-topics. The next most common operation is *Insert*, accounting for 32.69%, which is also a necessary action to make the contextual transition of the summary more coherent. The remaining three operations,



Figure 6: The proportion of predefined operations used by LLMs.

Delete, *Move*, and *Modify*, have lower proportions, suggesting that LLMs prioritize topic-level operations over sentence-level operations when enhancing structural rationality.

5.3.3 Analysis of Content Succinctness Enhancement

We also predefine three types of operations: *Modify*, *Merge*, and *Delete* in Cont_Enhance. We analyze the proportions of the three operations to clarify the modification strategy used by LLMs. The result of GPT-4 is shown in Figure 6 (b). We can find that *Modify* operation accounts for the highest proportion at 53.27%, primarily involving modifications to make sentences more concise. Besides, *Merge* operation accounts for 32.71%, which is used to merge different sentences to remove redundant information. Finally, *Delete* operation is also widely used, accounting for 14.02%, which deletes the whole redundant sentence.

6 Related Work

583

585

586

587

592

593

594

599

601

603

604

610

611

612

614

615

617

618

621

625

6.1 **Multi-Document Scientific Summarization** Multi-Document Scientific Summarization (MDSS) involves consolidating scattered information from multiple papers. Previous studies can be categorized into extractive, abstractive and LLM-based methods. Extractive methods are commonly used in the early stages, which select off-the-shelf sentences to form the summary (Hoang and Kan, 2010; Hu and Wan, 2014; Wang et al., 2018). With the advancement of deep neural networks, abstractive methods have rapidly become the dominant approach to MDSS (Chen et al., 2021, 2022; Wang et al., 2022; Moro et al., 2022; Wang et al., 2023a), which generate summaries from scratch, bringing better coherence and readability. Despite their advantages, current task setting and constructed datasets (Lu et al., 2020; Chen et al., 2022) lead to a significant gap between existing research on MDSS and practical applications. Recently, LLMs have brought new solutions to MDSS by leveraging the powerful zero-shot learning and in-context learning (Brown et al., 2020) ability. These LLM-based methods (Haman and Školník, 2023; Huang and Tan, 2023; Agarwal et al., 2024; Martin-Boyle et al., 2024) can tackle MDSS task via flexible instructions without the need for large amounts of data. However, these methods fail to provide a systematic and comprehensive evaluation of the performance of LLMs on MDSS, resulting in the shortcomings of LLMs in addressing MDSS remaining unknown, which is the objective of this paper.

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

667

668

669

670

671

672

673

674

675

676

6.2 Prompting Methods based Text Generation

LLMs exhibit a new ability of learning merely from a few demonstrations in the context, called In-Context Learning (ICL) (Brown et al., 2020; Dong et al., 2022), which brings a novel task-solving paradigm for text generation from the perspective of prompting methods. Recently, a plenty of prompting methods have been proposed to unleash more capabilities of LLMs via Chain-of-Thought (Radhakrishnan et al., 2023; Zhang et al., 2023a; Wang et al., 2023b), content plan (Narayan et al., 2021; Creo et al., 2023; You et al., 2023), iterative refinement (Zeng et al., 2023; Zhang et al., 2023b; Madaan et al., 2024), and problem decomposition (Sun et al., 2023; Khot et al., 2022). Our work differs from these prompting methods by designing prompts with Chain-of-Thought and fine-grained sentence-level operators, which ensures the modifications made by LLMs are specific, controllable and traceable, thereby contributing to a better solution for MDSS task.

7 Conclusion

In this paper, we redefine MDSS task from the perspective of practical applications, and construct a new dataset ComRW. Then, we conduct a comprehensive evaluation of the performance of LLMs on this newly defined task, and find that the summaries generated by LLMs suffer from three major deficiencies: low coverage of reference papers, disorganized structure, and high redundancy. To mitigate these deficiencies, we propose an Iterative Introspection based Refinement (IIR) method, which uses prompts equipped with Chain-of-Thought and fine-grained operators to treat LLMs as evaluators and generators to improve summary quality, respectively. Both automatic and human evaluations demonstrate that the proposed IIR method effectively alleviates these issues, resulting in higherquality summaries. Our IIR method also provides inspiration for utilizing LLMs to tackle MDSS task effectively with prompting methods.

729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 758 759 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775

776

777

778

779

781

782

Limitations

678The limitations of this paper are twofold: (1) The679constructed dataset ComRW has only 60 instances,680which cannot support more explorations of LLMs681based MDSS from the perspective of practical ap-682plications, such as instruction tuning based meth-683ods or parameter-efficient fine-tuning. (2) Our pro-684posed IIR method is somewhat complex and inflex-685ible, involving separated evaluation and regenera-686tion steps to handle different deficiencies of sum-687maries generated by LLMs, which requires great688effort in task decomposition and prompt designing.689Therefore, more flexible and efficient prompting690methods deserve exploration in the future.

References

693

701

702

703

704

706

710

711

712

713

714

715

716

717

718

719

720

721

723

724

725

727

728

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Shubham Agarwal, Issam H Laradji, Laurent Charlin, and Christopher Pal. 2024. Litllm: A toolkit for scientific literature review. *arXiv preprint arXiv:2402.01788*.
- AI Anthropic. 2024. Claude 3.5 sonnet model card addendum. *Claude-3.5 Model Card*, 3:1–8.
- Abhijeet Awasthi, Sunita Sarawagi, Rasna Goyal, Sabyasachi Ghosh, and Vihari Piratla. 2019. Parallel iterative edit models for local sequence transduction. 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 4260–4270.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Rui Yan, Xin Gao, and Xiangliang Zhang. 2022. Target-aware abstractive related work generation with contrastive learning. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 373– 383.
- Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Xiangliang Zhang, Dongyan Zhao, and Rui Yan. 2021.
 Capturing relations between scientific papers: An abstractive model for related work section generation.
 In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the

11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6068–6077.

- Cheng-Han Chiang and Hung-Yi Lee. 2023. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631.
- Aldan Creo, Manuel Lama, and Juan C Vidal. 2023. Prompting llms with content plans to enhance the summarization of scientific articles. *arXiv preprint arXiv:2312.08282*.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*.
- Michael Haman and Milan Školník. 2023. Using chatgpt to conduct a literature review. *Accountability in research*, pages 1–3.
- Cong Duy Vu Hoang and Min-Yen Kan. 2010. Towards automated related work summarization. In *Coling 2010: Posters*, pages 427–435.
- Yue Hu and Xiaojun Wan. 2014. Automatic generation of related work sections in scientific papers: an optimization approach. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1624–1633.
- Jingshan Huang and Ming Tan. 2023. The role of chatgpt in scientific communication: writing better scientific review articles. *American journal of cancer research*, 13(4):1148.
- Hayate Iso, Chao Qiao, and Hang Li. 2020. Fact-based text editing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 171–182.
- 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 Computing Surveys*, 55(12):1–38.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.

783

793

802

804

810

811

812

813

814

815

816

817

818

819

820

821

824

825

828

829

830 831

833

837

- Yutong Li, Lu Chen, Aiwei Liu, Kai Yu, and Lijie Wen. 2024. Chatcite: Llm agent with human workflow guidance for comparative literature summary. *arXiv preprint arXiv:2403.02574*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024.
 Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. G-eval: Nlg evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522.
- Yixin Liu, Budhaditya Deb, Milagro Teruel, Aaron Halfaker, Dragomir Radev, and Ahmed Hassan. 2023b. On improving summarization factual consistency from natural language feedback. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15144–15161.
- Yao Lu, Yue Dong, and Laurent Charlin. 2020. Multixscience: A large-scale dataset for extreme multidocument summarization of scientific articles. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8068–8074.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. Advances in Neural Information Processing Systems, 36.
- Anna Martin-Boyle, Aahan Tyagi, Marti A Hearst, and Dongyeop Kang. 2024. Shallow synthesis of knowledge in gpt-generated texts: A case study in automatic related work composition. *arXiv preprint arXiv:2402.12255*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1906–1919.

Gianluca Moro, Luca Ragazzi, Lorenzo Valgimigli, and Davide Freddi. 2022. Discriminative marginalized probabilistic neural method for multi-document summarization of medical literature. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 180–189. 838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

- Shashi Narayan, Yao Zhao, Joshua Maynez, Gonçalo Simões, Vitaly Nikolaev, and Ryan McDonald. 2021. Planning with learned entity prompts for abstractive summarization. *Transactions of the Association for Computational Linguistics*, 9:1475–1492.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *Preprint*, arXiv:2203.02155.
- Letian Peng, Yuwei Zhang, Zilong Wang, Jayanth Srinivasa, Gaowen Liu, Zihan Wang, and Jingbo Shang. 2024. Answer is all you need: Instruction-following text embedding via answering the question. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 459–477, Bangkok, Thailand. Association for Computational Linguistics.
- Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiūtė, et al. 2023. Question decomposition improves the faithfulness of model-generated reasoning. *arXiv preprint arXiv:2307.11768*.
- Machel Reid and Graham Neubig. 2022. Learning to model editing processes. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3822–3832.
- Timo Schick, A Yu Jane, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2022. Peer: A collaborative language model. In *The Eleventh International Conference on Learning Representations*.
- Simeng Sun, Yang Liu, Shuohang Wang, Chenguang Zhu, and Mohit Iyyer. 2023. Pearl: Prompting large language models to plan and execute actions over long documents. *arXiv preprint arXiv:2305.14564*.
- Simone Teufel. 2010. The structure of scientific articles: Applications to summarisation and citation indexing.
- Pancheng Wang, Shasha Li, Shenling Liu, Jintao Tang, and Ting Wang. 2023a. Plan and generate: Explicit and implicit variational augmentation for multidocument summarization of scientific articles. *Information Processing & Management*, 60(4):103409.

982

983

984

985

986

987

988

989

990

991

992

946

947

948

Pancheng Wang, Shasha Li, Kunyuan Pang, Liangliang He, Dong Li, Jintao Tang, and Ting Wang. 2022. Multi-document scientific summarization from a knowledge graph-centric view. In Proceedings of the 29th International Conference on Computational Linguistics, pages 6222–6233.

896

901

902

903 904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

925

926

927

928

929

930

931

933

934

935

937

938

942

945

- Yiming Wang, Zhuosheng Zhang, and Rui Wang. 2023b. Element-aware summarization with large language models: Expert-aligned evaluation and chain-ofthought method. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8640–8665.
- Yongzhen Wang, Xiaozhong Liu, and Zheng Gao. 2018. Neural related work summarization with a joint context-driven attention mechanism. In *Proceedings* of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 1776–1786.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Wang You, Wenshan Wu, Yaobo Liang, Shaoguang Mao, Chenfei Wu, Maosong Cao, Yuzhe Cai, Yiduo Guo, Yan Xia, Furu Wei, et al. 2023. Eipe-text: Evaluation-guided iterative plan extraction for long-form narrative text generation. *arXiv preprint arXiv:2310.08185*.
- Qi Zeng, Mankeerat Sidhu, Hou Pong Chan, Lu Wang, and Heng Ji. 2023. Meta-review generation with checklist-guided iterative introspection. *arXiv* preprint arXiv:2305.14647.
- Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023a. Extractive summarization via chatgpt for faithful summary generation. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023b.
 Summit: Iterative text summarization via chatgpt.
 In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10644–10657.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

A Dataset Construction and Analysis

According to the new task definition in Section 2, we first construct a new dataset ComRW. The construction process of ComRW is introduced below.

A.1 Dataset Construction

Target Papers Selection We first select target papers from ACL 2024, EMNLP 2024, and NAACL 2024, which ensures the publication dates of the papers are after the cut-off dates of the LLMs, thus

avoiding data contamination. Papers from these top conferences adhere to academic writing conventions and provide thorough reviews of references, thus having high-quality related work sections.

We manually select 60 papers as target papers. Their related work sections exhibit clear structure, moderate length, and appropriate number of references, rendering them suitable as gold summaries for our task.

Reference Papers Collection Then we identify all the references from the related work section and automatically download them using Google Scholar⁵. For references that cannot be downloaded automatically, we manually retrieve them using school library resources. This ensures that no reference paper is missed.

Content Extraction After gathering all the target papers and reference papers, we utilize PDFMINER⁶ to convert all downloaded papers from PDF to TXT format. We also develop a section extraction tool to automatically extract contents of different sections and save them in JSON files.

A.2 Dataset Analysis

Statistical Analysis The constructed dataset ComRW contains 60 instances and the statistical information of ComRW is shown in Table 5. On average, each instance includes 15.3 reference papers. The input document contains an average of 69,725.13 words, while the gold summary has an average of 477.3 words. Although ComRW has only 60 instances, the strong few-shot learning and in-context learning capabilities of LLMs enable the dataset to support a reasonable assessment of LLMs' performance on MDSS task.

Compared with previous MDSS datasets like Multi-Xscience (Lu et al., 2020), TAD (Chen et al., 2022) and TAS2 (Chen et al., 2022), ComRW significantly surpasses them in terms of the average number of reference papers, input words, and summary words. Furthermore, an analysis of the proportion of novel *n*-grams in the gold summary that do not appear in the input documents indicates that ComRW, by using the full text of papers as input, can greatly reduce the proportion of new unigrams and bigrams in the summary, thereby avoiding the problem of intrinsic hallucination. Thus,

⁵https://scholar.google.com/

⁶https://pypi.org/project/pdfminer/

Table 5: Statistical information of ComRW and other MDSS datasets.

Dataset	# Test Set	# Input Words	# Summary Words	# Reference Papers	Novel Unigrams	Novel Bigrams
Multi-Xscience	5,093	778.08	116.44	4.42	42.33%	81.75%
TAD	5,000	845	191	5.17	43.58%	83.29%
TAS2	5,000	788	126	4.8	42.62%	82.03%
ComRW	60	69,725.13	477.3	15.3	5.78%	36.58%

our dataset enables a more objective assessment of model performance.

993

995

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1011

More Analyses on ComRW Figure 7 illustrates the distribution of the number of reference papers and sub-topics in each instance for ComRW dataset. It can be observed that the number of reference papers is roughly distributed evenly between 9 and 21. Moreover, each instance in ComRW dataset contains 1 to 5 sub-topics, with an average of 2.55 sub-topics. Particularly, instances containing 2 sub-topics are the most common, with 27 intances, followed by 20 instances containing 3 sub-topics. How to effectively identify and organize reference papers according to different sub-topics will be a significant challenge to MDSS models.



Figure 7: Distribution of the number of reference papers and the number of sub-topics for ComRW.

B Prompt Design for LLMs

We use one-shot prompting (1-shot) to interact with LLMs. Given the limited context window of gpt-3.5-turbo-0125 with only 16,385 tokens, we take the Abstract, Introduction, and Conclusion1012section of each paper as input. The input of other1013LLMs is consistent with gpt-3.5-turbo-0125.1014The prompt template is shown in Figure 10.1015

1016

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

C Compared Models Setting

Since our ComRW dataset contains only 30 in-1017 stances, it lacks sufficient data for training BART 1018 and EDITSum from scratch. To address this, we 1019 consider an alternative method to generate sum-1020 maries by BART and EDITSum. Considering that 1021 current large-scale MDSS datasets, such as Multi-1022 Xscience, are constructed at the paragraph level, 1023 we first segment the ComRW dataset into indi-1024 vidual paragraphs and identify reference papers 1025 of each paragraph. The modified dataset is de-1026 noted as ComRW-Para. Then, we train BART and 1027 EDITSum on Multi-Xscience training dataset and 1028 choose the best-performing models according to 1029 their performance on Multi-Xscience validation 1030 dataset. Subsequently, we apply the trained models to generate summaries on ComRW-Para. The generated summaries are then organized in order 1033 to serve as section-level predictions for BART and 1034 EDITSum on ComRW. 1035

D Human Evaluation

We conduct human evaluation to assess the quality of summaries generated by LLMs comprehensively. We refer to the human evaluation settings from Li et al. (2024), and take into account the definitions, content and structure requirements of a well-written related work, and then set the following six aspects for human evaluation:

- Critical Analysis (CA): Whether the generated summary include proper analysis of the strengths and weaknesses of reference papers.
- **Structural Rationality** (**SR**): Whether the summary is organized by sub-topics in a coherent and structured manner, rather than simply listing different reference papers.
- **Grammatical Fluency** (**GF**): Whether the summary is fluent, with no obvious grammatical errors.

Table 6: Human evaluation of LLMs and other models.

Model	CA	SR	GF	CS	RC	CC
EDITSum	2.233	2.567	2.6	4.133	-	-
Llama-3.1-8B	3.1	3.467	4.633	3.167	61.82%	4.333
DeepSeek-v3	5.6	5.633	5.733	3.633	78.28%	4.567
Claude 3.5	5.333	5.167	5.733	4.133	72.10%	4.833
GPT-3.5	4.167	4.3	5.067	4.333	70.02%	4.567
GPT-4	5.6	5.333	5.767	3.933	73.53%	4.833

• Content Succinctness (CS): Whether the summary is concise, does not contain repetition or lengthy information, or information that is irrelevant to the topics discussed in the target paper.

1054

1055

1056

1057

1058

1059

1060

1063

1064

1066

1067

1068

1070

1071

1072

1073

1074

1075

1076

1079

1080

1081

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1098

- **Reference Coverage (RC)**: Does the summary include all the provided reference papers without any omissions.
- **Content Consistency** (**CC**): Whether the content of the summary is consistent with the input target paper and reference papers.

For Reference Coverage, the result can be calculated automatically, thus requiring no human involvement. For Reference Coverage and Content *Consistency*, we only conduct evaluation on these two aspects for summaries generated by LLMs, because EDITSum is trained on Multi-Xscience, and during training, all citation markers are normalized, rendering human evaluation infeasible for these two aspects. Regarding Content Consistency, we ask the evaluators to rank GPT-3.5, GPT-4, Claude 3.5, DeepSeek-v3, and Llama-3.1-8B from 1 (best) to 5 (worst). Models ranked 1, 2, 3, 4 and 5 receive scores of 5, 4, 3, 2, and 1 respectively. If the evaluators consider that different summaries have the same quality, they can assign them the same rank. For instance, if the rankings are 1, 2, 3, 3 and 5, then scores are 5, 4, 3, 3 and 1 respectively. For aspects other than Reference Coverage and Content Consistency, we ask the evaluators to rank all the six models from 1 (best) to 6 (worst) with scores ranging from 6 to 1 accordingly.

We ramdomly sample 10 instances from ComRW dataset for human evaluation and invite three graduate students majoring in natural language processing to conduct human evaluation. The final score is the average score of the three evaluators.

The result of human evaluation is shown in Table 6. We conclude the following four observations: (1) LLMs outperform EDITSum in aspects such as Critical Analysis, Structural Rationality, and Grammatical Fluency. This indicates that although LLMs perform worse than EDITSum in automatic evaluation, they are capable of generating better



Figure 8: Statistical result of topic aggregation degree of different LLMs.

summaries in terms of human evaluation. (2) All LLMs, except for Llama-3.1-8B, achieve closely matched scores across all aspects. DeepSeek-v3 rates highest in Critical Analysis, Structural Rationality and Reference Coverage, while GPT-4 performs best in Critical Analysis, Grammatical Fluency and Content Consistency. (3) For Reference Coverage, it can be observed that all LLMs struggle to include all the provided reference papers in the summaries. The highest Reference Coverage is only 78.28%, indicating that 21.72% of the reference papers are still omitted. The result underscores the urgency to address this issue when utilizing LLMs to address MDSS task.

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

E Deficiencies of Summaries Generated by LLMs

In this section, we provide detailed analysis on the disorganized structure and high redundancy deficiencies of LLMs.

E.1 Disorganized Structure

To qualitatively and quantitatively analyze the Structural Rationality of summaries, we first define the concept of **Topic Aggregation Degree (TA)** \mathcal{T} as follows:

$$\mathcal{T} = \frac{1}{|S|} \sum_{i=1}^{|S|} n_i / t_i \tag{1}$$

where S means the summary set, n_i and t_i denote the number of reference papers and sub-topics in the *i*-th generated summary, respectively.

Intuitively, TA measures the average number of reference papers contained within each sub-topic in the summary. This reflects the ability of a summarization model to organize reference papers into different sub-topics, where the higher the value, the stronger the ability. To count the number of subtopics, we use one-shot prompting to employ GPT-4 as the sub-topic extractor to automatically identify different sub-topics in the summary. Prompt

{
" <sentence_1>": "The exploration of text embeddings and their applications in various natural language processing (NLP) tasks has been a focal point of recent research "</sentence_1>
* <sentence_2>". "Our work, "Answer is All You Need: Instruction-following Text Embedding via Answering the Question," Introduces INBEDDER, a novel approach that leverages abstractive question answering to generate text embeddings based on user instructions.",</sentence_2>
" <sentence_3>": "This section reviews relevant literature to position our contributions within the broader context of text ambending and instruction following models."</sentence_3>
* <sentence_4>": "Early attempts at text clustering, such as those surveyed by @cite_1, laid the groundwork for understanding text data's inherent structures.".</sentence_4>
" <sentence_5>": "These methodologies, while foundational, often lacked the ability to adapt to specific user instructions or queries, a gap our work aims to bridge by providing more contextually relevant embeddings."</sentence_5>
"SENTENCE_b>": "The development of dense passage retrieval systems, as demonstrated by @cite_2, marked a significant advancement in retrieving relevant text passages for open-domain question answering.",
"-SENTENCE_7>": "Our approach builds on this foundation by not only retrieving relevant information but also encoding it in a way that aligns with specific user instructions, thereby enhancing the utility of text embeddings for
specialized tasks.", " <sentence_8>": "The introduction of Sentence-BERT (@cite_3) and SimCSE (@cite_4) represented major strides</sentence_8>
in generating semantically meaningful sentence embeddings.", " <sentence_9>": "These models, however, primarily focus on capturing general semantic relationships without</sentence_9>
explicit consideration for user-defined instructions.", " <sentence 10="">" "INBEDDER extends these models' capabilities by incorporating instruction-following</sentence>
mechanisms, thus enabling the generation of embeddings that are tailored to specific tasks as defined by the user,". * <sentence_11>*."cent works like those by @cite_§ and @cite_§ have shown the effectiveness of contrastive learning and large-scale datasets in improving text embeddings. While these approaches have significantly advanced</sentence_11>
the field, they often do not account for the nuanced requirements posed by instruction-following tasks.", " <sentence_12>": "Our model, by contrast, is specifically designed to interpret and follow user instructions, thereby offering a more targeted approach to embedding operation "</sentence_12>
" <sentence_13>": "The concept of instruction tuning, as explored in @cite_9 and @cite_10, closely aligns with our</sentence_13>
work.", " <sentence_14>": "These studies highlight the importance of aligning model outputs with user intentions, a principle that is central to INBEDDER "</sentence_14>
* <sentence 15="">*: "However, our approach distinguishes itself by focusing on the generation of text embeddings through the lens of abstractive question answering, thereby offering a novel methodology for instruction-based text</sentence>
embedding.", " <sentence_16>": "Furthermore, the advancements in large language models (LLMs), as discussed in @cite_11</sentence_16>
and @ctte_12, provide a valuable context for our work.", " <sentence_17>": "While these models have demonstrated remarkable capabilities, their application in instruction-</sentence_17>
following tasks remains an area ripe for exploration. [*] , " <sentence_18>": "INBEDDER leverages the strengths of LLMs while introducing a unique mechanism for generating instruction-specific text embeddings.",</sentence_18>
" <sentence_19>": "In summary, while existing literature has laid a solid foundation in text embedding and instruction-following models, our work introduces a novel approach that leverages abstractive question answering to</sentence_19>
generate embeddings that are not only semantically rich but also aligned with user-defined instructions.", " <sentence_20>": "By doing so, INBEDDER addresses a critical gap in the literature, offering a new pathway for</sentence_20>
the development of user-oriented embedding models."

Figure 9: An example of the summary generated by LLMs (The target paper is from Peng et al. (2024)).

of the sub-topic extractor is shown in Appendix H.1. Through preliminary experiments, we find that GPT-4 can effectively identify different sub-topic groups in the summary, making it a reliable sub-topic extractor.

Then we use the sub-topic extractor to count TA of different LLMs and the gold summary, and show the result in Figure 8. Notably, the average TA of the gold summary is 5.87, indicating that the reference papers are effectively organized and summarized into different sub-topics, which is a necessary attribute for a well-written related work. In contrast, the average TA of the summaries generated LLMs is only $1.41 \sim 2.25$. This suggests that most sub-topics are supported by only one or two reference papers, or in some cases, no sub-topics at all, resulting in a simple enumeration of reference papers.

To illustrate this, we present an example of the summary generated by GPT-4 in Figure 9. For the convenience of showing the text fragments belonging to different reference papers, the summary in Figure 9 is divided into sentences and displayed in JSON format, where "*<SENTENCE_?*" represents the sentence identifier, and citation markers are highlighted in green shading. From the figure, we can see that the summary generated by GPT-4 simply introduces the reference papers in the order of input, without summarizing a clear topic structure. In fact, the two reference papers "@*cite_11*" and "@*cite_12*" in sentence "*<SEN*-

Algorithm 1 Structural Rationality Enhancement based on Iterative Introspection of LLMs

Input: Target Paper \mathcal{T} , Reference Papers \mathcal{D} ,
Draft from last step S_0 , Evaluator $E(\cdot)$,
Generator $G(\cdot)$, Predefined Operations \mathcal{C} =
$\{Modify, Delete, Insert, Move, Merge\}$
Output: Draft after <i>n</i> steps of Structural Rational-
ity Enhancement \mathcal{S}_n
1: for $i = 1$ to n do
2: Obtain feedbacks and suggestions $g \leftarrow$
$E(\mathcal{T}, \mathcal{D}, \mathcal{S}_{i-1})$, where $g \in \mathcal{C}$

3: Refined draft
$$S_i \leftarrow G(\mathcal{T}, \mathcal{D}, S_{i-1}, g)$$

4: **end for**

TENCE_16>" belong to the category of "*instruction tuning*", which can be described together with the reference paper "@*cite_9*" and "@*cite_10*" in sentence "*<SENTENCE_13>*". This indicates that existing LLMs, even the most powerful ones like GPT-4, have obvious shortcomings in organizing sub-topics in MDSS task.

E.2 High Redundancy

The summary generated by LLMs also exhibits high redundancy, manifested in the following two aspects: (1) **Repetition of introducing own work**. Taking the summary in Figure 9 as an example, in "*<SENTENCE_2>*" and "*<SENTENCE_19>*", the contribution of the target paper is redundantly expressed as "*introduces a novel approach that leverages abstractive question answering to generate text embeddings based on user instructions*". (2) **Generation of unnecessary title information**, as shown in "*<SENTENCE_2>*" in Figure 9.

F Predefined Text Editing Operations

The five types of predefined text editing operations1187used in Structureal Rationality Enhancement is1188shown in Table 7. And the five types of predefined1189text editing operations used in Content Succinct-1190ness Enhancement is shown in Table 8.1191

G Experimental Details of IIR

The experiments of IIR are also conducted on the1193ComRW dataset. We use the summaries generated1194by LLMs of Section 3 as the initial draft. Addi-1195tionally, we set the number of iteration steps n for1196Structural Rationality Enhancement to 3 based on1197preliminary experiment.1198

Operation Type	Instruction Template
Modify	"Modify the sentence <sentence_?> to include information"</sentence_?>
Delete	"Delete the sentence <sentence_?>"</sentence_?>
Insert	"Insert a new sentence about between the position of sentence <sentence_n> and <sentence_m>"</sentence_m></sentence_n>
Move	"Move sentence <sentence_?> before sentence <sentence_n>, then slightly Modify sentence <sentence_?> and <sentence_n> to make them contextual coherent"</sentence_n></sentence_?></sentence_n></sentence_?>
Merge	"Merge different sub-themes,, into a unified theme by putting their sentences together, then slightly revise the sentences of the theme to make them contexutal coherent and reduce fragmentation"

Table 7: Predefined text editing operations for Structural Rationality Enhancement

Table 8: Predefined text editing operations for Content Succinctness Enhancement

Operation Type	Instruction Template
Modify	"Modify the sentence <sentence_?> to exclude information about"</sentence_?>
Delete	"Delete the sentence <sentence_?>"</sentence_?>
Merge	"Merge different sentences <sentence_?>,,<sentence_?> into a single</sentence_?></sentence_?>
	sentence <sentence_?> to make them more concise."</sentence_?>

H Prompt Templates

1199

1203

1204

1206

1207

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

In this section, we list the prompt templates usedthroughout this paper.

1202 H.1 Prompt for Sub-topic Extractor

The prompt for our sub-topic extractor is shown in Figure 11 and Figure 12.

1205 H.2 Prompt for Key Aspects Extractor

The prompt for our Key Aspects Extractor is shown in Figure 13.

H.3 Prompt for Reference Paper Supplement

The prompt for Reference Paper Supplement is shown in Figure 14.

H.4 Prompt for Structural Rationality Enhancement

The prompt for structural rationality evaluator is shown in Figure 15 and Figure 16. The prompt for structural rationality generator is shown in Figure 17.

H.5 Prompt for Content Succinctness Enhancement

The prompt for content succinctness evaluator is shown in Figure 18. And the prompt for content

succinctness generator is shown in Figure 19 and Figure 20.

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

I Case Study

We provide a case study to clearly demonstrate the effects of the three steps of IIR in improving the summary quality. Figure 21, Figure 22, Figure 23, and Figure 24 correspond to the initial draft, the summary after Reference Paper Supplement, the summary after Structural Rationality Enhancement, and the summary after Content Succinctness Enhancement, respectively. We also summarize the modifications made by the three steps of IIR in Table 9.

Comparing Figure 21 and Figure 22, we can see 1234 that Reference Paper Supplement step can effec-1235 tively identify the missing reference papers in the 1236 initial draft and add them into the summary. Com-1237 paring Figure 22 and Figure 23, we can see that the draft after Reference Paper Supplement merely 1239 lists the reference papers in the summary with an in-1240 coherent context and dispersed sub-topics. For this 1241 reason, our Structural Rationality Enhancement 1242 step inserts transitional sentences between differ-1243 ent sub-topics to make the transition smoother and 1244 merges different sub-topics effectively to enhance 1245 the inherent cohesion and organizational coherence of the summary. Comparing Figure 23 and Figure 1247

Table 9: Modifications of different steps of IIR.

Step	Modification
Reference Paper	• Insert a new sentence <sentence_17>, describing reference paper @cite_5</sentence_17>
Supplement	❷ Insert a new sentence <sentence_18>, describing reference paper @cite_8</sentence_18>
	• Insert a new sentence <sentence_19>, describing reference paper @cite_9</sentence_19>
Structural	• Insert a new sentence about transition from traditional methods to neural network
Rationality	based methods before sentence <sentence_9></sentence_9>
Enhancement	Modify sentence <sentence_9> to make contextual conherence</sentence_9>
	Merge different sub-topics of <sentence_10><sentence_19> into a</sentence_19></sentence_10>
	unified sub-topic "neural network based method"
Content	• Delete the title information of sentence <sentence_2></sentence_2>
Succinctness	Delete sentence <sentence_6></sentence_6>
Enhancement	• Merge different sentences: <sentence_7> and <sentence_8>, and simplify</sentence_8></sentence_7>
	the description of @cite_1
	Delete sentences <sentence_20> and <sentence_21></sentence_21></sentence_20>

1248 24, it can be found that our Content Succinctness Enhancement step can effectively eliminate redun-1249 dant information and irrelevant content from the summary, thereby enhancing the conciseness of the generated summary.

1250 1251 1252

Imagine you are a scientific researcher and you are writing an academic paper. You have already completed the Abstract section of the target paper and have already collected the reference papers that should be included in the related work section. Now your task is to write the related work section of the target paper. Please read the target paper and the reference papers carefully, and generate the related work section according to the following steps:
#Step 1: Read the target paper and understand the main content of this paper precisely.
#Step 2: Read the reference papers one by one and identify the relationship of each reference paper and the target paper. Figure out the reason why the reference papers should be cited in the related work section. And summarize the reference papers in academic and concise manner.
#Step 3: Make sure the generated related work section fulfill the following objectives: (1) situates your work within the broader scholarly community - connects your work to the broader field and shows that your work has grown organically from current trends; (2)illustrates a "gap" in previous researches; (3) if needed, shows how you achieve the improvement compared with previous researches.
The input will be given in the following JSON format:
<pre>{ "Target Paper": { "Title": xxxx, "Abstract":xxxx, "Introduction":xxxx, "Conclusion":xxxx }, "@cite_1": { "@cite_1": { "Title": xxxx, "Abstract":xxxx, "Abstract":xxxx, "Conclusion":xxxx, "Gente_n": { "Title": xxxx, "Abstract":xxxx, "Introduction":xxxx, "Conclusion":xxxx, "Abstract":xxxx, "Introduction":xxxx, "Abstract":xxxx, "Introduction":xxxx, "Conclusion":xxxx, "Abstract":xxxx, "Introduction":xxxx, "Introduction":XXX "Introdu</pre>
"Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion". "Reference Papers" contains multiple key-value pairs, where each key is a unique citation identifier (e.g., "@cite_1",, "@cite_n"), and each value is an object representing a reference paper. For each reference paper object, the meta information of the paper is provided, including "title", "abstract", "introduction", and "conclusion".
In the above input format, "@cite_1" "@cite_n" should be the citation markers of the corresponding references, which means when you cite one reference paper, you should use "@cite_?" to represent the corresponding reference paper.
Please also remember not to leave out any given reference.
Now I will give the input as follows:

Figure 10: Prompt template for One-shot Prompting.

```
You are an expert paper reviewer. You need to list the thematic groups of the related work section.
The related work will be given in the following JSON format:
 "<SENTENCE_1>": xxxx,
"<SENTENCE_2>": xxxx,
"<SENTENCE_3>": xxxx,
}
The output should be in the following JSON format:
۲
"thematic groups":
 {
    "theme_identifier": ["<SENTENCE_?>",...,"<SENTENCE_?>"],
    "theme_identifier": ["<SENTENCE_?>",...,"<SENTENCE_?>"],

  "theme_identifier": ["<SENTENCE_?>",...,"<SENTENCE_?>"],
3
}
"thematic groups" should be a JSON object, with several key-value pairs, where the key is themetic identifier and the value
is the list of the corresponding draft sentences identifier "<SENTENCE_?>".
I will first show you an example input and output:
Input:
{ "<SENTENCE_1>": "The development of effective word representations is a cornerstone of progress in natural language by capturing semantic and
processing (NLP), enabling systems to better understand and process human language by capturing semantic and
syntactic nuances.
  "<SENTENCE 2>":
                        "Early approaches to word representation often treated words as atomic units, ignoring the rich
morphological structure that many languages exhibit."
   SENTENCE 3>": "This limitation has spurred research into more sophisticated models that can account for the internal
structure of words, leading to significant improvements in various NLP tasks.",
    <SENTENCE_4>": "One line of research has focused on leveraging morphological information to enhance word
representations."
   <SENTENCE_5>": "For instance, the work by @cite_1 introduces a novel model that constructs representations for
morphologically complex words from their constituent morphemes, combining recursive neural networks (RNNs) with
neural language models to account for contextual information.",
   '<SENTENCE_6>": "This approach has shown to outperform existing word representations on word similarity tasks,
highlighting the importance of morphological awareness in word representation.",

"<SENTENCE_7>": "Similarly, @cite_4 presents a scalable method for integrating compositional morphological

representations into vector-based probabilistic language models, demonstrating substantial reductions in perplexity and
improvements in translation tasks for morphologically rich languages."
   SENTENCE 8>": "Another significant advancement in the field has been the adoption of character-level models, which
offer a way to mitigate the out-of-vocabulary (OOV) problem by composing word representations from smaller units.
  I<SENTENCE_9>": "The work by @cite_2 describes a neural language model that relies solely on character-level inputs,
employing a convolutional neural network (CNN) and a highway network over characters to produce word-level
predictions
    "<SENTENCE_10>": "This model achieves state-of-the-art performance on several languages, underscoring the
sufficiency of character inputs for language modeling.",

"<SENTENCE_11>": "@cite_5 further explores this direction by introducing a model that constructs vector

representations of words by composing characters using bidirectional LSTMs, achieving impressive results in language
modeling and part-of-speech tagging, especially in morphologically rich languages.",
"<SENTENCE_12>": "The exploration of character n-grams as a means to represent words and sentences has also
yielded promising results."
  <SENTENCE 13>": "@cite_3 introduces CHARAGRAM embeddings, which represent textual sequences through character
n-gram count vectors followed by a nonlinear transformation.",
   <SENTENCE_14>": "This simple yet effective approach surpasses more complex architectures based on character-level
RNNs and CNNs, setting new benchmarks on several similarity tasks."
  "<SENTENCE_15>": "In addition to these developments, the field has seen efforts to enrich word embeddings with
morpho-syntactic information.",
"<SENTENCE_16>": "@cite_7 presents a graph-based semi-supervised learning method for generating morpho-syntactic
lexicons, which, when used as features, improve performance in downstream tasks like morphological tagging and
probabilistic framework, where morphological priors help improve embeddings for rare or unseen words.",
"<SENTENCE_18>": "The integration of character-level information for part-of-speech tagging has been further explored
by @cite_6, which proposes a deep neural network that combines word-level and character-level representations for
enhanced accuracy in English and Portuguese."
```

Figure 11: Prompt for sub-topic extractor



I will give you the full text of an academic paper. You need to extract as much information as possible about the objective, motivation, method, experimental result, conclusion, advantage, and limitation of the paper.
The input paper will be given in the following JSON format, with five keys "title", "abstract", "introduction", "conclusion", and "other sections", which refer to the title, the Abstract section, the Introduction section, the Conclusion section and other sections, respectively. The values are the corresponding contents:
{ "title": xxxx, "abstract": xxxx, "introduction": xxxx, "conclusion": xxxx, "other sections": xxxx }
The output should also be in JSON format as follows: { "objective": (string) representing the objective of the paper, "motivation": (string) representing the motivation behind the paper, "method": (string) representing the method or approach used in the paper, "experimental result": (string) representing the results obtained in the paper, "conclusion": (string) representing the conclusion of the paper, "advantages": (string) describing the advantages or strengths of the paper, "limitations": (string) describing the limitations or weaknesses of the paper } Now I will give you the input:

Figure 13: Prompt for Key Aspects Extractor

You are a human evaluator and paper reviser. You will be given a target paper and some reference papers cited by the target paper, along with a draft related work section. Now you need to first judge whether the draft includes all the reference papers I have provided to you. If there are some reference papers not included in the draft, you need to regenerate the related work to include these missing references.
I will provide you with the draft related work, the target paper, and the reference papers in the following JSON format: $\{$
"Target Paper": {
"title": xxxx, "abstract": xxxx, "introduction": xxxx, "conclusion": xxxx }
"Reference Papers": {
"Total citation identifiers": [@cite_1, , @cite_n], "@cite_1": {
"objective": xxxx, "motivation": xxxx, "method": xxxx,
"experimental result": xxxx, "conclusion": xxxx,
"advantages": xxxx, "limitations": xxxx
}, "@cite_p":
@cite_ni. { "objective": yyyy
"motivation": xxxx, "method": xxxx,
"experimental result": xxxx, "conclusion": xxxx,
"advantages": xxxx, "limitations": xxxx
}, "Draft Related Work": xxxx, }
"Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion". "Reference Papers" are also structured as a JSON object, including "Total citation identifiers", which is a list that contains all the citation identifiers for all referenced papers (@cite_1,, @cite_n). And Each identifier (@cite_1,, @cite_n) is also a JSON object that represents an individual reference paper. For each reference paper object "@cite_n", the meta information of the paper is provided, including "objective", "motivation", "method", "experimental result", "conclusion", "advantages", and "limitations".
You need to solve this task step by step according to the following steps:
(1) Count the number of input reference papers N by counting the items of "Total citation identifiers";
(2) Count the number of cited reference papers M in the draft related work;
(3) if N > M, it means the draft related work fails to cite all the input reference papers; Then you should regenerate the related work to add all the missing reference papers.
(4) Remember that you should not simply concatenate the missing reference papers after the draft, but rather identify the relationship between the missing reference papers and the target paper, and put the missing reference papers to suitable position to make the related work contextual coherent. If the relationship is stated in the draft, then you should put the missing reference paper to the corresponding reasonable position. Otherwise, you should start a new paragraph to introduce the missing reference papers.
(5) if $N = M$, it means all the reference papers have been cited; Then you need to do nothing.
You should only output the refined related work as well as your modification operations towards the draft. The output should also be in JSON format as follows:
"Refined Related Work": xxxx, "Modification Operations": xxxx, }
I will first show you an example:



You are an expert paper reviewer. You need to evaluate the structure clarity of the related work draft written for a target paper and provide your operable instructions for improvements. Besides the related work, you will also be provided with information about the target paper as well as information about the reference papers it cites. The target paper and the reference papers as well as the related work draft are given in the following JSON format: { "Target Paper": {
"title": xxxx, "abstract": xxxx, "introduction": xxxx, "conclusion": xxxx },
"Reference Papers": {
"Total citation identifiers": [@cite_1, ..., @cite_n], "@cite_1": { "objective": xxxx, "motivation": xxxx, "method": xxxx, "experimental result": xxxx, "conclusion": xxxx, "advantages": xxxx, "limitations": xxxx }, ... "@cite_n": {
"objective": xxxx,
"motivation": xxxx, "method": xxxx, "experimental result": xxxx, "conclusion": xxxx, "advantages": xxxx, "limitations": xxxx }, "Related Work Draft": {
"<SENTENCE_1>": xxxx,
"<SENTENCE_2>": xxxx,
"<SENTENCE_3>": xxxx, ; } } "Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion". "Reference Papers" are also structured as a JSON object, including "Total citation identifiers", which is a list that contains all the citation identifiers for all referenced papers (@cite_1, ..., @cite_n). And Each identifier (@cite_1, ..., @cite_n) is also a JSON object that represents an individual reference paper. For each reference paper object "@cite_n", the meta information of the paper is provided, including "objective", "motivation", "method", "experimental result", "conclusion", "advantages", and "limitations". "Related Work Dwarfs" is the related work draft is which the keys ("CENTENCE 1", "CENTENCE 1", "CENTENCE 1", "CENTENCE 1", "CONTENCE 1", "CONTEN "Related Work Draft" is the related work draft, in which the keys ("SENTENCE_1", ... "SENTENCE_n") represent the sentences of the draft in order. You need to evaluate the structure clarity of the related work draft and give your instruction step by step according to the following steps: 1. Read the given target paper and all the reference papers, and make note of their contents. 2. Read the related work draft of the target paper. 3. List the thematic flows of the related work draft, and then check if the draft is well-written and well-organized. 4. Identify whether the organization of the reference papers is fragemented and loose. 5. Identify whether the draft contains abrupt transitions between sentences or themes. 6. Check whether there is unreasonable discourse organization in the draft. For example, the introduction of the target paper generally comes after the discussion of reference papers, rather than before introducing the reference papers.

Figure 15: Prompt for structural rationality evaluator

You should first generate the high-level thematic flows of the draft, and then point out the unreasonable text organization using sentence keys " <sentence_?>", then you should give your instructions on how to improve structure clarity.</sentence_?>
Remember when you give your instructions, you should use the following five pre-defined operations (Remove, Delete, Insert, Move_and_Modify, and Merge_and_Modify):
(1) Modify the sentence <sentence_?> to include information</sentence_?>
(2) Delete the sentence <sentence_?>.</sentence_?>
(3) Insert a new sentence about between the position of sentence <sentence_n> and <sentence_m>.</sentence_m></sentence_n>
(4) Move sentence <sentence_?> before sentence <sentence_n>, then slightly Modify sentence <sentence_?> and <sentence_n> to make them contextual coherent.</sentence_n></sentence_?></sentence_n></sentence_?>
(5) Merge different sub-themes,, into a unified theme by putting their sentences together, then slightly revise the sentences of the theme to make them contexutal coherent and reduce fragmentation.
Remember that you should only give one instruction that deals with the most prominent problem. And do not suggest delete operation on any sentence including citation identifier "@cite_n".
The output should be in the following JSON format:
{ "thematic flows":
{ "thematic name": [" <sentence_?>",,"<sentence_?>"], "thematic name": ["<sentence_?>",,"<sentence_?>"],</sentence_?></sentence_?></sentence_?></sentence_?>
 "thematic name": [" <sentence_?>",,"<sentence_?>"]</sentence_?></sentence_?>
}, "most prominent problem in text organization": xxxx, "instructions": xxxx
}
"thematic flows" should be a JSON object, with several key-value pairs, where the key is themetic name and the value is the list of the corresponding draft sentences keys " <sentence_?>".</sentence_?>
"most prominent problem in text organization": refers to the most prominent problem in text organization. There should be only one problem.
"instructions": refers to the operation from the pre-defined operations, which is used to deal with the problem.
I will first show you an example input and output: {example}

Figure 16: Prompt for structural rationality evaluator (Continued)

You are a scientist. Now you are writing the related work section of a target paper. You have already completed the related work draft and sent it to an expert reviewer for review. The reviewer reviewed your draft carefully and gave his feedback on the structure clarity of your draft and gave the instructions on how to improve the structure clarity and coherence. You need to revise your related work draft based on the target paper, the reference papers it cites, the draft, as well as the instructions from the reviewer. Please make sure you read and understand the instructions carefully. Please refer to the provided information while revising. The input includes four parts: the target paper, including its title, abstract section, introduction section and conclusion section.
 the reference papers cited by the target paper, including the objective, motivation, method, experimental result, conclusion, advantage, and limitation of each reference paper summarized by experts. (3) the related work draft (4) the feedback from the reviewer. The input information will be given in the following JSON format. Input: { "Target Paper": "title": xxxx, "abstract": xxxx, "introduction": xxxx, "conclusion": xxxx "Reference Papers": {
"Total citation identifiers": [@cite_1, ..., @cite_n], "@cite_1": objective": xxxx, "motivation": xxxx, "method": xxxx, "experimental result": xxxx, "conclusion": xxxx, "advantages": xxxx, "limitations": xxxx }, ... "@cite_n": {
"objective": xxxx, "motivation": xxxx, "method": xxxx, "experimental result": xxxx, "conclusion": xxxx, "advantages": xxxx, "limitations": xxxx }, "Related Work Draft": { "<SENTENCE_1>": xxxx, "<SENTENCE_2>": xxxx, "<SENTENCE_3>": xxxx, }, "Feedback From the Reviewer": { "thematic flows": {
 "thematic name": ["<SENTENCE_?>",...,"<SENTENCE_?>"],
 "thematic name": ["<SENTENCE_?>",...,"<SENTENCE_?>"], "thematic name": ["<SENTENCE_?>",...,"<SENTENCE_?>"] },
"most prominent problem in text organization": xxxx, "instructions": xxxx } } "Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion". "Reference Papers" is also structured as a JSON object, including "Total citation identifiers", which is a list that contains all the citation identifiers for all referenced papers (@cite_1, ..., @cite_n). And Each identifier (@cite_1, ..., @cite_n) is also a JSON object that represents an individual reference paper. For each reference paper object "@cite_n", the meta information of the paper is provided, including "objective", "motivation", "method", "experimental result", "conclusion", "advantages", and "limitations".

Figure 17: Prompt for structural rationality generator

```
You are an expert paper reviewer. You need to evaluate the succinctness of the related work draft written for a target paper. Besides
 the related work, you will also be provided with information about the target paper as well as information about the reference papers
 it cites
 The target paper and the reference papers as well as the related work draft are given in the following JSON format:
 {
"Target Paper":
  {
"title": xxxx,
    "abstract": xxxx
    "introduction": xxxx,
    "conclusion": xxxx
  },
"Reference Papers":
  {
"Total citation identifiers": [@cite_1, ... , @cite_n],
   "@cite_1":
   {
"objective": xxxx,
"motivation": xxxx,
    "method": xxxx,
    "experimental result": xxxx,
   "conclusion": xxxx,
"advantages": xxxx,
"limitations": xxxx
   },
    "@cite_n":
   {
"objective": xxxx,
"motivation": xxxx,
    "method": xxxx,
    "experimental result": xxxx,
   "conclusion": xxxx,
"advantages": xxxx,
    "limitations": xxxx
    "Related Work Draft":
   {
"<SENTENCE_1>": xxxx,
"<SENTENCE_2>": xxxx,
"<SENTENCE_3>": xxxx,
;;
;
;
"Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion".
"Reference Papers" are also structured as a JSON object, including "Total citation identifiers", which is a list that contains all the
citation identifiers for all referenced papers (@cite_1, ..., @cite_n). And Each identifier (@cite_1, ..., @cite_n) is also a JSON object
that represents an individual reference paper. For each reference paper object "@cite_n", the meta information of the paper is
provided, including "objective", "motivation", "method", "experimental result", "conclusion", "advantages", and "limitations".
"Related Work Draft" is the related work draft, in which the keys ("SENTENCE_1", ... "SENTENCE_n") represent the sentences of the
draft in order.
 draft in order.
 Please evaluate the succinctness of the related work draft step by step according to the following steps:
 1. Read the given target paper and all the reference papers, and make note of their contents.
     Read the related work draft of the target paper.
 2.
3. Check succinctness of citation: you should check whether the introduction of individual reference paper includes too much details.
In general, the citing of a reference paper usually focuses on a particular aspect of "objective", "motivation", "method", "experimental
result", and "conclusion", rather than multiple aspects. The particular aspect should be the most relevant aspect to the target paper.

    So If you find the introduce to a reference includes more than one aspect, then you should point out this problem.
    Check succinctness of target paper: you should check the statements about introduction of own work in the draft to identify whether these statements contain too much redundant information. In general, in the related work, the authors should situate their

own work in the context of reference papers and claim their contribution concisely. Other redundant information or irrelevant information should be removed.
 5. Check sentence by sentence to identify whether it includes paper title. If so, then you should point out this problem.
 Your output should be in the following JSON format:
 {
"Succinctness Problem": xxxx,
 Where the value of "Succinctness Problem" is the problems about the succinctness of the draft.
 I will first show you an example input and output:
 {example}
```

Figure 18: Prompt for content succinctness evaluator

```
You are a scientist. Now you are writing the related work section of a target paper. You have already completed the related
work draft and sent it to an expert reviewer for review. The reviewer reviewed your draft carefully and gave his feedback
on the succinctness aspect of your draft. You need to revise your related work draft based on the target paper, the
reference papers it cites, the draft, as well as the feedback from the reviewer. Please make sure you read and understand
the feedback carefully. Please refer to the provided information while revising.
The input includes four parts:
(1) the target paper, including its title, abstract section, introduction section and conclusion section.
(2) the reference papers cited by the target paper, including the objective, motivation, method, experimental result,
conclusion, advantage, and limitation of each reference paper summarized by experts.
(3) the related work draft
(4) the feedback from the reviewer.
The target paper and the reference papers as well as the related work draft are given in the following JSON format:
{
"Target Paper":
 {
"title": xxxx,
  "abstract": xxxx,
"introduction": xxxx,
  "conclusion": xxxx
 "Reference Papers":
 {
"Total citation identifiers": [@cite_1, ... , @cite_n],
 "@cite_1":
  objective": xxxx,
  "motivation": xxxx,
  "method": xxxx,
  "experimental result": xxxx,
  "conclusion": xxxx,
"advantages": xxxx,
  "limitations": xxxx
 },
 ...
"@cite_n":
  {
"objective": xxxx,
  "motivation": xxxx,
  "method": xxxx,
  "experimental result": xxxx,
  "conclusion": xxxx,
"advantages": xxxx,
  "limitations": xxxx
  },
"Related Work Draft":
  "<SENTENCE_1> : xxxx,
"<SENTENCE_2>": xxxx,
"<SENTENCE_3>": xxxx,
  },
"Feedback From the Reviewer": xxxx,
}
"Target Paper" includes four key-value pairs: "title", "abstract", "introduction", and "conclusion".
"Reference Papers" are also structured as a JSON object, including "Total citation identifiers", which is a list that contains
all the citation identifiers for all referenced papers (@cite_1, ..., @cite_n). And Each identifier (@cite_1, ..., @cite_n) is
also a JSON object that represents an individual reference paper. For each reference paper object "@cite_n", the meta
information of the paper is provided, including "objective", "motivation", "method", "experimental result", "conclusion",
 "advantages", and "limitations"
"Related Work Draft" is the related work draft, in which the keys ("SENTENCE_1", ... "SENTENCE_n") represent the
sentences of the draft in order.
"Feedback From the Reviewer" includes the feedback from the reviewer on succinctness aspect of the draft.
You should improve the succinctness of the related work draft while ensuring all critical information are accurately maintained and ensure the contextual coherence. Use the information provided in "Target Paper" and "Reference Papers"
to achieve a concise vet comprehensive revision.
```

Figure 19: Prompt for content succinctness generator

You can use the following three types of operations to revise the draft (Modify, Delete, and Merge):
(1) Modify the sentence <sentence_?> to exclude information about aspect.</sentence_?>
(2) Delete the sentence <sentence_?>.</sentence_?>
(3) Merge different sentences <sentence_?>,,<sentence_?> into a single sentence <sentence_?> to make them more concise.</sentence_?></sentence_?></sentence_?>
Rememeber when you revise the related work, the following principles should be followed:
(1) Do not delete a sentence easily, unless you think it's absolutely necessary.
(2) Do not exert delete operation on any sentence including citation identifier "@cite_n".
(3) Do not remove any citation identifier "@cite_n" when you modify a sentence or merge some sentences.
(4) Merge operation should be only exerted on different sentences that introduce the same reference paper or the target paper.
(5) when you delete one sentence, the contextual coherence cannot be damaged.
Your output should include (1) your actions on how to improve succinctness, (2) the revised related work. The output should be organized in the following JSON format:
[\] "Actions": { "1": xxxx.
"2": xxxx,
}, "Revised Related Work":
۲ " <sentence_1>":{"content": xxxx, "trajectory": xxxx}, "<sentence_2>":{"content": xxxx, "trajectory": xxxx},</sentence_2></sentence_1>
 " <sentence_n>":{"content": xxxx, "trajectory": xxxx} }</sentence_n>
Where the output JSON file should include two key-value pairs: "Actions" and "Revised Related Work": Where the output JSON file should include two key-value pairs: "Actions" and "Revised Related Work": The value of "Actions" is a JSON object, the key indicates the instruction index, the value refers to the instruction. The value of "Revised Related Work" is also a JSON object, including multiple key-value pairs, where each key represents a sentence from the original related work section, and each corresponding value is an object containing two keys: "content": This key contains the revised content of the sentence, addressing the succinctness problem described in the "Succinctness Problem" key. "trajectory": This key contains information about the revision, which should be from the above pre-defined operations.
I will first show you an example input and output: {example}

Figure 20: Prompt for content succinctness generator (Continued)

witnessed significant advancements through the integration of neural network architectures and the strategic utilization of syntactic information.", "<SENTENCE_2>": "Our work, \"A Position Encoding Convolutional Neural Network Based on Dependency Tree for Relation Classification,\" introduces a novel approach that leverages the syntactic structure of dependency parse trees to enhance relation classification.", "<SENTENCE 3>": "This section reviews relevant literature that has contributed to the development of relation classification methods, particularly those that incorporate dependency structures and neural network models.", "<SENTENCE 4>": "Early attempts to improve relation classification have focused on exploiting the syntactic structure provided by dependency trees.", "<SENTENCE_5>": "@cite_2 presents a pioneering approach that utilizes the shortest path in the dependency graph between two entities as a feature for relation extraction.", "<SENTENCE 6>": "This method underscores the importance of dependency paths in capturing the essential syntactic relations necessary for identifying semantic relationships between entities." "SENTENCE 7>": "Building on this concept, @cite 1 introduces an augmented dependency path (ADP) that combines the shortest dependency path with subtrees attached to this path.", "<SENTENCE 8>": "The study demonstrates how this enriched syntactic representation, processed through a dependency-based neural network. can achieve state-of-the-art results in relation classification. I<SENTENCE 9>": "The exploration of neural network models for relation classification has led to the development of various architectures designed to capture the complex features necessary for this task.", "<SENTENCE_10>": "@cite_3 and @cite_6 both employ convolutional neural networks (CNNs) but differ in their approach to feature extraction and representation.", "<SENTENCE_11>": "While @cite_3 proposes a ranking-based classification model that leverages pairwise ranking loss to enhance performance, @cite 6 focuses on extracting both lexical and sentence-level features directly from word tokens using a deep CNN architecture.", "<SENTENCE 12>": "These studies highlight the versatility of CNNs in processing textual information for relation classification." '<SENTENCE_13>": "Furthermore, the significance of capturing sequential and structural information along dependency paths is emphasized in @cite 4, which introduces a model based on long short-term memory networks (LSTMs) that operate along the shortest dependency paths (SDP-LSTM)." "<SENTENCE 14>": "This model effectively integrates heterogeneous information from different sources over the dependency paths, showcasing the potential of recurrent neural networks in understanding the syntactic and semantic nuances of language.", "<SENTENCE_15>": "In addition to neural network models, the role of semantic compositionality in relation classification is explored in @cite 7.", "<SENTENCE_16>": "This study presents a recursive neural network model that learns compositional vector representations for phrases and sentences, illustrating the importance of understanding the compositional meaning of language constructs for accurate relation classification.", "<SENTENCE_17>": "Our work builds upon these foundational studies by proposing a position encoding convolutional neural network (PECNN) that utilizes tree-based position features derived from dependency parse trees.", "<SENTENCE 18>": "Unlike previous approaches, our method redefines the context for relation classification by designing tree-based convolution kernels that capture both semantic and structural information provided by dependency trees.", "<SENTENCE 19>": "This innovative approach addresses the gap in existing research by enhancing word representations through the encoding of relative positions in dependency trees and by offering a more nuanced understanding of context in relation classification.", "<SENTENCE 20>": "In summary, our work is situated within a vibrant scholarly community that seeks to advance relation classification through the innovative use of syntactic information and neural network models.", "<SENTENCE_21>": "By introducing a novel method that emphasizes the structural and semantic richness of dependency parse trees, we contribute to the ongoing dialogue in the field and propose a new direction for future research." }

Figure 21: Case study: Initial draft

significant advancements through the integration of neural network architectures and the strategic utilization of syntactic information.". "<SENTENCE_2>": "Our work, \"A Position Encoding Convolutional Neural Network Based on Dependency Tree for Relation Classification," introduces a novel approach that leverages the syntactic structure of dependency parse trees to enhance relation classification.", "<SENTENCE 3>": "This section reviews relevant literature that has contributed to the development of relation classification methods, particularly those that incorporate dependency structures and neural network models.", "<SENTENCE 4>": "Early attempts to improve relation classification have focused on exploiting the syntactic structure provided by dependency trees.", '<SENTENCE_5>": "@cite_2 presents a pioneering approach that utilizes the shortest path in the dependency graph between two entities as a feature for relation extraction." "<SENTENCE 6>": "This method underscores the importance of dependency paths in capturing the essential syntactic relations necessary for identifying semantic relationships between entities.", "<SENTENCE_7>": "Building on this concept, @cite_1 introduces an augmented dependency path (ADP) that combines the shortest dependency path with subtrees attached to this path.", "<SENTENCE_8>": "The study demonstrates how this enriched syntactic representation, processed through a dependencybased neural network, can achieve state-of-the-art results in relation classification.", "<SENTENCE 9>": "The exploration of neural network models for relation classification has led to the development of various architectures designed to capture the complex features necessary for this task.", "<SENTENCE_10>": "@cite_3 and @cite_6 both employ convolutional neural networks (CNNs) but differ in their approach to feature extraction and representation.", "<SENTENCE_11>": "While @cite_3 proposes a ranking-based classification model that leverages pairwise ranking loss to enhance performance, @cite 6 focuses on extracting both lexical and sentence-level features directly from word tokens using a deep CNN architecture.", "<SENTENCE 12>": "These studies highlight the versatility of CNNs in processing textual information for relation classification." "<SENTENCE_13>": "Furthermore, the significance of capturing sequential and structural information along dependency paths is emphasized in @cite 4, which introduces a model based on long short-term memory networks (LSTMs) that operate along the shortest dependency paths (SDP-LSTM).", "<SENTENCE_14>": "This model effectively integrates heterogeneous information from different sources over the dependency paths, showcasing the potential of recurrent neural networks in understanding the syntactic and semantic nuances of language.", "<SENTENCE_15>": "In addition to neural network models, the role of semantic compositionality in relation classification is explored in @cite_7." "<SENTENCE 16>": "This study presents a recursive neural network model that learns compositional vector representations for phrases and sentences, illustrating the importance of understanding the compositional meaning of language constructs for accurate relation classification." accurate relation classification, ,
"<SENTENCE_17>": "The challenges of extracting semantic relationships amidst sparse data and entity recognition errors are addressed in @cite_5, which employs Maximum Entropy models to combine diverse lexical, syntactic, and semantic features, highlighting the potential for scalable solutions in complex relation classification scenarios.", <SENTENCE 18>": "Lastly, the establishment of a new task for multi-way classification of semantic relations between pairs of common nominals as part of SemEval-2010 is introduced in @cite_8, underscoring the community's interest in robust knowledge extraction and the importance of semantic relations in various NLP applications." "<SENTENCE_19>": "Additionally, @cite_9's development of a new kernel method for relation extraction emphasizes the ongoing need for robust methods capable of handling POS or parsing errors, particularly in challenging domains like biological corpora.", "<SENTENCE_20>": "Our work builds upon these foundational studies by proposing a position encoding convolutional neural network (PECNN) that utilizes tree-based position features derived from dependency parse trees.", "SENTENCE 21>": "Unlike previous approaches, our method redefines the context for relation classification by designing tree-based convolution kernels that capture both semantic and structural information provided by dependency trees.", "<SENTENCE_22>": "This innovative approach addresses the gap in existing research by enhancing word representations through the encoding of relative positions in dependency trees and by offering a more nuanced understanding of context in relation classification.", "<SENTENCE_23>": "In summary, our work is situated within a vibrant scholarly community that seeks to advance relation classification through the innovative use of syntactic information and neural network models.". "<SENTENCE 24>": "By introducing a novel method that emphasizes the structural and semantic richness of dependency parse trees, we contribute to the ongoing dialogue in the field and propose a new direction for future research." _____

Figure 22: Case study: Summary after Reference Paper Supplement

{
 "<SENTENCE_1>": "The burgeoning field of relation classification in natural language processing (NLP) has witnessed significant advancements through the integration of neural network architectures and the strategic utilization of syntactic information.", "<SENTENCE_2>": "Our work, \"A Position Encoding Convolutional Neural Network Based on Dependency Tree for Relation Classification," introduces a novel approach that leverages the syntactic structure of dependency parse trees to enhance relation classification.", SENTENCE 3>": "This section reviews relevant literature that has contributed to the development of relation classification methods, particularly those that incorporate dependency structures and neural network models." "<SENTENCE 4>": "Early attempts to improve relation classification have focused on exploiting the syntactic structure provided by dependency trees.", <SENTENCE_5>": "@cite_2 presents a pioneering approach that utilizes the shortest path in the dependency graph between two entities as a feature for relation extraction.", "<SENTENCE_6>": "This method underscores the importance of dependency paths in capturing the essential syntactic relations necessary for identifying semantic relationships between entities.", "<SENTENCE 7>": "Building on this concept, @cite 1 introduces an augmented dependency path (ADP) that combines the shortest dependency path with subtrees attached to this path.", SENTENCE 8>": "The study demonstrates how this enriched syntactic representation, processed through a dependency-based neural network, can achieve state-of-the-art results in relation classification.", "<SENTENCE_9>": "The transition from traditional methods to neural network models marks a pivotal evolution in the field, offering new perspectives and methodologies for tackling the complexities of relation classification.", "<SENTENCE_10>": "Advancements in relation classification methods have also been marked by the exploration of neural network models, which have been instrumental in understanding semantic compositionality and introducing new tasks and methods.", "<SENTENCE_11>": "@cite_3 and @cite_6 both employ convolutional neural networks (CNNs) but differ in their approach to feature extraction and representation, showcasing the versatility of CNNs in processing textual information for relation classification." "<SENTENCE_12>": "@cite_4 introduces a model based on long short-term memory networks (LSTMs) that operate along the shortest dependency paths (SDP-LSTM), effectively integrating heterogeneous information from different sources over the dependency paths.", <SENTENCE_13>": "@cite_7 presents a recursive neural network model that learns compositional vector representations for phrases and sentences, illustrating the importance of understanding the compositional meaning of language constructs for accurate relation classification.", "SENTENCE_14>": "Furthermore, @cite_5 addresses the challenges of extracting semantic relationships amidst sparse data and entity recognition errors by employing Maximum Entropy models to combine diverse lexical, syntactic, and semantic features.", "<SENTENCE_15>": "@cite_8 introduces a new task for multi-way classification of semantic relations between pairs of common nominals as part of SemEval-2010, underscoring the community's interest in robust knowledge extraction.", '<SENTENCE 16>": "@cite 9's development of a new kernel method for relation extraction emphasizes the ongoing need for robust methods capable of handling POS or parsing errors, particularly in challenging domains like biological corpora." "<SENTENCE_17>": "Our work builds upon these foundational studies by proposing a position encoding convolutional neural network (PECNN) that utilizes tree-based position features derived from dependency parse trees.", "<SENTENCE_18>": "Unlike previous approaches, our method redefines the context for relation classification by designing tree-based convolution kernels that capture both semantic and structural information provided by dependency trees.". "<SENTENCE 19>": "This innovative approach addresses the gap in existing research by enhancing word representations through the encoding of relative positions in dependency trees and by offering a more nuanced understanding of context in relation classification.". "<SENTENCE 20>": "In summary, our work is situated within a vibrant scholarly community that seeks to advance relation classification through the innovative use of syntactic information and neural network models.", I<SENTENCE_21>": "By introducing a novel method that emphasizes the structural and semantic richness of dependency parse trees, we contribute to the ongoing dialogue in the field and propose a new direction for future research." }

Figure 23: Case study: Summary after Structural Rationality Enhancement

{
 "<SENTENCE_1>": "The burgeoning field of relation classification in natural language processing (NLP) has witnessed
 "> the strategie utilization of syntactic
 "> the strategie utilization of syntactic significant advancements through the integration of neural network architectures and the strategic utilization of syntactic information " **"<SENTENCE** 2>": "Our work introduces a novel approach that leverages the syntactic structure of dependency parse trees to enhance relation classification.", SENTENCE 3>": "This section reviews relevant literature that has contributed to the development of relation classification methods, particularly those that incorporate dependency structures and neural network models." "<SENTENCE 4>": "Early attempts to improve relation classification have focused on exploiting the syntactic structure provided by dependency trees.", "<SENTENCE 5>": "@cite 2 presents a pioneering approach that utilizes the shortest path in the dependency graph between two entities as a feature for relation extraction, highlighting the significance of dependency paths in identifying semantic relationships.", "<SENTENCE 6>": "This method underscores the importance of dependency paths in capturing the essential syntactic relations necessary for identifying semantic relationships between entities.", "<SENTENCE_6>": "Building on this concept, @cite_1 introduces an augmented dependency path (ADP) that combines the shortest dependency path with subtrees attached to this path, achieving state-of-the-art results in relation classification through a dependency-based neural network.", "The study demonstrates how this enriched syntactic representation, processed through a <sentence 8>"· dependency-based neural network can achieve state-of-the-art results in relation classification "<SENTENCE 7>": "The transition from traditional methods to neural network models marks a pivotal evolution in the field, offering new perspectives and methodologies for tackling the complexities of relation classification.", "<SENTENCE 8>": "Advancements in relation classification methods have also been marked by the exploration of neural network models, which have been instrumental in understanding semantic compositionality and introducing new tasks and methods.", "SENTENCE 9>": "@cite_3 and @cite_6 both employ convolutional neural networks (CNNs), showcasing the versatility of CNNs in processing textual information for relation classification.", "<SENTENCE 10>": "@cite 4 introduces a model based on long short-term memory networks (LSTMs) that operate along the shortest dependency paths (SDP-LSTM), effectively integrating heterogeneous information from different sources." "<SENTENCE_11>": "@cite_7 presents a recursive neural network model that learns compositional vector representations for phrases and sentences, illustrating the importance of understanding the compositional meaning of language constructs for accurate relation classification.", '<ŠENTENCE_12>": "Furthermore, @cite_5 addresses the challenges of extracting semantic relationships amidst sparse data and entity recognition errors by employing Maximum Entropy models to combine diverse lexical, syntactic, and semantic features.", "<SENTENCE_13>": "@cite_8 introduces a new task for multi-way classification of semantic relations between pairs of common nominals as part of SemEval-2010, underscoring the community's interest in robust knowledge extraction.", "<SENTENCE_14>": "@cite_9's development of a new kernel method for relation extraction emphasizes the ongoing need for robust methods capable of handling POS or parsing errors, particularly in challenging domains like biological corpora." "<SENTENCE_15>": "Our work builds upon these foundational studies by proposing a position encoding convolutional neural network (PECNN) that utilizes tree-based position features derived from dependency parse trees.", "<SENTENCE 16>": "Unlike previous approaches, our method redefines the context for relation classification by designing tree-based convolution kernels that capture both semantic and structural information provided by dependency trees." "<SENTENCE_17>": "This innovative approach addresses the gap in existing research by enhancing word representations through the encoding of relative positions in dependency trees and by offering a more nuanced understanding of context in relation classification." relation classification through the innovative use of syntactic information and neural network models.". dependency parse trees, we contribute to the ongoing dialogue in the field and propose a new direction for future research " }

Figure 24: Case study: Summary after Content Succinctness Enhancement