Knowledge Editing through Chain-of-Thought

Anonymous ACL submission

Abstract

001

011

012

027

033

040

Knowledge Editing is a technique that updates large language models (LLMs) with new information to maintain their world knowledge. This approach avoids the need to rebuild the model from scratch, thereby addressing the high costs associated with frequent retraining. Among these, the in-context editing paradigm stands out for its effectiveness in integrating new knowledge while preserving the model's original capabilities. Despite its potential, existing in-context knowledge editing methods are often task-specific, focusing primarily on multi-hop QA tasks using structured knowledge triples. Moreover, their reliance on fewshot prompting for task decomposition makes them unstable and less effective in generalizing across diverse tasks. In response to these limitations, we propose EditCoT, a novel knowledge editing framework that flexibly and efficiently updates LLMs across various tasks without retraining. EditCoT works by generating a chain-of-thought (CoT) for a given input and then iteratively refining this CoT process using a CoT editor based on updated knowledge. We evaluate EditCoT across a diverse range of benchmarks, covering multiple languages and tasks. The results demonstrate that our approach achieves state-of-the-art performance while offering superior generalization, effectiveness, and stability compared to existing methods, marking a significant advancement in the field of knowledge updating¹.

1 Introduction

Large Language Models (LLMs) have shown remarkable performance across numerous NLP tasks in recent years. However, once an LLM has been fully trained, its parametric knowledge becomes fixed and constrained. Retraining the entire model to incorporate new information is both expensive and time-consuming. To address this challenge, the

¹Code and data available at: https://anonymous.4open.science/r/EditCoT-6E40 field of knowledge editing has emerged, offering techniques to modify or add knowledge efficiently while preserving most of the original model's parameters. Existing knowledge edit approaches can be broadly categorized into two types: parametric and non-parametric (Wang et al., 2024b). Parametric methods directly modify the model's weights to integrate new knowledge. While this approach enables precise editing, existing research indicates that such changes may negatively affect the model's original performance and interfere with previously learned knowledge that should remain unchanged (Gu et al., 2024b). In contrast, nonparametric methods, such as in-context editing, are increasingly popular because they utilize the incontext learning ability of LLMs without modifying the original model weights, thereby preserving the model's foundational capabilities and existing knowledge (Zheng et al., 2023).

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Most existing in-context knowledge editing frameworks are designed for multi-hop question answering (MQA) tasks using knowledge triples as the primary knowledge representation. Researchers construct multi-hop questions by linking multiple triples to evaluate these methods' performance (Zhong et al., 2023). For instance, Mello and PokeMQA guide LLMs in task decomposition with manually crafted examples, facilitating knowledge updates for sub-questions (Zhong et al., 2023; Gu et al., 2024a). RAE (Shi et al., 2024) constructs and then edits a knowledge graph, utilizing the retrieval and pruning of the knowledge graph to obtain contextually injected knowledge.

While effective for triple-based MQA tasks, these methods are overly specialized and struggle with other tasks and complex reasoning. For example, RAE is impractical for the LeKUBE(Wang et al., 2024b) legal dataset, where constructing a knowledge graph on its corpus is unfeasible. PokeMQA's strategy of using updated knowledge directly as answers to sub-questions works well

100

101

102

105

106

108

110

111

112

113

114

115

116

117

118

119

120

121

123

124

125

126

127

128

129

130

131

for triple-based tasks but fails with more complex tasks like DUNE (Akyürek et al., 2023). Furthermore, Mello and PokeMQA's reliance on few-shot prompting for task decomposition results in instability and reduce effectiveness across diverse tasks.

To address the limitations of existing in-context editing methods, we propose EditCoT, a framework that edits LLM's knowledge through chainof-thought (CoT). EditCoT is also an in-context editing method that does not change the parameters of LLMs used for reasoning, thereby preserving the original capacity of the backbone LLMs. Unlike previous methods that are task-specific or heavily reliant on knowledge triples, EditCoT iteratively refines the model's reasoning with retrieved updated knowledge. This flexible approach allows it to dynamically adapt its reasoning without requiring task-specific adjustments or predefined knowledge structures. Specifically, EditCoT starts by instructing the LLM to first generate an answer to the question, and then construct an initial CoT based on the question and answer. The original CoT is then iteratively refined by a CoT editor, which requires only a single training session without needing retraining when encountering new knowledge. The final answer is generated by prompting the LLM to inference based on the updated CoT.

We evaluate EditCoT on a variety of benchmarks, including the triple-based multi-hop QA dataset MQuAKE-CF-3k(Zhong et al., 2023), a free-form editing task dataset DUNE(Akyürek et al., 2023), and a domain-specific Chinese legal knowledge update dataset LeKUBE(Wang et al., 2024b). They cover different editing formats and tasks in both Chinese and English. Experimental results demonstrate that our method achieves state-of-the-art performance with better efficiency and robustness.

The contributions of this paper are as follows:

- We propose a novel framework, *EditCoT*, that edits the chain-of-thought of LLMs, enabling more effective and task-agnostic knowledge updates.
- We demonstrate our method performs well across various knowledge editing benchmarks with different editing formats and task types, exhibiting significantly better effectiveness, robustness, and generalizability compared to other baselines.

2 Related Work

2.1 Knowledge Editing

Knowledge editing methods modify or update knowledge within LLMs, categorized into para-

metric and non-parametric approaches (Wang et al., 2024b). ROME (Meng et al., 2022) is a parametric method, update the model's parameters by treating FFN layers as key-value storage for precise modifications. Non-parametric approaches leverage in-context learning (Brown et al., 2020), such as Mello (Zhong et al., 2023), which uses prompts for task decomposition, PokeMQA (Gu et al., 2024a), which enhances Mello's robustness with a scope detector, and RAE (Shi et al., 2024), which uses a knowledge graph for task-specific editing.

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

177

178

179

180

2.2 Chain-of-Thought

Chain-of-thought (CoT) has significantly enhanced LLMs by providing step-by-step reasoning for complex tasks. Wei et al. (2022) introduce CoT prompting, where manually supplied reasoning chains help LLMs generate multi-step solutions. Kojima et al. (2022) demonstrate phrases like "Let's think step by step" enable zero-shot CoT reasoning. Recent work has focused on enhancing the factuality of CoT. Wang et al. (2023b) propose a self-consistency decoding method to improve CoT reliability. Zhao et al. (2023) introduce Verify-and-Edit that involves editing the generated CoT. However, our approach significantly differs in some key aspects. The verification process of Verify-and-Edit relies on self-consistency, which is less effective for knowledge editing due to the LLM's high confidence in outdated information. Additionally, their method depends on the LLM's ability to learn from context, whereas our CoT editor provides more comprehensive adjustments. Further comparative analysis can be found in Appendix B.

3 Methodology

In this section, we present the methodology of the EditCoT framework in detail. We begin with an overview of our proposed approach, followed by a step-by-step description of the framework's workflow. We then explain the CoT editor training process in depth, concluding with a discussion of implementation details and dataset construction.

3.1 Methodology Overview

EditCoT aims to update LLMs by editing their CoT, integrating new factual knowledge while preserving their reasoning capabilities. Traditional parametric knowledge editing methods often require resource-intensive retraining or fine-tuning, which can result in the loss of existing knowledge. In contrast, EditCoT introduces an iterative editing



Figure 1: EditCoT updates LLMs by iteratively editing the CoT to integrate new knowledge while preserving its reasoning abilities. It generates an initial CoT, detects conflicts, and revises the CoT iteratively.

mechanism that efficiently incorporates updates by focusing solely on the CoT, ensuring that the model remains both accurate and logically consistent.

The framework leverages two components: an unedited version of the LLM responsible for generating the initial CoT, detecting conflict, and generating the final answer, and an editor capable of editing CoT based on the new knowledge. The following subsection details each step of EditCoT.

3.2 Workflow of EditCoT

181

182

183

185

190

191

192

193

194

195

196

199

200

201

Figure 1 illustrates the EditCoT framework's overall workflow, which is also described in pseudocode in Appendix A. The core steps are as follows:

Step 1: Generate Initial CoT by the Original LLM. At the beginning of the process, the unedited LLM generates an initial CoT in response to a query. This CoT outlines the sequence of reasoning steps leading to the model's final answer. Prioritizing the completeness of the reasoning path over the answer's accuracy, we first let the model provide the answer directly and then generate the CoT based on the question and answer.

203Step 2: Retrieve Step by Step and Detect Con-204flicts. Once the initial CoT is generated, we treat205each step of the CoT as a query and retrieve it from206the memory of edited knowledge. This memory207is typically organized in two ways: one where the208pre- and post-edit knowledge are paired together209(e.g., MQuAKE (Zhong et al., 2023)), and another

where only the post-edit knowledge is provided (e.g., DUNE (Akyürek et al., 2023)). For the former, we search the knowledge base for the pre-edit knowledge and then match the post-edit knowledge according to the corresponding relationship. For the latter, we directly search within the post-edit knowledge. Once updated knowledge is found, we proceed to the conflict detection stage. 210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

EditCoT utilizes the original LLM to detect knowledge conflicts by engaging in a three-class classification task: conflict, support, or irrelevant. Specifically, given a question, the initial CoT, and the retrieved knowledge, the classification is determined based on the probability of the prefix of the corresponding word output by the LLM. This process is crucial for identifying whether the original reasoning diverges from new, more accurate, or updated knowledge.

Step 3: Edit the chain-of-thought Once a knowledge conflict is detected, we input the question, new knowledge, and the old CoT into a CoT editor. This editor generates a new CoT based on the input. Then, our iterative process returns to Step 2, where this new CoT is treated as the old CoT for retrieval and conflict detection. The process continues until no conflicts are found, or the maximum number of iterations is reached.

Step 4: Get Final Answer Once the CoT has been fully edited and no further conflicts are detected, the final chain-of-thought is used to generate the final answer. The original LLM is responsible for this step, ensuring that the model's answer adheres to its existing inferential abilities, while also incorporating the new factual knowledge. We instruct the original LLM to pay more attention to the new CoT in the context since the new CoT often conflicts with the model's intrinsic knowledge.

3.3 CoT Editor Training

We design an automated approach to generate training data required for the CoT editor, where human annotation is unnecessary, except for the use of some human-posed questions. Our method does not rely on state-of-the-art LLMs like GPT-4(OpenAI et al., 2023) and is not tailored to specific datasets or tasks. The editor learns to edit CoT during training and can apply this skill to various tasks without retraining when encountering new knowledge. Additionally, the CoT editor is trained from the original LLM, ensuring that no knowledge distillation issues arise.



Figure 2: Diagram of the dataset construction: Our data source consists solely of questions from HotpotQA, without requiring human-annotated answers or related articles. In Part II, the blue sections represent the prefix used to guide generation, which have been generated in Part I, while red sections indicate newly generated reasoning paths.

263 264

267

271

273

275

276

279

285

3.3.1 Task Formulation

The editor's input includes a question Q, an original CoT, and new knowledge K_{new} that conflicts with the CoT. The goal is to generate a revised chain-of-thought CoT_{new} that incorporates new knowledge. Formally, this can be expressed as:

 $\operatorname{Editor}(Q, \operatorname{CoT}, K_{\operatorname{new}}) \to \operatorname{CoT}_{\operatorname{new}},$

3.3.2 Data Construction

To construct training data, we start by collecting high-quality, human-generated questions. In practice, we use the validation and test sets from the HotpotQA(Yang et al., 2018) dataset, which contains diverse multi-hop questions that often require multiple Wikipedia articles to answer. HotpotQA is selected for its diversity of question types, and multi-hop reasoning tasks generally yield more complex CoTs. After obtaining these questions, we use the following construction steps:

Part I: Retrieve, then Generate CoT and An-278 **swer** For each question, we first retrieve top-5 relevant Wikipedia paragraphs and then use a retrievalaugmented generation (RAG) method to generate the CoT and the corresponding answer. We consider this initial CoT and answer as a basic version, and in the following sections, these are highlighted as CoT and Answer for clarity. To facilitate steplevel separation, we prompt the model to indicate 286 each step with the [STEP] phrase.

Part II: Prefix-Guided CoT Generation Given a CoT with N steps generated in Part I, for each step, e.g. the kth step $(k = 1, 2, \dots, N)$, we take the first (k-1) steps as a prefix and prompt the model to generate the remaining steps and the final answer, marked as CoT k and Answer k respectively. It is important to note that this step does not use RAG, and the model is only prompted with the first (k-1) steps of the CoT as context. We finally obtain N CoTs in this part.

290

291

292

293

294

295

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

Part III: Conflict Detection We then compare each of the N CoTs generated in Part II to the original **CoT** and instruct the model to detect any knowledge conflicts. A data point is accepted into the final dataset only if: (1) there is a conflict in the final answers, and (2) there is a conflict at the first diverging step between the two CoTs.

For example, considering CoT k and CoT, the first diverging step is step k. If the original Step k and the new Step k conflict with each other, and there is also a conflict between Answer and Answerk, they are eligible to be added to the final dataset. This ensures that conflicting knowledge can be pinpointed at a specific step in the reasoning chain. Here we assume that the divergence occurs only once in most of the cases.

Part IV: Rewrite to Extract New Knowledge In this step, we rewrite the identified diverging step from the original **CoT** to explicitly extract the new knowledge. Specifically, for the example above, a complete training data entry is formed as a tuple:

(Q, CoT,	$K_{\text{new}}, \text{ CoT}_{\text{new}}) =$	31
(Question,	CoT k, rewrite(Step k), CoT).	32



Figure 3: A training data example. The black part of the two CoTs indicates a shared prefix.

We also apply regular expressions to filter out lowquality examples, such as those overly influenced by the retrieved documents in RAG. Figure 2 illustrates the overall process of dataset construction.

3.4 Implementation Details for Training

We implement conflict detection using a multiplechoice format. If the model assigns a probability of 0.6 or higher to a conflict between two answers or reasoning steps, we consider them inconsistent. During prefix-guided generation, to encourage the LLM to focus on the same topic at the diverging step, we provide the first word of the original CoT at the diverging step to the LLM. Prompt templates for dataset construction are detailed in subsection H.1. The dev and test sets of HotpotQA contain 14,810 questions. Using Meta-Llama-3-8B-Instruct(Dubey et al., 2024) and Qwen2.5-14B-Instruct(Yang et al., 2024; Team, 2024), we generate 459 and 316 training examples, respectively. Figure 3 illustrates a constructed data instance. During training, we calculate the loss based on CoT_{new} using cross-entropy loss.

4 Experimental Settings

4.1 Dataset and Metrics

We evaluate the performance of EditCoT and other baselines across the following three datasets:

347MQuAKE-CF-3k(Zhong et al., 2023)This348dataset is designed for knowledge editing using349Wikipedia-based knowledge triples, with queries350in the form of multi-hop questions. The evaluation351metric is Multihop-Accuracy, as in the original pa-352per: each instance contains 3 questions, and if any-353one is answered correctly, it is considered accurate.

We follow the settings of Mello and PokeMQA, using *Exact Match* to assess correctness.

354

355

356

357

358

360

361

362

363

364

366

367

368

369

370

371

373

374

375

376

377

378

379

380

381

382

383

387

388

389

390

391

392

394

395

396

397

398

399

400

401

402

DUNE(Akyürek et al., 2023) It includes various editing tasks, where edits are natural language instructions for the model to follow specific facts or requirements. It covers correcting errors, adding new information, etc. We conduct evaluations on 3 subsets: *Arithmetic Reasoning, Scientific Reasoning*, and *New Information*. The metric is *Accuracy*.

LeKUBE(Wang et al., 2024a) A Chinese legal knowledge editing dataset . We focus on two challenging tasks, where prior baselines show poor performance: *Multiple-Choice Questions of the Legal Scenario* and *True/False Questions of Change in Statute.* The metric is *Accuracy.* The statistics of the three benchmarks are detailed in Appendix C.

4.2 Models

We evaluate two open-source models: Meta-Llama-3-8B-Instruct(Dubey et al., 2024) and Qwen2.5-14B-Instruct(Team, 2024). Since the former does not support Chinese, we evaluate the latter on LeKUBE. Following LeKUBE's protocol, we first fine-tune the model on the STARD corpus(Su et al., 2024a), which contains Chinese laws and judicial interpretations, to ensure the model acquires Chinese legal knowledge before any updates.

We also assess GPT-4o-2024-08-06(OpenAI et al., 2023) on MQuAKE-3k-CF. Due to its proprietary nature, we test only a few non-parametric editing methods on this model. Following RAE(Shi et al., 2024), we randomly sample 300 data points to minimize costs. Meta-Llama-3-8B-Instruct is used as a proxy model for tasks the proprietary model cannot perform. In RAE, this proxy extracts and prunes knowledge graphs based on probabilities, while in EditCoT, it acts as a CoT editor.

4.3 Baselines

We compare EditCoT with a range of model editing techniques, including knowledge neuronbased (KN)(Dai et al., 2022), rank-one model editing (ROME)(Meng et al., 2022), and incontext methods such as Mello(Zhong et al., 2023), PokeMQA(Gu et al., 2024a), and RAE(Shi et al., 2024). We also include full-parameter fine-tuning (FT) and retrieval-augmented generation (Naive RAG)(Lewis et al., 2020) as baselines. To distinguish our approach from the RAG method, we evaluate two advanced variants, RAT(Wang et al., 2024c) and FLARE(Jiang et al., 2023), on the

MQuAKE and DUNE benchmarks. RAT incorporates multi-turn query rewriting, retrieval, and CoT 404 enhancements to address complex tasks effectively. 405 FLARE dynamically determines the necessity of 406 retrieval during model inference. These comparisons highlight the superiority of our approach in 408 performing knowledge editing tasks. 409

403

407

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

451

Notably, since LeKUBE is a Chinese dataset, we translate the official prompts for Mello and PokeMQA into Chinese to accommodate their sensitivity to language type. RAT and FLARE are not evaluated on LeKUBE because it is challenging to align their settings with those of the dataset. RAE is not applicable to datasets other than MQuAKE-CF-3k since it requires construction of knowledge graphs, which is impractical on these datasets. Details on the baselines are available in Appendix D.

4.4 Implementation Details

For the CoT editor, the data construction process is outlined in subsection 3.4. We use 80% of the filtered data generated by the respective model and train it for 1 epoch. To generate the initial answer, CoT, and final answer, and detect conflicts, we use 5-shot prompts styled after Mello and PokeMQA and apply this format across all tasks, as detailed in subsection H.2. For conflict detection, we instruct the model to classify outputs into three categories: "Contradict", "Support", or "Unrelated". We take the token with the highest probability as the final classification. While we use the [STEP] token during dataset construction, in actual inference, we allow the model to freely generate the CoT, with steps separated by newline characters.

For the EditCoT retrieval corpus, settings are provided in section 3.2. For datasets offering both preand post-edit knowledge (e.g., MQuAKE-CF-3k and LeKUBE), we retrieve from the pre-edit corpus and match with corresponding post-edit sentences. For datasets with only post-edit knowledge (e.g., DUNE), we directly retrieve from the post-edit corpus. We use the same retriever, Contriever(Izacard et al., 2022), as Mello in MQuAKE-CF and DUNE. For Chinese datasets, we use BM25(Robertson et al., 2009), which performs well in LeKUBE.

Experimental Results 5

5.1 Effectiveness

Generalizability. EditCoT demonstrates superior dataset generalization and model robustness compared to baselines. As shown in Table 1, EditCoT

achieves top performance on 6 out of 8 results and second-best on 2 out of 8 results in MQuAKE-3k-CF and DUNE, outperforming all baselines. It consistently excels across two LLMs, unlike other methods like Mello, PokeMQA, and RAE, which perform well on Llama-3-8B-Instruct but poorly on Owen2.5-14B-Instruct. Mello and PokeMOA's lack of robustness stems from their dependence on strict instructions and hand-crafted few-shot examples, which exhibit varying effectiveness across models and tasks. Additionally, safety-aligned models often reject external edits. To provide a detailed comparison between our method and the baselines, we include case studies in Appendix G.

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

Applicability to Vertical Domains and Blackbox Models. EditCoT is also effective in specialized vertical domains and proprietary models. Table 2 presents the evaluation results on LeKUBE, where EditCoT outperforms all baselines across two tasks, significantly surpassing other methods in the multiple-choice questions of the legal scenario. Furthermore, Table 3 reports GPT-4o's performance on MQuAKE-CF, where EditCoT also demonstrates competitive results, achieving secondbest overall. Notably, the proxy model in this experiment is Meta-Llama-3-8B-Instruct. When we compare the performance of the proxy model itself with the GPT-40, we find that, while RAE's accuracy improves from 54.1% to 59.7% with GPT-40, EditCoT achieves a larger relative improvement, rising from 35.4% to 45.0%. This suggests that EditCoT's success may depend more on the model's internal reasoning abilities, whereas RAE benefits primarily from leveraging external knowledge.

Challenges in General-Purpose In-Context Editing. Developing a broadly applicable incontext editing framework remains a significant challenge across both datasets and models. The three datasets we evaluate spanning multi-hop reasoning, generalized edits, domain-specific knowledge, and different languages, posing difficulty for current editing methods. Among the baselines, PokeMQA performs well on multi-hop questions with Llama-3-8B-Instruct, and simpler methods like Naive RAG excel in narrow settings. However, these methods perform poorly on other datasets and models. RAE excels on the MQuAKE-CF-3k, because it constructs structured knowledge graphs and accesses a broader range of external knowledge (e.g., Wikipedia), but fails to generalize beyond structured datasets. In contrast, EditCoT addresses these limitations, achieving the best or second-best

Table 1: Results on MQuAKE-CF-3k and three subsets of DUNE, with MQuAKE-CF-3k using the Multi-hop Accuracy metric and the others using Accuracy. The best performance for each dataset and model is highlighted in bold, and the second-best is underlined. A dash ("-") indicates that the method is not applicable to the given dataset. Llama and Qwen represent Llama-3-8B-Instruct and Qwen2.5-14B-Instruct, respectively.

Models	D	Datasats		Parametric Methods		RAG Methods		In-Context Editing				
wioueis	Datasets		FT	KN	ROME	Naive RAG	RAT	FLARE	Mello	PokeMQA	RAE	EditCoT(Ours)
	MQuA	KE-CF-3k	11.2	2.8	4.0	9.7	0.7	3.3	10.0	26.1	54.1	35.4
Llama		Arithmetic	74.1	80.4	80.5	84.0	70.0	40.1	73.7	83.8	-	90.7
	DUNE	Scientific	11.1	82.1	81.4	81.6	84.3	82.1	77.1	61.7	-	85.0
		New Info	30.3	66.6	69.1	89.2	77.0	81.6	89.8	58.7	-	91.3
	MQuA	KE-CF-3k	9.3	3.2	0.0	10.0	3.3	3.7	5.8	5.3	26.9	34.2
Qwen		Arithmetic	82.8	86.9	86.9	89.7	82.6	23.4	61.0	35.2	-	97.3
	DUNE	Scientific	13.8	75.3	75.5	70.7	86.0	83.2	37.0	7.7	-	86.3
		New Info	76.2	73.4	73.1	96.8	87.4	70.1	34.3	17.3	-	<u>93.8</u>

Table 2: Applicability to vertical domains: LeKUBE is a Chinese Legal Knowledge Editing Benchmark. The model used is Qwen2.5-14B-Instruct.

Subsets of LeKUBE	Fine-tuning	Naive RAG	KN	ROME	Mello	PokeMQA	EditCoT(Ours)
MCQ of the Legal Scenario	42.2	47.2	38.9	38.3	6.7	0.0	58.3
T/F Questions of Change in Statute	65.9	<u>69.2</u>	55.8	56.1	18.6	0.0	69.5

Table 3: Results on GPT-40 (a black-box LLM). Llama-3-8B-Instruct serves as a proxy LLM for RAE and Edit-CoT. The dataset is 300 samples from MQuAKE-CF.

	Naive RAG	Mello	PokeMQA	RAE	EditCoT(Ours)
Score	12.0	15.0	9.7	59.7	<u>45.0</u>

performance across most datasets and models, although it is not always optimal.

RAG methods underperform compared to EditCoT. There are two reasons: 1) RAG methods seldom address conflicts between parametric and contextual knowledge, which are common in knowledge editing. While RAT modifies CoT, Table 5.4 shows that EditCoT, using a trained editor, injects knowledge more effectively than a simple prompt. 2) Certain assumptions of RAG are incompatible with knowledge editing. FLARE dynamically decides on retrieval based on model uncertainty about generated words. However, in knowledge editing, the LLM is often confident about the old knowledge contained within its parameters.

5.2 Robustness

505

508

509

510

511

513

514

515

516

517

518

519

520

522

523

526

In this section, we evaluate EditCoT against other in-context knowledge editing methods, using Llama3-8B-Instruct from multiple perspectives.

Locality: Locality is a crucial metric in knowledge editing, assessing a method's ability to update specific knowledge without impairing untargeted information (Mitchell et al., 2022). We evaluate loTable 4: Locality test on the New Info (Locality) set of DUNE. EditCoT is the best in all methods.

B	efore Editing	Naive RAG	Mello	PokeMQA	EditCoT
Accuracy	65.2	34.3	58.0	45.6	59.7

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

cality using the New Information (Locality) subset of DUNE, where the retrieval scope is limited to the edited set (consistent with the main experimental setup), but test tasks pertain to knowledge outside this set. The results, presented in Table 4, are compared with the unedited base model's performance. While all editing methods show some decline in performance relative to the unedited model, Edit-CoT exhibits the best locality. RAG performs the worst, as it lacks conflict detection between documents and queries, resulting in poor handling of irrelevant information.

Performance with Different Batch Sizes: The performance of knowledge editing methods varies with the editing batch size. Here we denote the batch size as the number of questions in one batch. We test the sensitivity of Mello, PokeMQA, and EditCoT to batch sizes.² We evaluate batch sizes of 1, 10, 100, and 300 on a 300-sample dataset from MQuAKE-CF-3k. Results in Figure 5 indicate that all methods perform best with a single sample, with performance decreasing as size increases. However,

²We don't compare RAE because the locality of in-context editing is influenced by retrieval. RAE uses knowledge graph retrieval, which is not comparable to the other three methods.



Figure 4: Inference time vs. performance: A method positioned in the top-left corner indicates shorter inference time and higher performance. Due to the extensive horizontal span, the figure is divided in the middle.



Figure 5: Performance with different batch sizes on MQuAKE-CF. The metrics here is multi-hop accuracy.

EditCoT consistently outperforms the others across all batch sizes, showing the smallest performance degradation (around 20%) from size 1 to 300, compared to declines of over 30% for the others.

5.3 Efficiency

552

553

554

555

557

558

559

560

566

570

We investigate efficiency on 100 sampled questions from MQuAKE-CF-3k. By varying the maxiterations for Mello, PokeMQA, and EditCoT, we record their corresponding inference time and performance. Complete experimental details are in Appendix E. As shown in Figure 4, EditCoT achieves superior performance with the shortest inference time, outperforming baselines at similar computational costs. The efficiency advantage stems from requiring fewer iterations (3.65 average iterations for EditCoT vs. 4.33 average iterations for PokeMQA) and demonstrating stronger early performance within initial iterations. Further details on the impact of iteration limits are in Figure F. The average inference time for RAE and two advanced RAG methods exceeds 20 seconds, making their application in real-world tasks challenging.

Table 5: Impact of training settings on the performance of the CoT editor, evaluated on MQuAKE-CF.

# Epoch	1	2	3	4	5
Accuracy	40.3	40.7	40.3	40.0	40.0
# Training Data Accuracy	50 39.0	100 39.3	200 40.7	250 40.3	316(All) 40.3

5.4 Ablation Study

We conduct ablation studies using a dataset of 300 examples sampled from MQuAKE-CF-3k, with the model Qwen2.5-14B-Instruct.

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

Impact of Components We evaluate two components of our method. First, removing conflict detection before editing the CoT results in a 18.1% performance drop (33.0 vs. 40.3). Second, using the original model as the CoT editor without specialized training leads to an even more significant performance degradation (26.7 vs. 40.3). These results demonstrate the critical role of both components: conflict detection prevents contamination from irrelevant retrieved information, while training the editor enables effective utilization of model editing capabilities beyond simple prompting.

Impact of Training Configurations Table 5 illustrates the effect of training dataset sizes and epochs on the CoT editor. While increasing the number of epochs has little impact, performance stabilizes after expanding the training size to 200 examples, suggesting that the CoT editor requires a sufficient amount of diverse data to learn effectively. This supports our main experimental setup, where training for one epoch with 80% of generated data (252 examples) is a reasonable choice.

6 Conclusion

In this paper, we introduce **EditCoT**, a novel framework designed for efficiently updating LLMs through iterative CoT editing. EditCoT enhances the generalizability and robustness of knowledge editing across diverse tasks, without the need for retraining the LLM each time new knowledge is encountered. Our experiments, conducted on a variety of benchmarks, demonstrate that EditCoT consistently outperforms existing methods, proving its effectiveness in both general-purpose and domainspecific applications. This approach presents a flexible and efficient solution for continuous knowledge updates in LLMs, marking a significant advancement in the field of knowledge editing.

7 Limitations

612

621

623

631

632

633

634

635

641

645

651

657

659

663

EditCoT has certain limitations, particularly in its reliance on two distinct models. While this design choice does not increase inference time, it does result in higher GPU resource consumption, which can be challenging in environments with limited hardware resources. Additionally, the iterative nature of the CoT refinement process may add to the computational overhead.

> For the CoT editor, the extent to which the model size can be reduced and the potential for further compression have not been fully explored in this work. Future work could explore the possibility of smaller editors, aiming to reduce the computational burden without compromising performance. This could make EditCoT more accessible and practical for a wider range of applications and deployment scenarios.

References

- Afra Akyürek, Eric Pan, Garry Kuwanto, and Derry Wijaya. 2023. DUnE: Dataset for unified editing. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 1847–1861, Singapore. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Junjie Chen, Weihang Su, Zhumin Chu, Haitao Li, Qinyao Ai, Yiqun Liu, Min Zhang, and Shaoping Ma. 2024. An automatic and cost-efficient peer-review framework for language generation evaluation. *arXiv preprint arXiv:2410.12265*.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493– 8502, Dublin, Ireland. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-

664

665

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

ney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks,

728

729

736

737

739

740

741

743

744

746

747

749

751

753

755

770

771

772

773

774 775

776

777

778

779

780

781

782

784

785

790

791

Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models.

792

793

794

795

796

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

- Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ninghao Liu, Ruobing Wang, and Xin Wang. 2024a. PokeMQA: Programmable knowledge editing for multi-hop question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8069–8083, Bangkok, Thailand. Association for Computational Linguistics.
- Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024b. Model editing harms general abilities of large language models: Regularization to the rescue. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP* 2024, Miami, FL, USA, November 12-16, 2024, pages 16801–16819. Association for Computational Linguistics.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. *Trans. Mach. Learn. Res.*, 2022.

Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun,

Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 7969–7992. Association for Computational Linguistics.

853

867

871

873

876

878

879

883

884

886

891

894

896

900 901

902

903 904

905

906

907

908

909

910

911

- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022. Memorybased model editing at scale. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 15817–15831. PMLR.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik

Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-912 Lopes, Jonathan Gordon, Morgan Grafstein, Scott 913 Gray, Ryan Greene, Joshua Gross, Shixiang Shane 914 Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 915 Yuchen He, Mike Heaton, Johannes Heidecke, Chris 916 Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 917 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin 918 Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, 919 Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 920 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-921 woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-922 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, 923 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, 924 Christina Kim, Yongjik Kim, Jan Hendrik Kirch-925 ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, 926 Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-927 stantinidis, Kyle Kosic, Gretchen Krueger, Vishal 928 Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan 929 Leike, Jade Leung, Daniel Levy, Chak Ming Li, 930 Rachel Lim, Molly Lin, Stephanie Lin, Mateusz 931 Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, 932 Anna Makanju, Kim Malfacini, Sam Manning, Todor 933 Markov, Yaniv Markovski, Bianca Martin, Katie 934 Mayer, Andrew Mayne, Bob McGrew, Scott Mayer 935 McKinney, Christine McLeavey, Paul McMillan, 936 Jake McNeil, David Medina, Aalok Mehta, Jacob 937 Menick, Luke Metz, Andrey Mishchenko, Pamela 938 Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 939 Mossing, Tong Mu, Mira Murati, Oleg Murk, David 940 Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, 941 Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, 942 Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex 943 Paino, Joe Palermo, Ashley Pantuliano, Giambat-944 tista Parascandolo, Joel Parish, Emy Parparita, Alex 945 Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-946 man, Filipe de Avila Belbute Peres, Michael Petrov, 947 Henrique Ponde de Oliveira Pinto, Michael, Poko-948 rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-949 950 ell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, 951 Cameron Raymond, Francis Real, Kendra Rimbach, 952 Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-953 der, Mario Saltarelli, Ted Sanders, Shibani Santurkar, 954 Girish Sastry, Heather Schmidt, David Schnurr, John 955 Schulman, Daniel Selsam, Kyla Sheppard, Toki 956 Sherbakov, Jessica Shieh, Sarah Shoker, Pranav 957 Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, 958 Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin 959 Sokolowsky, Yang Song, Natalie Staudacher, Fe-960 lipe Petroski Such, Natalie Summers, Ilya Sutskever, 961 Jie Tang, Nikolas Tezak, Madeleine B. Thompson, 962 Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, 963 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-964 lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, 965 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, 966 Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, 967 CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-968 ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, 969 Clemens Winter, Samuel Wolrich, Hannah Wong, 970 Lauren Workman, Sherwin Wu, Jeff Wu, Michael 971 Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-972 ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong 973 Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 974 Zheng, Juntang Zhuang, William Zhuk, and Barret 975

976 977

978

979

- 981 982 984
- 989

987

991

992 993

- 995 996 997
- 998 999
- 1000
- 1001 1002

1004

- 1005
- 1007 1008

1009 1010

- 1012 1013
- 1015 1016

1017

1018 1019 1020

1021 1022

- 1027 1028
- 1029 1030

- Zoph. 2023. Gpt-4 technical report.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. Foundations and Trends® in Information Retrieval, 3(4):333-389.
 - Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 2024. Retrieval-enhanced knowledge editing in language models for multi-hop question answering.
- Weihang Su, Yiran Hu, Anzhe Xie, Qingyao Ai, Zibing Que, Ning Zheng, Yun Liu, Weixing Shen, and Yiqun Liu. 2024a. Stard: A chinese statute retrieval dataset with real queries issued by non-professionals.
- Weihang Su, Yichen Tang, Qingyao Ai, Changyue Wang, Zhijing Wu, and Yiqun Liu. 2024b. Mitigating entity-level hallucination in large language models. arXiv preprint arXiv:2407.09417.
- Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024c. Unsupervised real-time hallucination detection based on the internal states of large language models. In Findings of the Association for Computational Linguistics ACL 2024, pages 14379–14391, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Weihang Su, Changyue Wang, Anzhe Xie, Qingyao Ai, Yiran Hu, and Yiqun Liu. 2024d. Legalaid: A large language model for the chinese legal field. https: //github.com/oneal2000/LegalAID.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Changyue Wang, Weihang Su, Yiran Hu, Qingyao Ai, Yueyue Wu, Cheng Luo, Yiqun Liu, Min Zhang, and Shaoping Ma. 2024a. Lekube: A knowledge update benchmark for legal domain. In Proceedings of the 2024 Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region, pages 175–185.
- Changyue Wang, Weihang Su, Hu Yiran, Qingyao Ai, Yueyue Wu, Cheng Luo, Yiqun Liu, Min Zhang, and Shaoping Ma. 2024b. Lekube: A legal knowledge update benchmark.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, et al. 2023a. Easyedit: An easy-to-use knowledge editing framework for large language models. ArXiv preprint, abs/2308.07269.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. Self-consistency improves chain of thought reasoning in language models. In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net.

Zihao Wang, Anji Liu, Haowei Lin, Jiaqi Li, Xiaojian Ma, and Yitao Liang. 2024c. Rat: Retrieval augmented thoughts elicit context-aware reasoning in long-horizon generation. Preprint, arXiv:2403.05313.

1031

1032

1034

1035

1036

1038

1039

1041

1042

1043

1045

1046

1047

1048

1049

1050

1051

1052

1053

1055

1056

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1081

1082

1083

1084

1085

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38-45, Online. Association for Computational Linguistics.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Kegin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. ArXiv preprint, abs/2407.10671.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 2369-2380, Brussels, Belgium. Association for Computational Linguistics.
- Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023. Verify-and-edit: A knowledge-enhanced chain-of-thought framework. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5823-5840, Toronto, Canada. Association for Computational Linguistics.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong 1087 Wu, Jingjing Xu, and Baobao Chang. 2023. Can 1088 we edit factual knowledge by in-context learning? 1089

In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 4862-4876, Singapore. Association for Computational Linguistics.

1090 1091

1092 1093

1098

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

Zexuan Zhong, Zhengxuan Wu, Christopher Manning, Christopher Potts, and Danqi Chen. 2023. MQuAKE: Assessing knowledge editing in language models via multi-hop questions. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 15686–15702, Singapore. Association for Computational Linguistics.

A Pseudocode Description of EditCoT

The pseudocode description of EditCoT is shown in Algorithm 1.

Detailed Comparison with the B Verify-and-Edit Framework

Zhao et al. (2023) introduce a Verify-and-Edit framework that, on the surface, appears similar to our approach as it also involves editing the generated CoT. However, there are several key distinctions: (1) their verification process relies on the self-consistency, which works well for general QA tasks but is less applicable to knowledge editing since the model often shows high confidence in outdated information, leading to high consistency even when incorrect. We conducted tests on MQuAKE-CF-3k, focusing only on the triggering of editing and ignoring subsequent edit results. Using the self-consistency, edits were triggered in only 37% of cases on the knowledge editing dataset, while our method triggered edits in 95% of cases.; (2) their editing relies on the model's ability to learn from context, whereas in Table 5.4, we demonstrate that our proposed CoT editor significantly outperforms this self-editing approach in knowledge editing tasks; and (3) their method actually resembles Mello's (Zhong et al., 2023) strategy of decomposing tasks, whereas our approach modifies the entire CoT in each round of editing, allowing for more comprehensive adjustments.

С **Details of Benchmarks**

Table 6 shows the statistics of the three benchmarks, including the total number of edits and queries.

Details of Baselines D

We focus on In-Context Editing methods and also compare EditCoT with several traditional model editing techniques: 1136

Dataset	Subset	Edits	Queries
MQuAKE-CF-3k	-	2785	3,000
	Scientific Reasoning	223	1508
DUNE	Arithmetic Reasoning	184	1065
	New Information	200	1000
LeKUBE	MCQ of the Legal Scenario	180	180
	T/F Questions of Change in Statute	180	642

Table 6: Statistics on the number of edits and questions for the benchmarks

Fine-tuning (FT): A full-parameter fine-tuning Since we evaluate Instruct Models, approach. we train the model on question-answer pairs constructed from each dataset's editing instances.

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

Naive RAG (Lewis et al., 2020): Given an external knowledge base and a retriever, the Retrieval-Augmented Generation (RAG) framework injects relevant knowledge into the model by concatenating retrieved documents into the context, effectively augmenting the language model's ability to reason and reduce hallucinations (Su et al., 2024c; Chen et al., 2024; Su et al., 2024b,d). Following the settings from the original LeKUBE and DUNE papers, we use BM25(Robertson et al., 2009) as the retriever, a retrieval method based on lexical matching, and concatenate 3 and 1 retrieved documents, respectively. For the evaluation of MQuAKE-CF-3k, since the original paper does not test RAG, we still use BM25 as the retriever, setting it to concatenate 5 retrieved documents.

RAT(Wang et al., 2024c): A RAG method that refines the model-generated reasoning steps through multiple rounds of query rewriting and retrieval. We utilize the official prompt template and ultimately guided the model to derive the final answer based on the original question and the final reasoning steps. All other settings are identical to those in Naive RAG.

FLARE(Jiang et al., 2023): A dynamic RAG that determines whether to perform retrieval by assessing the uncertainty of the generated sentences during the generation process. We set the hallucination detection threshold at 0.12. All other settings remain the same as in Naive RAG.

KN(Dai et al., 2022): This method updates knowledge by identifying and editing "knowledge neurons" in the model. In our experiments, we set the prompt number n for neuron identification to 10, the knowledge attribution threshold to 0.2, and the probability of retaining shared neurons to 0.4.

ROME(Meng et al., 2022): A Rank-One Model 1177 Editing method that directly modifies key-value 1178 pairs in the FFN layers. For our experiments, we 1179

Algorithm 1 EditCoT Framework

1:	Input: Query q , Original Model M , Edited Model R	M_{edit} , Memory of Facts K, Maximum Iterations
	N, Mapping of old facts to new facts F	
2:	Output: Final Answer A_{final}	
3:	$A_{\text{old}} \leftarrow \text{GenerateAnswer}(M, q)$	▷ Generate initial answer
4:	$CoT_{old} \leftarrow GenerateCoT(M, q, A_{old})$	▷ Generate initial chain-of-thought
5:	for $i = 1$ to N do	C C
6:	for each step s in CoT_{old} do	
7:	$F_{\text{new}} \leftarrow \text{RetrieveFacts}(s, K)$	▷ Retrieve new fact from memory
8:	if F exists then	
9:	$F_{\text{new}} \leftarrow F[F_{\text{new}}]$	▷ Map old fact to new fact if applicable
10:	end if	
11:	if $F_{\text{new}} \neq \emptyset$ then	
12:	break	
13:	end if	
14:	end for	
15:	$is_Conflict \leftarrow Verify(M, q, CoT_{old}, F_{new})$	▷ Verify conflict between CoT and new fact
16:	if not <i>is_Conflict</i> then	
17:	break	
18:	end if	
19:	$CoT_{new} \leftarrow EditChain(M_{edit}, q, CoT_{old}, F_{new})$	▷ Edit chain-of-thought with new fact
20:	$CoT_{old} \leftarrow CoT_{new}$	▷ Update old CoT with new one
21:	end for	
22:	$A_{\text{final}} \leftarrow \text{GenerateAnswer}(M, CoT_{\text{new}})$	▷ Generate final answer based on edited CoT
23:	return A _{final}	

uniformly modify the fifth layer and calculate the loss at the final layer of the model. The weight decay is set to 1×10^{-3} .

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1204

Both KN and ROME methods are implemented using the EasyEdit(Wang et al., 2023a) library, which integrates several model editing techniques for ease of use. It is important to note that both methods rely on knowledge triples for editing. In the *Scientific Reasoning* and *New Information* subsets of DUNE, the edits consist of longer natural language sequences. We use GPT-40 to extract the relevant knowledge triples.

Mello(Zhong et al., 2023): Mello performs incontext editing by decomposing the problem into sub-questions and utilizing retrieval. We follow the official setting with a maximum of 4 retrieval rounds and adapt the prompts to Instruct Models. In English datasets, we use Contriever(Izacard et al., 2022) as the retriever, following the settings from the original paper. For the Chinese LeKUBE dataset, we translate the prompts into Chinese and use BM25 as the retriever.

PokeMQA(Gu et al., 2024a): Built on top of Mello, PokeMQA adds entity extraction and scope determination to refine the question understanding. We follow the official setting with a maximum of 5 rounds and use their pre-trained Scope-Detector. For DUNE and LeKUBE, where no pre-extracted entities are available, we instruct the model to extract entities in a similar format. Since the official Scope-Detector only supports English, we translate the input sentences for LeKUBE tasks into English before applying the detector.

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

During the GPT-40 evaluations, we observed difficulties in strictly following the prompt format. To mitigate this issue, we add the system prompt: "Follow the examples below."

Additionally, for Mello, PokeMQA and Edit-CoT, since LeKUBE evaluates changes in legal provisions, and other baselines can utilize both old and new laws, we concatenate the pre- and postupdate legal texts and provide them to the model after retrieval for fair comparison.

RAE(Shi et al., 2024): RAE constructs knowledge graphs for retrieval and leverages the model to assist in retrieving and pruning the graphs. However, for DUNE and LeKUBE, knowledge graphs similar to those in Wikipedia cannot be constructed, making RAE less generalizable to these tasks.

Additionally, the LLMs we used are imple-



Figure 6: Performance with different max iterations. A full retrieval-generation cycle is defined as one iteration.

mented by Huggingface Transformers library (Wolf et al., 2020).

E Details of Efficiency Experiments

1230

1231

1233

1234

1235

1236

1237

1238

1239

1240

1241

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255 1256

1257

1258

1259

1260

1261

1262

All efficiency experiments are conducted on a single NVIDIA A100 (40G) GPU.

RAE's official implementation uses Wikipedia API calls; hence, we exclude network latency from its reported inference time (including this would place RAE around 46 seconds on the x-axis). Despite this, RAE has a notably longer inference time due to multiple forward propagations for knowledge graph retrieval and pruning, particularly when selecting the next relation, which requires feeding all candidates into the model respectively.

The efficiency of EditCoT is attributed to fewer required iterations: while PokeMQA averages 4.33 iterations, EditCoT needs only 3.65, indicating fewer long-text generation steps. Here a full retrieval-generation cycle is defined as one iteration.

F Performance with different max iterations

We explore how varying the maximum number of iterations affects their performance. Figure 6 shows the results. PokeMQA and Mello require more iterations to reach optimal performance, with relatively low accuracy in early iterations. On the other hand, EditCoT outputs a complete chain-of-thought in each iteration, achieving decent performance after the first iteration.

G Case Studies

In Figure 7, we compare EditCoT with Mello and PokeMQA on two questions. The first question

from the Scientific Reasoning subset of DUNE 1263 shows that EditCoT initially generates an incor-1264 rect answer and CoT. However, after retrieval and 1265 conflict detection, it modifies the CoT to arrive at 1266 the correct solution. In contrast, Mello retrieves 1267 a relevant edit but fails to detect any conflict with 1268 the initial answer, missing the opportunity for cor-1269 rection. This demonstrates the strength of Edit-1270 CoT's conflict detection mechanism, which eval-1271 uates the entire reasoning chain against retrieved 1272 edits. PokeMQA struggles with both sub-questions, 1273 producing a repetitive decomposition and failing to 1274 trigger an edit. 1275

The second question from MQuAKE-CF-3k 1276 again illustrates EditCoT's ability to transition from 1277 an incorrect inference to the correct answer through 1278 CoT editing. Although both PokeMQA and Mello 1279 successfully retrieve and detect the relevant edit 1280 (changing the author of Harry Potter to Maurice 1281 Maeterlinck), they do not apply it. This indicates 1282 that models designed with safety in mind may re-1283 sist incorporating external knowledge in favor of 1284 their internal knowledge. EditCoT addresses this 1285 issue by decoupling conflict detection from the editing process and training a model specifically for 1287 CoT editing, leading to more effective reasoning 1288 updates. 1289

H Prompt Templates

1290

Since the models we are working with are all In-1291struct Models, the following prompt template uses1292"system", "user", and "assistant" to represent dif-1293ferent roles.1294

H.1 Prompt Templates of Dataset 1295 Construction 1296

First, we use RAG to generate CoT and the answer.	1297
The prompt template is as follows:	1298



Figure 7: A face-to-face comparison of our proposed EditCoT method with other baseline approaches. The first example comes from the Scientific Reasoning subset of DUNE. The second example comes from MQuAKE-CF-3k. The model here is Meta-Llama-3-8B-Instruct.

Dataset Construction - RAG Prompt	
User: { <i>Top-5 relevant paragraphs from</i> <i>Wikipedia</i> }	tioi
Instruction: You have access to background information that may assist in answering the question. Focus on reasoning through the problem step by step, keeping the ex- planation concise. Use the token [STEP] to start each thought step and the token [AN- SWER] to indicate the final answer. Avoid mentioning or referencing the background information directly in your reasoning.	
Template: [STEP] [STEP] [ANSWER]	
Question: { <i>A question from HotpotQA</i> }	

Second, we conduct prefix-guided CoT generation. The prompt template is as follows:

1300

User: Evaluate the relationship between the User: Please answer the following question following two sentences based on their facusing a chain-of-thought. Use the token tual content. Choose the most appropriate [STEP] to start each thought step, and the option from the following: token [ANSWER] to indicate the final answer. Keep each step brief and to the point A. The two sentences contain conflicting and keep the final answer concise. knowledge. B. The two sentences support or comple-Template: ment each other. [STEP] ... C. The two sentences are unrelated (no con-[STEP] ... flict, but no connection). [ANSWER] ... Sentence 1: {*CoT Step k*} Sentence 2: {*CoT Step k**} Question: {A question from HotpotQA} **Assistant:** [STEP] {*The first (k-1) steps of* Your choice: *the CoT from RAG*}

Finally, we let the model rewrite the key step:

Dataset Construction - Prompt of CoT Step

Conflict Detection

1308

1307

Then, we detect conflict in the final answers:

Dataset Construction - Prompt of prefix-

guided CoT generation

Dataset Construction - Prompt of Answer **Conflict Detection**

User: For the given question Q, and the two provided answers (Answer 1 and Answer 2), determine if the answers are semantically consistent. Choose the most appropriate option from the following:

A. Yes, they convey similar ideas or information.

B. No, they are different or contradictory.

Question Q: {A question from HotpotQA} Answer 1: {*The answer generated by RAG*} Answer 2: {The answer generated by prefixguided CoT generation}

Your choice:

And detect if there is a conflict at the first diverging step between the two CoTs:

Dataset Construction - Rewriting Prompt

User: Rewrite the following sentence in the style of a Wikipedia data, using formal and objective language. Only express the factual knowledge present in the sentence, without adding any extra information or inventing details. Your output must be one or more sentences, only needs to include direct results without extra words.

Input:

Sentence: {A CoT step that needs to *rewrite* }

Output:

1309

1310

H.2 Prompt Templates of EditCoT

The prompt template for CoT editor is as follows. 1311 Note that the template is the same during both the 1312 training and inference phases. 1313

1306

1	3	0	2

1303

{...}

EditCoT Framework - Prompt for CoT editor

System: You can edit the following chains of thought based on the new knowledge. User: Question: {*A user question*} Old Thoughts: {*A old CoT*} New knowledge: {*A sentence containing new knowledge*} Please give me the new chain-of-thought based on the new knowledge. Assistant: New Thoughts:

1314

1315 1316

1317

The prompt templates below are for the original LLM. First, we instruct it to generate an answer directly:

EditCoT	Framework - Prompt of Answering
directly	

User: Question: What is the capital of the country where Plainfield Town Hall is located?

Assistant: Answer: Washington, D.C.

User: Question: In which country is the company that created Nissan 200SX located?

Assistant: Answer: Japan

User: Question: Which continent is the country where the director of "My House Husband: Ikaw Na!" was educated located in?

Assistant: Answer: Asia

User: Question: Who is the spouse of the US president? Assistant: Answer: Jill Biden

User: Question: Who has ownership of the developer of the Chevrolet Corvette (C4)? **Assistant:** Answer: General Motors

The prompt template of get an original CoT from

User: Question: {*A user question*} **Assistant:** Answer:

a question and an answer is as follows:

1318

1319 1320 EditCoT Framework - Prompt of Getting the original CoT

User: Question: What is the capital of the country where Plainfield Town Hall is located?

Answer: Washington, D.C.

Please give the chain of thought based on the question and answer pairs above.

Assistant: Thoughts: Plainfield Town Hall is located in the country of the United States of America.

The capital of United States is Washington, D.C.

User: Question: In which country is the company that created Nissan 200SX located?

Answer: Japan

Please give the chain of thought based on the question and answer pairs above.

Assistant: Thoughts: Nissan 200SX was created by Nissan.

Nissan is located in the country of Japan.

User: Question: Which continent is the country where the director of "My House Husband: Ikaw Na!" was educated located in?

Answer: Asia

Please give the chain of thought based on the question and answer pairs above.

Assistant: Thoughts: The director of "My House Husband: Ikaw Na!" is Jose Javier Reyes.

Jose Javier Reyes was educated at De La Salle University.

De La Salle University is located in the country of Philippines.

Philippines is located in the continent of Asia.

User: Question: Who is the spouse of the US president? Answer: Jill Biden

Please give the chain of thought based on the question and answer pairs above.

Assistant: Thoughts: The US president is Joe Biden.

1321

The spouse of Joe Biden is Jill Biden.

User: Question: Who has ownership of the developer of the Chevrolet Corvette (C4)? Answer: General Motors Please give the chain of thought based on the question and answer pairs above. **Assistant:** Thoughts: The developer of Chevrolet Corvette (C4) is Chevrolet. Chevrolet is owned by General Motors.

User: Question: {*A user question*} Answer: {*An original answer*} Please give the chain of thought based on the question and answer pairs above. **Assistant:** Thoughts:

The prompt template of conflict detection is as follows. We use 5-shot examples and apply a scenario similar to CoT editor.

EditCoT Framework - Prompt of Conflict Detection

System: You can edit the following chains of thought based on the new knowledge. User: Question: Which country is the author of Moneyballä citizen of? Old Thoughts: The author of Moneyball is Michael Lewis. Michael Lewis is a citizen of United States of America. New knowledge: The author of Our Mutual Friend is Charles Dickens Please give me the new chain-of-thought based on the new knowledge. Assistant: Relevance of knowledge and chain-of-thought: Unrelated. **User:** Question: What language does the Shadow Home Secretary speak, write, or sign? Old Thoughts: The Shadow Home Secretary is Diane Abbott. Diane Abbott speaks the language of English. New knowledge: The Shadow Home Secretary is Diane Abbott Please give me the new chain-of-thought based on the new knowledge.

Assistant: Relevance of knowledge and chain-of-thought: Support.

User: Question: What is the notable work of the author of This Side of Paradise? Old Thoughts: The author of This Side of Paradise is F. Scott Fitzgerald.

F. Scott Fitzgerald is famous for The Great Gatsby.

New knowledge: The author of This Side of Paradise is Geoffrey Chaucer

Please give me the new chain-of-thought based on the new knowledge.

Assistant: Relevance of knowledge and chain-of-thought: Contradict.

New Thoughts: The author of This Side of Paradise is Geoffrey Chaucer.

Geoffrey Chaucer is famous for Shaman King.

User: Question: Who is the head of government of the country where Roy Ascott holds citizenship?

Old Thoughts: Roy Ascott is a citizen of United Kingdom.

The name of the current head of the United Kingdom government is Boris Johnson.

New knowledge: Roy Ascott is a citizen of United States of America

Please give me the new chain-of-thought based on the new knowledge.

Assistant: Relevance of knowledge and chain-of-thought: Contradict.

New Thoughts: Roy Ascott is a citizen of United States of America.

The name of the current head of the United States of America government is Ole Johan Vierdal.

User: Question: Who is the spouse of the performer of the song Don't Be Cruel? Old Thoughts: Don't Be Cruel was performed by Elvis Presley.

Elvis Presley is married to Priscilla Presley. New knowledge: Don't Be Cruel was performed by Elvis Presley

Please give me the new chain-of-thought based on the new knowledge.

Assistant: Relevance of knowledge and chain-of-thought: Support.

1327

User: Question: {*A user question*} Old Thoughts: {*An old CoT*}

New knowledge: {*A sentence containing new knowledge*} Please give me the new chain-of-thought based on the new knowledge. **Assistant:** Relevance of knowledge and chain-of-thought:

The prompt of getting final answer is as follows. We instruct the original LLM to pay more attention to the new CoT.

EditCoT Framework - Prompt of Getting Final Answer

System: Suppose the world knowledge has been updated. Answer the question based on the [New Thoughts] rather than your own knowledge. Don't analyze whether the reasoning of [New Thoughts] is correct. Output the most concise answer directly.

User: Question: In which location are the headquarters of the developer of Android Jelly Bean situated?

[New Thoughts]: Android Jelly Bean was developed by Google.

The headquarters of Google is located in the city of Googleplex.

Assistant: Answer from [New Thoughts]: Googleplex

User: Question: Who is the head of government in the city where Husky Energy's headquarters is located?

[New Thoughts]: The headquarters of Husky Energy is located in the city of Calgary.

The name of the current head of the Calgary government is Jyoti Gondek.

Assistant: Answer from [New Thoughts]: Jyoti Gondek

User: Question: Who is the head of state of the country of citizenship of Theodoros Angelopoulos?

[New Thoughts]: Theodoros Angelopoulos is a citizen of Greece.

The name of the current head of state in Greece is Aikaterini Sakellaropoulou.

Assistant: Answer from [New Thoughts]: Aikaterini Sakellaropoulou **User:** Question: Who is the head of government in the location where Vulcan Inc. is headquartered?

[New Thoughts]: The headquarters of Vulcan Inc. is located in the city of Seattle. The name of the current head of the Seattle government is Bruce Harrell.

Assistant: Answer from [New Thoughts]: Bruce Harrell

User: Question: Who is the head of government of the country where Greg Combet holds citizenship?

[New Thoughts]: Greg Combet is a citizen of Australia.

The name of the current head of the Australia government is Anthony Albanese. Assistant: Answer from [New Thoughts]: Anthony Albanese

User: Question: {*A question*} [New Thoughts]: {*A CoT*} Assistant: Answer from [New Thoughts]:

I Licensing

Qwen2.5-14B-Instruct is released under the Apache License 2.0. Meta-Llama-3-8B-Instruct is released under the META LLAMA 3 COMMU-NITY LICENSE. KN, ROME, FT, FLARE are released under the MIT license. Contriever is released under the CC BY-SA 4.0 License.

The datasets MQuAKE, LeKUBE, and DUNE are released under the MIT license. This paper's research objective is academic exploration, which aligns with the terms of this license. 1339

1340

1341

1342

1343

1344

1333

1332

1328

1329

1330