# **COMPKE: Complex Question Answering under Knowledge Editing**

**Anonymous ACL submission** 

# Abstract

Knowledge Editing-Efficiently modifying 001 the knowledge in large language models has 002 gathered great attention. Current bench-004 marks primarily use multi-hop question answering to assess and analyze newly in-005 jected or updated knowledge. However, we 006 argue that these benchmarks fail to effec-007 tively evaluate how well the updated models apply this knowledge in real-life scenar-010 ios, particularly when questions require complex reasoning, involving one-to-many rela-011 tionships or multi-step logical intersections. 012 To fill in this gap, we introduce a new benchmark, COMPKE: Complex Question Answer-014 ing under Knowledge Editing, which includes 11,924 complex questions that reflect real-016 life situations. We perform a comprehensive evaluation of four different knowledge edit-019 ing methods on COMPKE, and our results show that the performance of these methods varies between different models. For exam-021 ple, MeLLo achieves an accuracy of 39.47 on GPT-40-MINI but drops significantly to 3.83 on QWEN2.5-3B. We further analyze the rea-024 sons behind these results from both methodological and model perspectives. Our dataset will be publicly available on GitHub. 027

# 1 Introduction

029

031

032

037

039

041

Despite large language models (LLMs) are powerful in solving a wide range of real-world scenarios, they often generate outdated or incorrect knowledge (Wang et al., 2023b; Zhang et al., 2024b). Therefore, Knowledge Editing (KE), *i.e.*, updating the model's knowledge by avoiding expensive fine-tuning, has become an active research domain (Wang et al., 2023b; Zhang et al., 2024b).

To comprehensively evaluate the effectiveness of KE methods, a common approach is to assess whether the model can reproduce the newly injected knowledge, as demonstrated in ZsRE (Levy et al., 2017) and COUNTERFACT (Meng et al.,



Figure 1: (a) An example of a multi-hop question involving only one-to-one sequential step-by-step reasoning. (b) An example of a complex problem involving one-to-many knowledge mapping, logical operations, and conditional confirmation.

2022a). However, these benchmarks cannot determine whether the model genuinely utilizes the newly injected knowledge or simply memorizes and regurgitates it. MQuAKE (Zhong et al., 2023) addresses this issue through multi-hop question answering (MQA). An example in this regard is illustrated in Figure 1 (a), which shows a question: *"Who is the spouse of the president of U.S.?"* This question requires multiple reasoning steps: (*a*) identifying who is the current president of U.S.; and, (2) determining the president's spouse.

However, the evaluation dimensions of multihop questions remain too narrow to fully assess the model's ability to flexibly apply the newly integrated knowledge. This limitation manifests itself in three key areas: *(i) linear question structure*, the questions follow a strict pattern, resulting in an overly simplistic structure that can be arranged in a linear sequence. *(ii) one-to-one relations*, each fact triple in the sub-questions adheres strictly to a one-to-one relation, which fails to reflect real-world knowledge representations. In prac-

063

116

117

118

119

120

121

122

123

124

125

126

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

064tice, many facts involve one-to-many relation-065ships, such as "Who are the major shareholders of066a company?"—where a single subject is linked to067multiple entities. (iii) limited edit operations, the068knowledge edits are limited to substitutions, over-069looking more complex real-world modifications.

070

071

074

075

077

079

080

083

087

089

090

091

092

093

094

095

096

099

100

101

102

103

104

105

106

108

109

110

111

112

113

114

To address this gap, we propose a new benchmark for complex questions, *i.e.*, COMPKE: **Comp**lex Question Answering under **K**nowledge **E**diting. COMPKE, originally derived from Wikidata, comprises 11,924 complex questions. As shown in Figure 1 (b), compared to the multi-hop knowledge editing benchmark, COMPKE offers the following advantages:

(*i*) *Diverse structures*: Individual sub-questions in COMPKE are combined in multiple ways to create complex questions, integrating logical operations, conditional verification, and knowledge mapping.

(*ii*) **One-to-many relations**: The fact triples that formulate complex problems encompass both one-to-one and one-to-many relations.

(*iii*) *Expanded capabilities*: COMPKE systematically incorporates real-world knowledge updates, extending beyond simple substitutions to encompass additions and deletions.

We perform a comprehensive evaluation of major KE methods across five LLMs from different model families, including both open-source and closed-source architectures with varying parameter sizes. The results show that most methods demonstrate relatively low performance. Additionally, we analyze the effectiveness of each method across models with different parameter scales. Our findings suggest that parameter-based approaches are more effective for smaller models, whereas memory-based methods achieve better results in larger models with stronger reasoning capabilities. We summarize the key contributions of our work as follows:

- We introduce COMPKE, a novel KE benchmark that overcomes existing limitations by incorporating diverse question structures, one-to-many relations, and expanded edit types.
- We comprehensively evaluate major KE methods across five LLMs, uncovering significant differences in their ability to handle complex logical problems in diverse KE scenarios and providing an in-depth analysis of the underlying factors.

# 2 Related Work

**Knowledge Editing Benchmarks.** KE is a crucial research area for LLMs, enabling them to update information and adapt to evolving real-world queries. Various benchmarks have been established to assess the effectiveness of KE methods.

Early works like COUNTERFACT (Meng et al., 2022a) assess counterfactual updates, while ZsRE (Levy et al., 2017) and MzsRE (Wang et al., 2023c) extend evaluations to zero-shot and multilingual settings. ECBD (Onoe et al., 2023) examines whether newly injected facts can propagate reasoning across related entities. Easyedit (Wang et al., 2023a) propose an easy-to-use framework for LLMs that supports a variety of cutting-edge KE approaches. More recent works such as MQuAKE (Zhong et al., 2023), MQA-AEVAL (Ali et al., 2024) extend the evaluation to multihop reasoning under KE. TEMPLAMA (Zheng et al., 2023a) and ATOKE (Yin et al., 2023) explore the task of time-series knowledge editing, aiming to modify knowledge without affecting knowledge from other time periods. However, these benchmarks fall short in capturing realworld complexity, such as reasoning with one-tomany relations or combining entities via logical operations such as intersection and union.

Knowledge Graph Question Answering. There exist several complex question-answering datasets in the Knowledge Graph (KG) domain. ComplexQuestions (Bao et al., 2016) evaluates KGbased systems' ability to handle multi-constraint MetaQA (Zhang et al., 2018) is a queries. multi-hop dataset in the movie domain, incorporating both textual and audio data and requiring reasoning over up to three hops. ComplexWebQuestions (Talmor and Berant, 2018), built on the Freebase knowledge base, involves answering complex questions by reasoning across multiple web snippets. CR-LT-KGQA (Guo et al., 2024) focuses on commonsense reasoning and long-tail knowledge. Although complex questions have been extensively studied in the KG domain, they cannot be directly applied to the knowledge editing field due to two key challenges: (i) Omission of sub-questions and (ii) Knowledge dependency. A detailed explanation is provided in Appendix A.2.

**Knowledge Editing Methods.** We sub-divide existing research on KE into parameter-based and memory-based methods.

Parameter-based KE methods aim to directly 166 modify the model's internal parameters to reflect updated knowledge. For example, ROME (Meng et al., 2022a) and MEMIT (Meng et al., 2022b) focus on identifying and modifying parameters associated with specific knowledge, while Transformer-Patcher (Huang et al., 2023) edits facts by adding neurons. To reduce computational costs and prevent catastrophic forgetting, techniques such as: LoRA (Hu et al., 2021), Prompt Tuning (Shi and Lipani, 2024), and QLoRA (Dettmers et al., 2023) have been proposed.

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

198

199

200

201

202

203

205

206

207

209

211

212

213

214

215

However, after KE, these methods struggle with multi-hop and complex questions and cannot be applied to closed-source models like OpenAI GPTs, which are accessible only via APIs. Moreover, they are more computationally expensive than memory-based approaches.

Memory-Based methods store updates in external memory and retrieve them as needed during inference. For instance, SERAC (Mitchell et al., 2022) combines semi-parametric editing with retrieval augmented counterfactual models for efficient knowledge updates. GRACE (Hartvigsen et al., 2022) integrates adapters into LLMs and uses vector matching to modify knowledge entries. IKE (Zheng et al., 2023b) applies in-context learning with stored demonstrations for knowledge modification, MeLLo (Zhong et al., 2023) stores edited facts externally and utilizes prompts to incorporate edits during inference. PokeMQA (Gu et al., 2023) separates question decomposition and conflict detection using a two-stage programmable scope detector. GLAME (Zhang et al., 2024a) employs a knowledge graph module to enhance retrieval efficiency.

We observed, MeLLo and PokeMQA excel at multi-hop problems, therefore in our experiments, we use them as baselines to assess the generalization of memory-based methods to complex questions. We provide further details about related work in Appendix A.

#### **Preliminaries** 3

We use  $\mathcal{D} = \{(s, r, o)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$  to denote the set of knowledge triplets, where  $\mathcal{E}$  and  $\mathcal R$  denote the set of entities and relations respectively. Each triple (s, r, o) represents a knowledge instance, implying that the subject entity s and the object entity o are related by relation r. In order to represent one-to-many knowledge instances, we



Figure 2: An example of a complex question under knowledge editing. Knowledge editing occurs in the first sub-question, where the filming location of Christime is modified from { Los Angeles} to {San Francisco, Los Angeles, New York}.

expand the original definition of knowledge instance to  $(s, r, \mathcal{O})$ , where  $\mathcal{O} = \{o_1, o_2, \cdots\}$  is a set of object entities, e.g., (Avatar, actors\_are, {Worthington, Saldana,  $\cdots$ }).

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

# 3.1 Complex Questions

Building on the example in the introduction, we formally define the complex questions to be explored in this paper. A brief recap of multi-hop question answering (MQA) and MQA under KE is provided in the Appendix B.1. We define a complex question Q as a question that could be represented as a graph-like reasoning structure, i.e.,  $Q = (\mathbf{S}, \mathbf{L})$ , where  $\mathbf{S} = \{S_1, S_2, \cdots\}$  represents a set of *intermediate entities* and  $\mathbf{L} = \{L_1, L_2, \cdots\}$ denotes a set of *reasoning links*. Each  $S_i \in \mathbf{S}$  is a set of entities, *i.e.*,  $S_i = \{s_1, \dots\}$ , used to represent one-to-one and one-to-many knowledge instances. Each  $L_i \in \mathbf{L}$  is a reasoning link. Note that unlike the relation typically used in knowledge graphs (used to map one entity  $s_i$  to another  $s_i$  via the relation r), reasoning links offer extended operations by allowing conditional confirmation and logical operations, which are formally explained below.

Reasoning Links. We categorize the reasoning links into two distinct categories:

(a) Knowledge-related Links: These links facilitate entity traversal along the link, e.g., given a set of entities  $S_i \in \mathbf{S}$ , a reasoning link may be used to obtain the next step entities  $S_i \in \mathbf{S}$ . We subdivide these links into:

For  $S_i$ , we consider (i) Knowledge Mapping.

248a knowledge mapping link as a mapping to the249set of adjacent entities  $S_j = \bigcup_{s \in S_i} A_r(s)$ , where250 $A_r(s) = \{s' \mid (s, r, s') \in \mathcal{D}\}$  represents the enti-251ties related to s via relation r.

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

270

271

273

274

275

276

277

278

279

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

(ii) Condition Confirmation. Given r and s', this link aims to identify a set of entities  $S_j = \{s \in S_i \mid C(s, r, s') = \text{True}\}$  conforming to the condition  $[C(s, r, s') \in D]$ , which is used to examine whether s can obtain s' via r.

(b) Logical Links: Given a set of intermediate entities  $\{S_1, S_2, \dots, S_n\} \in \mathbf{S}$ , these reasoning links perform logical operations among the elements of  $S_i$ . Specifically, we use logical links for the following operations:

(i) Intersection. The intersection operation identifies adjacent entities shared across all sets, *i.e.*,  $S_j = \bigcap_{k=1}^n S_k$ .

(*ii*) Union. The union operation computes the set of adjacent entities, including all entities from any of the sets, *i.e.*,  $S_j = \bigcup_{k=1}^n S_k$ .

**Example 1.** An example of a complex question with reasoning links is illustrated in Figure 2. It shows the question: "Where were the movies Christine and Pacific Heights filmed?". The intermediate entities are  $S_1$ ={Christine};  $S_2$ ={Pacific Heights};  $S_3$ = {Los Angeles};  $S_4$ ={San Franciso, Los Angeles}; and  $S_5$ ={Los Angeles}. The reasoning operations are:  $L_1 : S_1 \xrightarrow{\text{filming}_{at}} S_3$ ;  $L_2 : S_2 \xrightarrow{\text{filming}_{at}} S_4$ ; followed by  $L_3$  : logical operation on  $S_3$  and  $S_4$  ( $S_3 \cap S_4$ ) to obtain the final answer, *i.e.*,  $S_5$  = {Los Angeles}.

**Complex Question Answering under KE.** We use  $e = (s, r, \mathcal{O} \rightarrow \mathcal{O}')$  to represent knowledge editing for one-to-many instances showing that  $\mathcal{O}$  is updated to  $\mathcal{O}'$ . The task assumes that the language model has access to original knowledge base  $\mathcal{D}$ . Given a batch of edits  $\mathcal{E} = \{e_1, e_2, \cdots\}$ , the knowledge to be deleted as  $\mathcal{D}_{del}^{\mathcal{E}} = \{(s_i, r_i, \mathcal{O}_i) \mid e_i \in \mathcal{E}\}$ , and the newly added knowledge as  $\mathcal{D}_{add}^{\mathcal{E}} = \{(s_i, r_i, \mathcal{O}_i') \mid e_i \in \mathcal{E}\}$ , the end-goal is to update the LLM's knowledge by  $\mathcal{D}'$ , define as:  $\mathcal{D}' = (\mathcal{D} - \mathcal{D}_{del}^{\mathcal{E}}) \cup \mathcal{D}_{add}^{\mathcal{E}}$ . Finally, updated knowledge  $\mathcal{D}'$  is used to generate the final answer to the complex question Q.

# 4 COMPKE

While complex questions are common in reallife scenarios, they remain underexplored in LLM question answering under KE. We argue that existing benchmarks predominantly focus on linear multi-hop questions, limiting their effectiveness in evaluating complex queries. To bridge this gap, we propose COMPKE: **Comp**lex Question Answering under **K**nowledge **E**diting. In the following, we provide a brief overview of COMPKE followed by a detailed flow of the process.

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

# 4.1 Dataset Construction

**Overview.** The workflow of our data construction process, illustrated in Figure 3, follows six key steps. First, we extract factual triples from Wikidata. Next, we select relevant relations and sample triples corresponding to those relations. In the third step, these triples are combined into complex questions with diverse reasoning structures, and edits are introduced at appropriate positions within the questions. This is followed by the introducing counterfactual modifications in step four and the filtering of conflicting instances in step five to maintain consistency. Finally, in step six, the structured questions are converted into natural language. Further details on each step are provided below.

**Step 1: Collecting Relation Templates.** In step 1, we carefully select 37 relations from Wikidata's *List of Properties*, through a two-step process. First, we identify essential one-to-many relations (e.g., family-child, book-authors, movie-actors) for one-to-many knowledge mapping. Next, we incorporate one-to-one relations (e.g., country-capital, person-hometown) that capture fundamental entity attributes, enabling one-to-one knowledge mapping and conditional confirmation. In addition, we prioritize relations commonly encountered in everyday scenarios to enhance the practical utility of the data set for the relevance of the real world. The full list of relation templates used in COMPKE is provided in Appendix Table 9.

Step 2: Sampling Facts. After selecting relation templates, we build the knowledge base  $\mathcal{D}$ with a focus on commonly known rather than obscure knowledge. Using the collected relation templates, we sample single-hop knowledge triples from Wikidata and rank them by access frequency, prioritizing the most frequently accessed triples. To refine this selection, we employ GPT-3.5-TURBO-INSTRUCT to filter out knowledge the model cannot recall. The finalized knowledge base  $\mathcal{D}$  serves as the basis for generating complex questions.



Figure 3: The construction process of COMPKE.

**Step 3: Constructing Complex Questions.** We observe that complex questions often follow structured reasoning patterns, as shown in Figure 2, where knowledge mapping is followed by logical operations (e.g., intersection). To systematically collect these reasoning structures, we first manually curate a high-quality subset of complex questions as seed examples. We then extract their underlying reasoning structures by removing intermediate entities, forming reusable templates. These templates are instantiated with real-world facts from  $\mathcal{D}$  to generate specific complex questions. The process begins with the random initialization of the leaf nodes, followed by the iterative identification of intermediate entities using logical operations or knowledge of  $\mathcal{D}$ , continuing until all entities are fully determined.

346 347

348

349

351

354

355

357

360

361

363

364

365

367

371

To ensure the practical relevance of instantiated questions, we filter out cases that meet the following criteria: *(i)* questions with no valid answer, *(ii)* questions yielding an empty set of intermediate entities, and *(iii)* cases where entities involved in logical operations are of incompatible types. For illustration, exemplar relational structures and their corresponding instantiated complex questions are provided in Appendix (Figure 10).

Step 4: Introducing Counterfactual Edits. To 372 simulate knowledge edits, we construct counter-373 factual knowledge updates. For each complex 374 question, we randomly select knowledge mappings with knowledge of the form: (s, r, O) and 376 introduce an edit e = (s, r, O'). Unlike previ-377 ous benchmarks that only involve one-to-one relations with edits limited to entity replacement, our dataset introduces edits that involve one-to-380 many relations, leading to more complex edits. To 381 clearly represent the changes in a fact triple, we 382 define three basic operations, each of which can be combined to form an edit: 384

Е	(Christine, filming location, { Los Angeles}→ {San Francisco, Los Angeles, New York})
Q	<ul> <li>i)Where were the movies Christine and Pacific Heights filmed?</li> <li>ii)In which locations were both the movie Christine and Pacific Heights filmed?</li> <li>iii)What were the filming locations for both the movie Christine and Pacific Heights?</li> </ul>
$\mathcal{A} \\ \mathcal{A}^*$	{Los Angeles} {San Francisco, Los Angeles}
$\mathcal{T}$ $\mathcal{T}^*$	(Christine, filming location, {Los Angeles}) (Pacific Heights, filming location, {San Francisco, Los Angeles}) (Christine, filming location, {San Francisco, Los Angeles, New York}) (Pacific Heights, filming location, {San Francisco, Los Angeles})
$\mathcal{L} \\ \mathcal{L}^*$	{Los Angeles} ∩ {San Francisco, Los Angeles} = {Los Angeles} {San Francisco, Los Angeles, New York} ∩ {San Francisco, Los Angeles} ={San Francisco, Los Angeles}

Table 1: A case from COMPKE, illustrating the components involved in question editing. Here,  $\mathcal{E}$  represents the edit,  $\mathcal{Q}$  is the natural language question,  $\mathcal{A}$  and  $\mathcal{A}^*$  denote the answers before and after editing respectively.  $\mathcal{T}$  and  $\mathcal{T}^*$  are the sets of fact triples before and after editing, which form the complex question. Additionally,  $\mathcal{L}$  and  $\mathcal{L}^*$  indicate the logic operations applied to the question before and after editing.

(*i*) Addition:  $\mathcal{O}_{add} = \mathcal{O}' \setminus \mathcal{O}$ , where  $\mathcal{O}_{add}$  represents the set of newly added entities;

385

386

387

388

389

390

391

392

393

394

395

396

397

399

400

401

402

403

<u>(*ii*) Deletion:</u>  $\mathcal{O}_{del} = \mathcal{O} \setminus \mathcal{O}'$ , where  $\mathcal{O}_{del}$  represents the set of removed entities;

(*iii*) Retention:  $\mathcal{O}_{ret} = \mathcal{O} \cap \mathcal{O}'$ , where  $\mathcal{O}_{ret}$  represents the set of retained entities.

**Example 2.** An example of an edit to change the management of the "Microsoft" may be expressed as: (Microsoft, managers\_are, {John, Smith, Dave}  $\rightarrow$  {Smith, Eden, Keyes}), which involves deleting {John}, retaining {Smith}, and adding {Eden, Keyes}.

Step 5: Filtering Conflicting Edits. Since the counterfactual edits in Step 4 are introduced randomly, for a batch of edits  $\mathcal{E} = \{e_1, e_2, ...\}$  there may be edits corresponding to different cases where  $e_i = (s_i, r_i, \mathcal{O}_i \to \mathcal{O}_i^*)$  and  $e_j = (s_j, r_j, \mathcal{O}_j \to \mathcal{O}_j^*)$ , with  $s_i = s_j$  and  $r_i = r_j$ , but  $\mathcal{O}_i^* \neq \mathcal{O}_j^*$ . This indicates the presence of 404 conflicting facts within the batch. Simultaneously
405 introducing such conflicts can severely compro406 mise the validity evaluation. To mitigate this, we
407 identify and group conflicting cases, and then ran408 domly select only one to retain.

Step 6: Phrasing in Natural Language. Build-409 410 ing on steps 1-5, we construct complex questions involving edits, each comprising multiple 411 fact triples. To facilitate evaluation by the tar-412 get LLMs, these questions must be translated 413 into natural language. For each reasoning struc-414 ture defined in Step 3, we first manually curate 415 eight high-quality examples. Then, using GPT-416 40-mini, we generate three natural language ques-417 418 tions for each structured question. Further details on constructing the dataset are provided in the Ap-419 pendix C. 420

# 4.2 Dataset Summary

421

441

442

443

444

445

Table 2 presents the dataset distribution across 422 two dimensions: Edit\_num and Step\_num. 423 Edit\_num represents the number of triples 494 edited in a complex question. In COMPKE, most 425 cases involve editing a single triple, followed by 426 cases with two edits. Step\_num denotes the 427 number of reasoning steps required to answer the 428 complex question, with 3-step questions being 429 the most prevalent, followed by 4-step and 5-step 430 questions respectively. 431

**Example 3.** Table 1 presents a detailed example 432 from COMPKE, illustrating a complex question 433 constructed by combining two sub-questions with 434 an intersection operation. We assume that the edit-435 ing takes place in the first sub-question (i.e., Chris-436 tine's filming locations are updated from {Los 437 Angeles} to {San Francisco, Los Angeles, New 438 York}), leading to the addition of San Francisco in 439 the final answer. 440

# 5 Experiments

In this section, we conduct a comprehensive evaluation of recent knowledge editing methods in COMPKE, assessing them from three aspects: whether newly added knowledge can be recalled,

#Edits	1	2	3	4	5	Total
Edit_num	9,697	1,118	998	103	8	11,924
Step_num	200	424	5,770	2,949	2,581	11,924

Table 2: Statistical Results of COMPKE.

whether existing knowledge is retained, and overall accuracy. We also analyze how the performance of different methods changes when the edit batch size (*i.e.*, the number of edits performed at once) increases. Additionally, by case studies, we observe several interesting phenomena, including overfitting in parameter-based methods, model collapse when increasing edit batch size, and the omission phenomenon in memory-based methods. 446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

# 5.1 Experimental Settings

Language Models. We conduct experiments using five different target LLMs corresponding to three model families. For open source models, we select LLAMA-3.1-8B-INSTRUCT (Abhimanyu Dubey et al., 2024), QWEN2.5-3B-INSTRUCT (Team, 2024), QWEN2.5-7B-INSTRUCT (Team, 2024). For closed source models, we select GPT-3.5-TURBO and GPT-4O-MINI (Achiam et al., 2023).

**Baselines.** For performance comparison, we use the best performing methods for MQA under KE as baselines. These include parameterbased variants: ROME (Meng et al., 2022a), and MEMIT (Meng et al., 2022b); and memorybased variants: MeLLo (Zhong et al., 2023), and PokeMQA (Gu et al., 2023). Since GPT-3.5-TURBO and GPT-4O-MINI can only be accessed through APIs, parameter-based knowledge editing methods cannot be applied to them.

**Evaluation Metrics.** We use the following metrics for evaluation:

(*i*) Augment Accuracy (Aug): The number of newly introduced entities added to the answer list after the knowledge edit that are correctly identified compared to the original list.

(*ii*) Retain Accuracy (Ret): The number of entities present in both the original and edited answer lists, indicating the model's ability to preserve unmodified knowledge.

(*iii*) Accuracy (Acc): The average of Aug and Ret, offering a holistic measure of the model's accuracy in answering complex questions under KE. The detailed mathematical formulations of these metrics are provided in Appendix D.3.

**Example 4.** Following the example in Figure 2, the final answer changes from {Los Angeles} before editing to {San Francisco, Los Angeles} after editing. Aug evaluates whether the model correctly includes the newly added entity, {San Francisco}, while Ret assesses its ability to retain

Model	Method	1-edited			100-edited			3000-edited		
		Aug	Ret	Acc	Aug	Ret	Acc	Aug	Ret	Acc
	ROME	12.61	17.91	15.26	4.80	4.40	4.60	0.82	1.59	1.21
OWEN 2 5 2D	MEMIT	20.99	23.86	22.43	7.80	6.73	7.27	1.52	3.75	2.64
QWEN2.5-3D	MeLLo	5.40	2.25	3.83	3.06	3.39	3.23	0.69	2.00	1.35
	PoKeMQA	4.26	1.85	3.06	2.85	1.38	2.12	0.71	0.62	0.67
	ROME	22.82	25.09	23.96	7.50	7.98	7.74	0.73	0.98	0.86
OWEN 2 5 7B	MEMIT	29.40	27.72	28.56	24.11	24.80	24.46	1.88	2.05	1.97
QWEN2.5-7D	MeLLo	17.78	13.38	15.58	10.35	17.32	13.84	8.98	12.59	10.79
	PoKeMQA	15.59	11.41	13.50	8.17	13.67	10.92	5.04	9.15	7.10
	ROME	7.44	24.84	16.14	1.50	1.14	1.32	0.56	0.61	0.59
II AMA 21 0D	MEMIT	4.90	33.22	19.06	5.00	29.27	17.14	5.03	29.20	17.12
LLAMA-3.1-0D	MeLLo	14.06	17.95	16.00	9.17	17.84	13.51	8.98	14.17	11.58
	PoKeMQA	11.40	15.10	13.25	8.87	16.85	12.86	7.45	12.73	10.09
CDT 2 5 TUDDO	MeLLo	49.21	44.88	47.05	37.10	44.09	40.60	32.61	38.58	35.60
GP 1-3.5-10KB0	PoKeMQA	23.20	25.15	24.18	21.47	23.28	22.38	20.20	22.20	21.20
CPT 40 MINI	MeLLo	22.07	25.19	23.63	20.31	23.62	21.96	18.75	22.14	20.45
GF 1-40-MINI	PoKeMQA	36.60	42.33	39.47	35.42	41.35	38.39	28.36	35.02	31.69



Figure 4: Variation of *Accuracy* (Acc) across QWEN2.5-3B, QWEN2.5-7B, and LLAMA-3.1-8B models with varying edit numbers. Results for GPT-3.5-TURBO and GPT-40-MINI are provided in Appendix E.2.

the existing entity, {Los Angeles}, which appears both before and after editing. Acc, as the average of Aug and Ret, measures the model's effectiveness in integrating new knowledge while preserving existing information.

**Experiment Setup.** We conduct experiments on varying scales of knowledge edits, *i.e.*, using a batch of k-edits at a time with k = $\{1, 100, 1000, 3000\}$ . To ensure a fair comparison with existing memory-based methods, we use the decomposition examples of complex questions for MeLLo and PokeMQA, as prompts. Additional details on the experimental setting are provided in Appendix D.

#### 510 5.2 Experimental Results

496

497

498

499

500

501

502

503

504

505

506

507

508

509

511The experimental results are summarized in Ta-512ble 3. In general, MeLLo achieves the highest per-513formance in the 1-edit setting on GPT-3.5-Turbo,514yielding Aug score = 49.21. When comparing

different approaches, memory-based methods perform poorly on smaller models (*e.g.*, QWEN2.5-3B) due to their reliance on the instructions' following and reasoning capabilities. In contrast, parameter-based methods are more effective for smaller models, but suffer substantial performance degradation as the edit batch size increases. In the following, we provide a detailed analysis of these findings.

**Batch Editing (#***k***-edits).** Figure 4 illustrates the accuracy of the four methods on QWEN2.5-3B, QWEN2.5-7B, and LLAMA-3.1-8B with an increase in the number of edits. Performance variations for GPT-3.5-TURBO and GPT-40-MINI are provided in Appendix E.2.

We observe that memory-based methods experience a gradual decline in performance as the number of edits (k) increases. In contrast, parameterbased methods degrade more rapidly, especially when the number of edits exceed a certain thresh-

old. Notably, when  $k \ge 100$ , the model loses coherence, producing inconsistent responses and generating irrelevant outputs, as detailed in the Appendix Figure 9.

Smaller Models. For smaller models, such as 539 QWEN2.5-3B, memory-based methods perform 540 poorly compared to parameter-based methods. 541 This can be attributed to two key factors: (i) 542 Smaller models have limited instruction-following 543 capabilities and struggle to adhere to the required 544 format for response planning. (ii) During the 545 problem-solving process, these models fail to ef-546 fectively integrate their internal knowledge with 547 external edits, making it difficult to address dif-548 ferent sub-questions. 549

550

551

552

553

554

555

556

557

559

560

561

562

563

564

565

566

567

568

570

571

572

573

574

575

576

577

578

A notable example is the baseline model: PokeMQA, which relies heavily on instructionfollowing capabilities and performs poorly on both LLAMA-3.1-8B and QWEN2.5-3B. This highlights the importance of an effective decomposition mechanism that does not depend on strong instruction follow-up abilities, particularly for smaller models, as it plays a crucial role in overall performance.

**Overfitting of parameter-based methods.** Our experiments reveal that parameter-based methods perform remarkably well on models with smaller parameter sizes. For example, in the Qwen2.5-3B (1-edited) setting, MEMIT achieves a significantly higher accuracy score of 22.43, compared to 3.83 for MeLLo. This result is unexpected, as prior research suggests that memory-based methods generally exhibit better generalization than parameterbased approaches. To investigate this discrepancy, we conducted a detailed case study and found that MEMIT's high accuracy is primarily driven by model overfitting.

Specifically, after injecting modified knowledge, the model consistently outputs the newly introduced information whenever it encounters related questions, even in contexts where it is not appropriate. The example in Figure 8 provides a detailed explanation of why this phenomenon leads to a higher augmentation bias.

579Omission Phenomenon. We also analyze the per-580formance of MeLLo using the original prompts581provided with the model implementation. We ob-582serve that it leads to omission phenomenon in583the decomposition for complex questions, *i.e.*, the584MeLLo's decomposition plan skips certain steps,585specifically the logical intersection part. Under-



Figure 5: Performance comparison of MeLLo and PoKeMQA on the MQuAKE-T, MQuAKE-CF-3k, and COMPKE datasets on GPT-40-MINI, with COMPKE presenting more challenging than previous datasets.

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

612

613

614

615

616

617

618

619

lying justification in this regard is the fact that the conditional confirmation operations, *e.g.*, logical intersection, does not appear in the multi-hop questions. This showcases that the generalization of decomposition operation through prompt examples is insufficient, highlighting the essence of incorporating examples similar to the question being decomposed. An example illustration in this regard is provided in Appendix Table 8.

**Comparision with other Datasets.** We select two popular KE datasets (MQuAKE-T and MQuAKE-CF-3k) for comparison with our dataset. Using GPT-4o-mini as the test model, we evaluate the performance of two methods, MeLLo and PoKeMQA, on these datasets (detailed information about MQuAKE and its evaluation metric are shown in appendix D.1 and D.3) and compare them with COMPKE. The results are illustrated in Figure 5. Both methods exhibit lower accuracy on COMPKE than on MQuAKE datasets, indicating that COMPKE presents a greater challenge then previous ones.

# 6 Conclusion

In this paper, we introduce the concept of complex questions in the context of knowledge editing and propose a new benchmark, COMPKE. Through a comprehensive evaluation of various knowledge editing methods on COMPKE, we find that existing approaches struggle when dealing with complex question scenarios. We analyze the cause of these limitations and suggest that future work can leverage our dataset and evaluation framework to develop more robust and generalizable knowledge editing methods.

716

717

718

719

720

721

722

723

#### Limitations 620

622

623

624

625

626

627

629

630

631

632

633

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

663

664

665

This work poses following limitations: 621

- In COMPKE, edits are randomly introduced through counterfactual modifications, which may result in discrepancies from actual/realworld modifications.
  - The fact triples in COMPKE are restricted to one-to-one and one-to-many relations, excluding many-to-many and many-to-one relationships.

#### **Ethics Statement**

This work directly deals with updating the capability and/or editing the knowledge of large models. It has the potential for abuse, such as adding poisonous misinformation, malicious content, bias 634 etc. Keeping in view these concerns, we highlight 635 this work must not be used under critical settings. 636

# References

- Abhinav Jauhri Abhimanyu Dubey et al. 2024. The llama 3 herd of models. ArXiv, abs/2407.21783.
- 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.
  - Muhammad Asif Ali, Nawal Daftardar, Mutayyaba Waheed, Jianbin Oin, and Di Wang. 2024. Mgakeal: Multi-hop question answering under knowledge editing for arabic language.
  - Junwei Bao, Nan Duan, Zhao Yan, Ming Zhou, and Tiejun Zhao. 2016. Constraint-based question answering with knowledge graph. In Proceedings of COLING 2016, the 26th international conference on computational linguistics: technical papers, pages 2503-2514.
  - Yuchen Cai, Ding Cao, Rongxi Guo, Yaqin Wen, Guiquan Liu, and Enhong Chen. 2024. Editing knowledge representation of language model via rephrased prefix prompts.
  - Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. 2023. Evaluating the ripple effects of knowledge editing in language models.
  - Fernanda De La Torre, Cathy Mengying Fang, Han Huang, Andrzej Banburski-Fahey, Judith Amores Fernandez, and Jaron Lanier. 2024. Llmr: Real-time prompting of interactive worlds using large language models. In Proceedings of the CHI Conference on Human Factors in Computing Systems, pages 1–22.

- Jingcheng Deng, Zihao Wei, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. 2024. Everything is editable: Extend knowledge editing to unstructured data in large language models.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms.
- Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. 2024. A primer on the inner workings of transformer-based language models. arXiv preprint arXiv:2405.00208.
- Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ninghao Liu, Ruobing Wang, and Xin Wang. 2023. Pokemqa: Programmable knowledge editing for multi-hop question answering. arXiv preprint arXiv:2312.15194.
- Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024. Model editing harms general abilities of large language models: Regularization to the rescue.
- Willis Guo, Armin Toroghi, and Scott Sanner. 2024. Cr-lt-kgqa: A knowledge graph question answering dataset requiring commonsense reasoning and longtail knowledge. arXiv preprint arXiv:2403.01395.
- Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024. Model editing at scale leads to gradual and catastrophic forgetting. arXiv preprint arXiv:2401.07453.
- Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2022. Aging with grace: Lifelong model editing with discrete key-value adaptors. ArXiv, abs/2211.11031.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2024a. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. Advances in Neural Information Processing Systems, 36.
- Peter Hase, Thomas Hofweber, Xiang Zhou, Elias Stengel-Eskin, and Mohit Bansal. 2024b. Fundamental problems with model editing: How should rational belief revision work in llms? arXiv preprint arXiv:2406.19354.
- Bing He, Mustaque Ahamad, and Srijan Kumar. Reinforcement learning-based counter-2023.misinformation response generation: A case study of covid-19 vaccine misinformation. In Proceedings of the ACM Web Conference 2023, WWW '23, page 2698-2709, New York, NY, USA. Association for Computing Machinery.
- Cheng-Hsun Hsueh, Paul Kuo-Ming Huang, Tzu-Han Lin, Che Wei Liao, Hung-Chieh Fang, Chao-Wei Huang, and Yun-Nung Chen. 2024. Editing the mind of giants: An in-depth exploration of pitfalls of knowledge editing in large language models. In

724

725

- 757 758 762 763 764 765
- 766 767 768
- 769 770
- 771 772 773
- 774
- 777

- Findings of the Association for Computational Linguistics: EMNLP 2024, pages 9417-9429, Miami, Florida, USA. Association for Computational Linguistics.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.
- Wenyue Hua, Jiang Guo, Mingwen Dong, Henghui Zhu, Patrick Ng, and Zhiguo Wang. 2024. Propagation and pitfalls: Reasoning-based assessment of knowledge editing through counterfactual tasks.
- Baixiang Huang, Canyu Chen, Xiongxiao Xu, Ali Payani, and Kai Shu. 2024. Can knowledge editing really correct hallucinations? arXiv preprint arXiv:2410.16251.
  - Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformerpatcher: One mistake worth one neuron.
  - Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), pages 333-342, Vancouver, Canada. Association for Computational Linguistics.
- Jun-Yu Ma, Jia-Chen Gu, Zhen-Hua Ling, Quan Liu, and Cong Liu. 2024. Untying the reversal curse via bidirectional language model editing.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. Advances in Neural Information Processing Systems, 35:17359–17372.
- Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2022b. Massediting memory in a transformer. In The Eleventh International Conference on Learning Representations.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. ArXiv. abs/2206.06520.
- Kento Nishi, Maya Okawa, Rahul Ramesh, Mikail Khona, Hidenori Tanaka, and Ekdeep Singh Lubana. 2024. Representation shattering in transformers: A synthetic study with knowledge editing.
- Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. 2024. What does the knowledge neuron thesis have to do with knowledge? arXiv preprint arXiv:2405.02421.
- Yasumasa Onoe, Michael J. Q. Zhang, Shankar Padmanabhan, Greg Durrett, and Eunsol Choi. 2023. Can lms learn new entities from descriptions? challenges in propagating injected knowledge.

Hao Peng, Xiaozhi Wang, Chunyang Li, Kaisheng Zeng, Jiangshan Duo, Yixin Cao, Lei Hou, and Juanzi Li. 2024. Event-level knowledge editing.

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

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

- Amit Rozner, Barak Battash, Lior Wolf, and Ofir Lindenbaum. 2024. Knowledge editing in language models via adapted direct preference optimization. arXiv preprint arXiv:2406.09920.
- Zhengxiang Shi and Aldo Lipani. 2024. Dept: Decomposed prompt tuning for parameter-efficient finetuning.
- Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 641-651.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu, Yong Jiang, Pengjun Xie, et al. 2024. Knowledge mechanisms in large language models: A survey and perspective. arXiv preprint arXiv:2407.15017.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bo Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Yuansheng Ni, Guozhou Zheng, and Huajun Chen. 2023a. Easyedit: An easy-to-use knowledge editing framework for large language models. ArXiv, abs/2308.07269.
- Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, et al. 2023b. Knowledge editing for large language models: A survey. arXiv preprint arXiv:2310.16218.
- Weixuan Wang, Barry Haddow, and Alexandra Birch. 2023c. Retrieval-augmented multilingual knowledge editing.
- Wanli Yang, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, and Xueqi Cheng. 2024. The butterfly effect of model editing: Few edits can trigger large language models collapse.
- Yunzhi Yao, Ningyu Zhang, Zekun Xi, Mengru Wang, Ziwen Xu, Shumin Deng, and Huajun Chen. 2024. Knowledge circuits in pretrained transformers. arXiv preprint arXiv:2405.17969.
- Xunjian Yin, Jin Jiang, Liming Yang, and Xiaojun Wan. 2023. History matters: Temporal knowledge editing in large language model.
- Mengqi Zhang, Xiaotian Ye, Qiang Liu, Pengjie Ren, Shu Wu, and Zhumin Chen. 2024a. Knowledge graph enhanced large language model editing.

Ningyu Zhang, Yunzhi Yao, Bo Tian, Peng Wang, 829 830 Shumin Deng, Meng Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, 831 832 Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, 833 Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiao-Jun Zhu, Jun Zhou, and Huajun 834 Chen. 2024b. A comprehensive study of knowl-835 edge editing for large language models. ArXiv, 836 abs/2401.01286. 837

838

839

840

841

842

843

844

845 846

847

848

849

850

851

852

853

854

855 856

857

- Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander Smola, and Le Song. 2018. Variational reasoning for question answering with knowledge graph. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023a.
  Can we edit factual knowledge by in-context learning? In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 4862–4876, Singapore. Association for Computational Linguistics.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023b. Can we edit factual knowledge by in-context learning? *ArXiv*, abs/2305.12740.
- Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions. *arXiv preprint arXiv:2305.14795*.

# A Related Work

859

860

861

862

864

865

866

867

868

870

871

872

873

874

876

877

878

879

880

882

883

884

885

886

889

890 891

892

893

894

895

896

897

900

901

902

903

904

905

906

907

908

# A.1 More detailed Related Work

Besides benchmarks, many researchers in recent years have explored knowledge editing from various perspectives. There is a type of research that aim to understand the working mechanisms of knowledge editing techniques, such as the relationship between model parameter localization and editing (Wang et al., 2024; Niu et al., 2024; Hase et al., 2024a,b; Ferrando et al., 2024; Gupta et al., 2024; Yao et al., 2024). For example, causal tracing does not effectively indicate the optimal editing location (Hase et al., 2024a), and some researchers have also employed computation graph to uncover the specific impacts on the model's internal behavior of knowledge editing (Yao et al., 2024). Another line of research focuses on enhancing the effectiveness of knowledge editing in specific scenarios (Rozner et al., 2024; Ma et al., 2024; De La Torre et al., 2024; Huang et al., 2024; Deng et al., 2024; Peng et al., 2024; Cai et al., 2024). For instance, bidirectional relationship modeling has been proposed to address consistency issues in bidirectional models (Ma et al., 2024), while real-time knowledge editing methods have been developed to adapt to dynamic environments where knowledge evolves frequently (De La Torre et al., 2024). Additionally, this paper focuses on exploring knowledge editing in the context of complex logical reasoning. Also some studies focus on addressing the side effects of knowledge editing techniques (Hsueh et al., 2024; Gu et al., 2024; He et al., 2023; Hua et al., 2024; Yang et al., 2024; Cohen et al., 2023; Nishi et al., 2024).

# A.2 Drawbacks of KGQA Datasets

Although complex questions have been extensively studied in the KG domain, they cannot be directly applied to the knowledge editing field due to two key challenges:

(*i*) Omission of sub-questions. These datasets do not explicitly provide sub-questions of complex questions. For example, ComplexQuestions(Bao et al., 2016) only includes only the question and its final answer, while ComplexWebQuestions(Talmor and Berant, 2018) provides only a SPARQL statement for each complex question. However, KE requires modifications at the subquestion level. Without explicitly defined subquestions, introducing targeted edits becomes im-

practical.	909
(ii) Knowledge dependency. These data sets do not	910
require models to rely on LLMs' intrinsic knowl-	911
edge to generate answers, while KE heavily relies	912
on model's internal knowledge. Directly adopting	913
these datasets risks introducing unlearned knowl-	914
edge into the evaluation, leading to unreliable an-	915
swers regardless of editing success. In construct-	916
ing COMPKE, we mitigate this by filtering out	917
knowledge instances that cannot be recalled by the	918
model.	919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

# **B** Additional Preliminaries

#### **B.1** Multi-hop Question Answering

A multi-hop question can be represented as  $s_1 \xrightarrow{r_1} s_2 \cdots \xrightarrow{r_{n-1}} s_n$ , continuously mapping one entity to another. For example, consider the question "Who is the spouse of president of U.S.", it an be represented as U.S.  $\xrightarrow{\text{president is}}$  Donald Trump  $\xrightarrow{\text{spouse is}}$  Melania Trump.

# B.2 Multi-hop Question Answering under KE.

We use  $e = (s, r, o \rightarrow o')$  to represent a knowledge edit indicating that the object entity of subject s with relation r is updated from o to o'. This task is to solve multi-hop questions under a batch of knowledge edits  $\mathcal{E} = \{e_1, e_2, \dots\}$ .

# B.3 MQA with Complex Question Answering.

We consider the previously studied linear multihop questions as a special case of complex questions involving continuous mapping of entity through a series of relational links, forming a oneway graph chain:  $S_1 \xrightarrow{L_1} S_2 \xrightarrow{L_2} \cdots \xrightarrow{L_{n-1}} S_n$ , where *n* represents the number of reasoning hops. Note that compared to complex questions, here the intermediate set  $S_i$  only encompasses a single entity, and  $L_i$  only covers one-to-one relation mapping.

# **C** COMPKE (Additional Details)

Figure 3 shows the process by which we construct complex question. Figure 10 gives some examples of the structures in COMPKE and the corresponding decomposition methods. Table 7 gives the SPARQL which we used to sample facts from WikiData. Table 6 presents the prompt used for converting structured triples into natural lan-

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

991

992

993

994

954 955

# 956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

# **D** Additional Experimental Settings

tion counts across triplets in COMPKE.

guage. Figure 6 displays the distribution of rela-

# D.1 MQuAKE

The existing data MQUAKE includes two datasets: MQUAKE-CF-3K, which is based on counterfactual editing, and MQUAKE-T, which is based on real-world changes. These datasets cover k-hop questions ( $k \in \{2, 3, 4\}$ ), each associated with one or more edits. Statistics are presented in Table 4.

Datasets	#Edits	2-hop	3-hop	4-hop	Total
	1	513	356	224	1,093
	2	487	334	246	1,067
MQUAKE-CF-3K	3	-	310	262	572
	4	-	-	268	268
	All	1,000	1,000	1,000	3,000
MQUAKE-T	1	1,421	445	2	1,868

Table 4: Statistics of the MQUAKE dataset.

#### D.2 Baselines

**ROME.** ROME by Meng et al. (2022a) uses a locate-then-edit paradigm. For a specific knowl-edge editing, ROME employs causal tracing to pin-point the exact layer of the MLP module within the Transformer model architecture that encodes the paticular factual association. Then it will perform a rank-one modification on the identified layer.

974MEMIT. MEMIT by Meng et al. (2022b) is an975evolution of ROME to transcend the inherent lim-976itation that ROME can only edit a single fact at a977time. At a time, MEMIT can identify and modify978multiple layers in a single pass, allowing for the979simultaneous editing of numerous facts.

MeLLo. MeLLo by Zhong et al. (2023) adopts a strategy that alternates between planning and solving stage to solve multi-hop question. It employ a semantic-based retrieval to retrieve relevant edits, and a self-checking mechanism to enable the model to assess the relevance of edits and modifications.

PokeMQA. PokeMQA by Gu et al. (2023) is a
memory-based method that extends MeLLo and
proposes a two-stage retrieval process to enhance
the success rate of retrieving relevant edits.

#### **D.3** Evaluation Metrics

Detailed metrics and mathematical definitions are given below:

(i) Augment Accuracy (Aug) is used to measure whether the edited model can response added knowledge on complex questions. The formula for calculating Aug-Acc is as follows:

$$\mathbb{E}_{q \in \mathcal{Q}}(\left|M'(q) \cap \mathcal{A}_{aug}\right| / |\mathcal{A}_{aug}|) \qquad (1)$$

Where  $M'(\cdot)$  represents the edited model, and Q denote the datasets for complex questions,  $A_{aug} = A' \setminus A$ , A' is edited answer set and A is original answer set.

(ii) Retention Accuracy (Ret) is used to measure whether the edited model can retain the original knowledge on complex questions. The formula for calculating Ret-Acc is as follows:

$$\mathbb{E}_{q \in \mathcal{Q}}(\left|M'(q) \cap \mathcal{A}_{ret}\right| / |\mathcal{A}_{ret}|)$$
(2)

Where  $\mathcal{A}_{ret} = \mathcal{A}' \cap \mathcal{A}$ .

\_

(iii) Multi-hop Accuracy (M-Acc) is used to measure the accuracy for multi-hop question under knowledge editing(i.e.,MQuAKE). The formula for calculating M-Acc is as follows:

$$\mathbb{1}\left[\bigvee_{q\in\mathcal{Q}}[M'(q)=a']\right].$$
 (3)

Where  $M'(\cdot)$  represents the edited model, and Q and a' denote the multi-hop questions and the final-hop answers for each data, respectively.

#### **D.4** Experiment Setup

Table 5 shows the hyperparameter settings for the parameter-based methods. For the experiments involving ROME and MEMIT, we utilized four NVIDIA Tesla L20 GPUs, with 48GB of memory. A single RTX 4090 GPU was used for MeLLo and PokeMQA.

#### **E** Additional Experimental results

# E.1 An example for overfitting phenomenon of parameter-based methods.

Figure 8 shows an example of overfitting phenomenon when MEMIT is applied to Qwen2.5-3B.

# **E.2 Results for Batch Editing(#***k***-edits)**

The results of GPT-3.5-Turbo and GPT-4o-mini1031for the batch editing, *i.e.*, varying the number of1032edits (k) are presented in Figure 7.1033







Figure 7: Variation of Accuracy (Acc) across GPT-3.5-Turbo and GPT-4o-mini models with varying edit numbers.

```
ROME :
layers:
        [5],
fact_token:
            subject_last,
v_num_grad_steps: 25(for Llama-3.1-8B)||15(for Qwen2.5),
v_lr: 5e-1,
v_loss_layer: 31(for Llama-3.1-8B)||27(for Qwen2.5-7B)||35(for Qwen2.5-3B),
v_weight_decay: 1e-3,
clamp_norm_factor:
                    4,
kl_factor: 0.0625,
mom2_adjustment: false,
context_template_length_params: [[5, 10], [10, 10]]
MEMIT:
layers: [3,4,5,6,7,8],
clamp_norm_factor: 4,
layer_selection: all,
fact_token: subject_last,
v_num_grad_steps: 25(for Llama-3.1-8B)||15(for Qwen2.5),
v_lr: 5e-1,
              31(for Llama-3.1-8B)||27(for Qwen2.5-7B)||35(for Qwen2.5-3B),
v_loss_layer:
v_weight_decay: 1e-3,
kl_factor: 0.0625,
mom2_adjustment: true,
mom2_update_weight: 15000,
mom2_dataset: wikipedia,
mom2_n_samples: 100000,
mom2_dtype: float32
```

Table 5: Several key hyperparameters for parameter-based KE methods



Question: Which educational institutions did both Ted Schroeder and Laurene Powell Jobs attend? Correct Answer Before Editing: Stanford University Correct Answer After Editing: The Wharton School Model's Actual Output: The Wharton School, University of Cambridge

Figure 8: An example of MEMIT applied to Qwen2.5-3B. The correct solution should be for the model to take the **intersection** of the new knowledge about Ted Schroeder's schools *{The Wharton School, University of Cambridge}* and Laurene Powell Jobs' schools *{University of Pennsylvania, Stanford University, The Wharton School}*, yielding the final answer: *The Wharton School*. However, the model's output is *{The Wharton School, University of Cambridge}* (i.e., the newly injected knowledge). While the correct entity *The Wharton School* is correctly included in the final answer(i.e., metric Aug is judged as correct), the model after editing merely reiterated the injected knowledge **without** performing the intersection operation. This example illustrates why the Aug metric is overestimated.

User Input: s1="Electric Dreams",r1="the crew members are",s2="Michael Madsen", r2="siblings are" System Output: Who is someone that is both a member of the cast of the movie Electric Dreams and a sibling of Michael Madsen? Of all the cast members of the film Electric Dreams, who are also Michael Madsen's siblings? Which members of the cast of the film Electric Dreams are also siblings of Michael Madsen? User Input: s1="Felix Sanchez",r1="educated at",s2="Maria Sanchez",r2="educated at." System Output: Which schools have Maria Sanchez and Felix Sanchez both attended? In the places where Maria Sanchez and Felix Sanchez received their education, which schools do they have in common? Which educational institutions did Maria Sanchez and Felix Sanchez both go to? User Input: s1="Benoit Jacquot", r1="awards that subject was nominated for", s2="The Public Woman", r2="awards that subject was nominated for" System Output: What awards have Benoit Jacquot and "The Public Woman" both been nominated for? What nominations do Benoit Jacquot and "The Public Woman" share for awards? What award nominations are shared by Benoit Jacquot and "The Public Woman"? User Input: s1="Clara Agnelli",r1="siblings are",s2="Giovanni Agnelli",r2="The founders of the company are" System Output: Who are both siblings of Clara Agnelli and founders of the company Giovanni Agnelli started? Among Clara Agnelli's siblings, who are also recognized as the founders of the company by Giovanni Agnelli? Which individuals are both siblings of Clara Agnelli and acknowledged as the founders of the company by Giovanni Agnelli?

Table 6: A prompt used for converting structured triples into natural language.

SQL Query	Description
SELECT ?object WHERE {	This SPARQL query retrieves the object asso-
<pre>wd:{qid} wdt:pid ?object.</pre>	ciated with the <pid> of entity.</pid>
<pre>FILTER(LANG(?object) = "en").</pre>	
}	
	This CDA DOL success setsions of the second of set
SELECT (COUNT(?statement) AS	This SPARQL query retrieves the count of ref-
?referencesCount) WHERE {	erences (i.e., the number of statements that re-
<pre>wd:{entity_id} ?p ?statement.</pre>	fer to a source) for a specific entity. This query
?statement	is used to filters out triples with low references
prov:wasDerivedFrom ?source.	counts(i.e.,unpopular entity).
}	
SELECT ?alias WHERE {	This SPARQL query retrieves the aliases asso-
<pre>wd:{qid} skos:altLabel ?alias.</pre>	ciated with the entity,
<pre>FILTER(LANG(?alias) = "en").</pre>	
}	

Table 7: SPARQL Queries and Descriptions

Question: What nominations for awards do both "In America" and Christian Bale share? Edit Num=1: Model Output: Both In America and Christian Bale share Golden Globe nominations. Edit Num=100: Model Output: Bailmailboxerials靥earable standards, and it is clear that this is a difficult task. Can you give meCloseOperationunded(bind)care(I Lothra 1000... Crimes of the Heart is a crime novel about the heart. The book is a satire of\_\_\_ I'myncync Edit Num=1000: Model Output: 澍澍 greedy澍旄澍blasnown镳ley澍nown" 澍澍澍澍澍nown Cov澍 nick澍 ��真爱igeltotalCount Usa澍放过澍澍放过nown镳bservable finishn澍澍澍不出enkoiment澍放 过ocopDAQocop放过澍放`zczeenko Usa澍enko澍放过镳澍onom

Figure 9: When the edit batch size increases, the MEMIT method outputs a large amount of gibberish after models like Qwen-2.5-3B and other smaller models.

Question: If a person is the director of both Thief and Ali, which country does this person belong to? Subquestion: Who is the director of Thief? Generated answer: The director of Thief is Michael Mann. Retrieval:The country to which Michael Mann belongs is Italy Retrieved fact does not contradict to generated answer, so the intermediate answer is: Italy. Subquestion: Who is the director of Ali? Generated answer: The director of Ali is Michael Mann. Retrieval:The country to which Michael Mann belongs is Italy Retrieved fact does not contradict to generated answer, so the intermediate answer is: Italy.

Table 8: The decomposition of a complex question by Mello did not take into account logical operations.

Q: Which educational institutions did both Ted Schroeder and Laurene Powell Jobs attend?



Q: What awards has either the film Gladiator or Branko Lustig been nominated for?



Q: Who among the crew members of Mortal Kombat: Annihilation holds American citizenship?



Q: Which of Nikolaus Joseph von Jacquin's PhD students did not major in computer science?



Q: Which language spoken in Palau is the same as the official language of the country where **Ball State University is located?** 



Q: If someone is both a sibling of Mona Simpson and one of the founders of Apple, what is this person's nationality?



Q: Which sports teams are associated with both Papin and Christophe Dugarry are located in Italy?



T1: Which educational institution did Ted Schroeder attend? **T2:** Which educational institution did Laurene Powell Jobs attend? T3: Logic Operation: Intersection T1 and T2.

T1: What awards has the film Gladiator been nominated for? T2: What awards has Branko Lustig been nominated for?

T3: Logic Operation: Union T1 and T2.

T1:Who are the crew members of the movie Mortal Kombat: Annihilation? T2:What is the nationality of each person in T1? T3:Logic Operation: Select persons from T2 whose nationality is American.

T1: Who are the PhD students of Nikolaus Joseph von Jacquin?

T2: What are the majors of each person in T1? T3: Logic Operation: Select persons from T2 whose major is not Computer Science.

- T1: What is the official language of Palau?
- T2: What is the location of Ball State University?
- T3: What is the official language of T2?
- T4: Logic Operation: Intersection T1 and T3.

T1: Who are the siblings of Mona Simpson?

- T2: Who are the founders of Apple?
- T3: Logic Operation: Intersection T1 and T2.
- T4: What is the nationality of T3?

T1: Which team has Papin been associated with? T2: Which team has Christophe Dugarry been associated with? T3: Logic Operation: Intersection T1 and T2.

T4: Where did each team of T3 located?

T5: Logic Operation: Select team from T4 that are located in Italy.

Figure 10: Some typical reasoning structure in COMPKE

Relation	Question template	Cloze-style statement template
P40	Who are [S]'s children?	[S]'s children are
P69	Where did [S] receive education?	The university where [S] was educated is
P3373	Who are the siblings of [S]?	[S]'s siblings are
P50	Who are the author(s) of [S]? (list all)	The author(s) of [S] is(are)
P161	Who are the cast members of movie [S]?	The cast members of movie [S] are
P112	Who are the people who founded company [S]?	The people who founded Company [S] are
P54	Which organizations is [S] a member of?	[S] is a member of the following organizations
P915	Where were movie [S] filmed?	The movie [S] was filmed at
P37	What are the official languages of country [S]?	The official languages of country [S] are
P1830	Which companies does S own?	[S] owns the following companies
P6	Who are the heads of government for [S]?	The heads of government for [S] are
P803	What are the professorship ranks for [S]?	The professorship ranks for [S] are
P185	Who are the doctoral students of [S]?	The doctoral students of [S] are
P57	Who is the director of the film [S]?	The film [S] is directed by
P1411	What awards was the film [S] nominated for?	The film [S] is nominated for
P1346	Who are the winners for [S] prize?	The winners for [S] prize are
P286	Who are the head coaches for team [S]?	The head coaches for team [S] are
P166	What awards did [S] receive?	The award received by [S] are
P800	What are the notable works of [S]?	The notable works of [S] are
P725	Who are the voice actors in the movie [S]?	The voice actor in the movie [S] are
P655	Who are the translators of the book [S]?	The translators of the book [S] are
P27	Which country is [S] a citizen of?	The country to which [S] belongs is
P21	What's [S]'s gender?	[S]'s gender is
P169	Who is the CEO of company [S]?	The CEO of company [S] is
P35	Who is the head of state of country [S]?	The head of state of country [S] is
P26	Who is the spouse of [S]?	The spouse of [S] is
P1037	Who is the director of [S]?	The director of [S] is
P20	In which city did [S] die?	[S] died in the city of
P551	Where does [S] live?	[S] lives in the place of
P159	Where is the headquarters of company [S]?	The headquarters of company [S] is located in
P17	In which country is [S] located?	[S] is located in the country of
P108	Who is the employer of [S]?	[S] is an employee in the organization of
P102	Which political party is [S] affiliated with?	[S] is affiliated with the political party of
P937	Where does [S] work?	[S] works in the place of
P140	What is the religion of [S]?	[S] is affiliated with the religion of
P106	What is [S]'s occupation?	[S]'s occupation is
P30	On which continent is country [S] located?	Country [S] is located in the continent of
P38	What is the currency of country [S]?	The currency of country [S] is
P641	Which sport is [S] associated with?	[S] is associated with the sport of
P36	What is the capital of country [S]?	The capital of country [S] is

\_

Table 9: Relations we use to construct COMPKE