# MQUAKE-REMASTERED: MULTI-HOP KNOWLEDGE EDITING CAN ONLY BE ADVANCED WITH RELIABLE EVALUATIONS

#### Anonymous authors Paper under double-blind review

#### ABSTRACT

Large language models (LLMs) can give out erroneous answers to factually rooted questions either as a result of undesired training outcomes or simply because the world has moved on after a certain knowledge cutoff date. Under such scenarios, knowledge editing often comes to the rescue by delivering efficient patches for such erroneous answers without significantly altering the rests, where many editing methods have seen reasonable success when the editing targets are simple and direct (e.g., "what club does Lionel Messi currently play for?"). However, knowledge fragments like this are often deeply intertwined in the real world, making effectively propagating the editing effect to non-directly related questions a practical challenge (to entertain an extreme example: "What car did the wife of the owner of the club that Messi currently plays for used to get to school in the 80s?"). Prior arts have coined this task as *multi-hop knowledge editing* with the most popular dataset being MQUAKE, serving as the sole evaluation benchmark for many later proposed editing methods due to the expensive nature of making knowledge editing datasets at scale. In this work, we reveal that up to 33% or 76% of MQUAKE's questions and ground truth labels are, in fact, corrupted in various fashions due to some unintentional clerical or procedural oversights. Our work provides a detailed audit of MQUAKE's error pattern and a comprehensive fix without sacrificing its dataset capacity. Additionally, we benchmarked almost all proposed MQUAKE-evaluated editing methods on our post-fix dataset, MQUAKE-REMASTERED. We observe that many methods try to overfit the original MQUAKE by exploiting some dataset idiosyncrasies of MQUAKE. We provide a guideline on how to approach such datasets faithfully and show that a simple, minimally invasive approach can bring excellent editing performance without such exploitation.

034 035

037

038

005 006

007

008 009 010

011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

#### 1 INTRODUCTION

039 Given the widespread public-facing popularity of various Large Language Model-powered (LLM) 040 products (Zhao et al., 2023; Yang et al., 2024b), even an occasional user has likely experienced 041 LLMs giving out erroneous answers to factually rooted, knowledge-intensive questions. While why 042 LLMs would hallucinate such kind of misinformation is complex and still an open problem noisy training data, model bias, out-of-distribution questions, or even simply because the world 043 has moved on after a certain knowledge cutoff date, all likely contributed their fair share to this 044 rather undesired character of LLMs (Huang et al., 2023; Zhang et al., 2023)— under a practical 045 context, *knowledge editing* is often considered the go-to remedy by delivering efficient patches for such erroneous answers without significantly altering the LLM's output on unrelated queries, 047 nor undergoing another extensive pretraining or finetuning section (Sinitsin et al., 2020; Mitchell 048 et al., 2022).

With the growing need for more credible and trustworthy LLMs, a vast amount of LLM-specific knowledge editing methods have been proposed, and many of them have seen reasonable success in addressing simple and direct editing targets. For example, most modern knowledge editing methods can reliably edit the answer of "What club does Lionel Messi currently play for?" from "Paris Saint-Germain" to "Inter Miami CF" and therefore correctly reflecting the occupation status of Messi (Zhong et al., 2023).

#### 054 1.1 MULTI-HOP KNOWLEDGE EDITING POSES PRACTICAL SIGNIFICANCE AND NON-TRIAL 055 CHALLENGES. 056

However, due to the intertwined nature of different knowledge fragments, a small change in one 057 knowledge fragment can produce ripple-like effects on a vast amount of related questions (Zhong 058 et al., 2023; Cohen et al., 2023). It is often a non-trivial challenge to efficiently propagate the editing effect to non-directly related questions with proper precision and locality. E.g., — as an intentionally 060 extreme case — "What car did the wife of the owner of the club that Messi currently plays for used 061 to get to school in the 80s?" Many knowledge-edited LLMs can still struggle while being fully 062 aware of Messi's abovementioned club transfer (Zhong et al., 2023).

063 Prior arts have realized the practical significance of being able to edit such complex/non-direct 064 questions upon a certain knowledge update, as different knowledge fragments are almost always 065 deeply entangled with each other in the real world (Zhong et al., 2023; Cohen et al., 2023; Wei 066 et al., 2024). Meanwhile, exhausting all potential combinations of questions related to one or a few 067 updated knowledge fragments is impractical. Even if it is feasible, this poses high operational costs 068 and a repeated effort would be required should Messi ever opt to transfer again.

069 Intuitively, a practical knowledge editing method needs to produce correct answers to relevant fac-070 tual questions with only a few updated knowledge fragments available. This task has been coined 071 as multi-hop knowledge editing, with the founding, most popular, and only publicly available re-072 flective dataset to date being MQUAKE by Zhong et al. (2023); serving as the sole evaluation 073 backbone for many proposed modern editing methods due to the expensive nature of making 074 counterfactual and temporal datasets at such a scale (>10,000 cases provided, see Table 7).

075 1.2 UNFORTUNATELY, MQUAKE IS FLAWED DUE TO UNINTENTIONAL CLERICAL AND 076 PROCEDURAL ERRORS — WE FIXED/REMADE IT AND RE-BENCHMARKED ALMOST ALL 077 PROPOSED MULTI-HOP KNOWLEDGE EDITING METHODS.

While MOUAKE is the founding dataset of multi-hop knowledge editing tasks and very much 079 brings life to this vital subject, through a comprehensive audit, we reveal that up to 33% or 76%of MQUAKE questions and ground truth labels are, in fact, corrupted in various fashions due 081 to some unintentional clerical or procedural errors; which inevitably cast doubts on the effec-082 tiveness of developed methods evaluated on MQUAKE. The issues with MQUAKE are significant 083 and growing, especially as MQUAKE becomes a widely adopted dataset in the editing community. 084 Given its importance for building more reliable LLMs — a critical aspect of NLP development -085 we present our work to advance multi-hop knowledge editing with the following contributions:

- A comprehensive audit of MQUAKE: We are the first to present a comprehensive audit of 087 the existing errors within MQUAKE (Zhong et al., 2023), bringing awareness to the knowledge 088 editing community regarding this popular dataset with significant task importance attached.
  - Fix/remake MQUAKE to MQUAKE-Remastered: We present the only available fix/remake that not only patches all discovered errors, and done so without sacrificing the intended intensity and capacity of the original MQUAKE whenever possible.
- 092 • Extensively re-benchmark of almost all existing multi-hop knowledge editing methods: Given the currently existing reports based upon the original MQUAKE are flawed reflections 094 of such proposed methods' capability, we additionally re-benchmark almost all existing multi-hop knowledge editing methods that are available against our MQUAKE-REMASTERED datasets.
- 096 · Present a faithful yet beyond SOTA pilot method for future multi-hop knowledge editing development. We observe that many proposed multi-hop knowledge editing methods intentionally or unintentionally overfit the original MQUAKE dataset by applying data-specific operations that 098 are largely unique to the MQUAKE dataset family. We provide guidance on how to approach 099 these datasets faithfully and additionally show that a simple, minimally invasive method with no 100 such overfitting operations can also achieve excellent editing performance. 101
- 102 103

078

090

091

095

2 PRELIMINARY

- 104
- 105 2.1 BACKGROUND OF MQUAKE 106
- MQUAKE (Multi-hop Question Answering for Knowledge Editing) is a knowledge editing dataset 107 focusing on the abovementioned multi-hop question answering tasks proposed in Zhong et al.

(2023), where every case of MQUAKE is a multi-hop question made by a chain of single-hop subquestions. Specifically, MQUAKE is constructed based on the Wikidata:RDF dataset (Vrandečić & Krötzsch, 2014), which, in its rawest format, is a knowledge graph consisting 15+ trillion of Resource Description Framework (RDF) triples<sup>1</sup>. MQUAKE essentially builds a much more concise subgraph with only 37 manually elected common relations and top 20% of the most common entities, where a walk of {2, 3, 4}-hop on this subgraph can form a case (which is a chain of {2, 3, 4} single-hop subquestions connected together) in the MQUAKE dataset.

115 MQUAKE is presented as two sub-datasets: MQUAKE-CF and MQUAKE-T. The former focuses 116 on counterfactual tasks, while the latter on temporal changes. We highlight that there is also a 117 MQUAKE-CF-3K dataset, a subset of MQUAKE-CF that only contains 3,000 cases out of the 118 original 9171 cases. Authors of MQUAKE evaluate their proposed method, MeLLo (Zhong et al., 2023), upon this MQUAKE-CF-3K dataset; which then become an unspoken standard for the later 119 proposed multi-hop knowledge editing methods (Gu et al., 2024; Shi et al., 2024; Wang et al., 2024; 120 Anonymous, 2024; Cheng et al., 2024). Due to the popularity of this sub-sampled dataset, we 121 provide our error analysis mostly based on MQUAKE-CF-3K and MQUAKE-T in the following 122 §3. For interested readers, we additionally provide the same error analysis upon the full MQUAKE-123 CF in the Appendix B.3. We also collect the dataset statistics in Table 7. 124

125

127

#### 126 2.2 EVALUATING USING MQUAKE

Datasets like MQUAKE-CF and MQUAKE-CF-3K are often tested under varying "editing intensities," based on the number of cases considered "edited." This simulates different levels of deviation between the model's learned knowledge and the newly edited information. This approach is effective because strong knowledge editing methods should handle both large-scale updates and smaller, more localized edits, ensuring that the changes do not interfere with unrelated knowledge.

In its original paper, MQUAKE-CF-3K is evaluated when {1,100,1000,3000} of its 3,000 cases are edited, similarly, MQUAKE-T is evaluated when {1,100,500,1868} of its 1,868 cases being edited, forming an experiment report like Table 6. This kind of report granularity is also adopted by the majority of later proposed multi-hop knowledge editing methods, either in full (Anonymous, 2024) or in spirit with different subsample settings (Gu et al., 2024; Wang et al., 2024; Shi et al., 2024; Cheng et al., 2024; Mengqi et al., 2024). In this work, we report at an even finer level of granularity for maximum cross-reference potentials.



<sup>&</sup>lt;sup>1</sup>https://www.wikidata.org/wiki/Property:P10209

#### 162 AUDITING MQUAKE 3

163 In this section, we present a comprehensive audit of the error pattern that existed in MQUAKE-CF-164 3K and MQUAKE-T (Zhong et al., 2023). We specifically note that our audit is there to provide 165 a better understanding to the knowledge editing community, especially when digesting methods 166 evaluated on these datasets. Our audit is not to discredit the contribution of MOUAKE, or any 167 of the proposed methods evaluated on MOUAKE. We recognize that no dataset can be perfect, 168 especially when it is intrinsically hard to collect large-scale counterfactual and temporal datasets.

169 170

171

#### 3.1 INTRA CONTAMINATION BETWEEN EDITED CASES AND UNEDITED CASES

As discussed in §2.2, having a gradual evaluation coverage from a few to all cases being edited 172 like Table 6 makes sense as an evaluation granularity. However, one critical issue is that  $k \in$ 173  $\{1, 100, 1000, 3000\}$ -edited cases (supposed MQUAKE-CF-3K) are randomly sub-sampled from 174 the 3,000 total cases. Thus, there is no guarantee that the k-edited cases and (3000-k) unedited 175 cases would require two disjoint sets of knowledge and, therefore, risk contamination. 176

For a concrete example, consider the following two multi-hop questions from MQUAKE-CF-3K 177 illustrated in Figure 1. When case 482 is selected as an edited case, the edited fact in case 482 would 178 contaminate the unedited case 126 since both questions would ask for "The citizenship of Kamal 179 Haasan" and the corresponding edited fact would be retrieved. This leads to the model generating 180 an answer in conflict with MQUAKE-CF-3K's label, causing inaccurate experiment readings. See 181 Appendix B.1 for a detailed walk-through. 182

We further note the above-illustrated contamination is not a cherry-picked fluke, but rather a wild-183 spread error. Here, we sample {1,100,1000,2000,3000}-editing targets from MQUAKE-CF-3K 184 using random seed 100, and find the following error statistics in Table 1. 185

186 Table 1: Error statistics of MQUAKE-CF-3K and MQUAKE-T (Zhong et al., 2023) in terms edited 187 cases contaminating unedited cases. k-edited means k cases out of the total dataset are edited. 188

# of Contominated		1	MQUAKE-0	CF-3K			MQ	UAKE-T	
# of Contaminated	1-edit	100-edit	1000-edit	2000-edit	3000-edit	1-edit	100-edit	500-edit	1868-edit
Cases	0	2,013	1,772	910	0	29	1421	1327	0
Subquestions	0	2,706	3,075	1,664	0	29	1421	1327	0

192 193

189 190 191

194 195

196

197

198

It is observable from Table 1 that even a small number of edited cases will cause concerningly large contamination to unedited cases and subquestions, where 67% and 76% of all cases from MQUAKE-CF-3K and MQUAKE-T are contaminated with just 100 cases being edited, introducing a significant distortion to the reported experiment results.

199 We also note that this contamination decreases as the number of edited cases (k-edit) increases, 200 but it's simply a result of fewer unedited cases being available for contamination as k grows. For example, in the extreme case of 3000-edit, there is no contamination between edited and unedited 201 cases because all cases are edited. However, 3000-edit has the highest level of contamination within 202 edited cases, which we explore further in §3.2. 203

204 205

#### 3.2 INNER CONTAMINATION BETWEEN DIFFERENT EDITED CASES

206 Contamination might also happen among multiple edited cases because a certain subquestion pre-207 sented in different edited cases can be edited in some but unedited in others<sup>3</sup> as illustrated in Figure 2. 208 Similar to §3.1, an edited fact from case 1968 would alter the answer to an unedited hop in the edited 209 case 1570. So, the model would generate an answer in conflict with the dataset ground truth label, 210 causing inaccurate experiment readings. See Appendix B.2 for a detailed walk-through.

211 This type of contamination is, once again, universally visible in MQUAKE, as shown in Table 2; 212 which is very much a flipped version of Table 1. With k-edit growing, there are more edited cases, 213

214

<sup>2</sup>We note that in Zhong et al. (2023), "k-edit" means only k of edited cases are evaluated, without any unedited cases. We evaluated both to better reflect the locality of different knowledge editing methods. 215

<sup>&</sup>lt;sup>3</sup>Note, an edited case does not require all of its subquestions being edited, but merely one of it (Table 7)

216 thus more edited-to-edited contamination. Notably, under the 3000-edit tasks, almost one-third 217  $(998/3000, \approx 33\%)$  of the evaluated cases are contaminated, which again introduces distortion to 218 the reported experiment results. We omit the report on MQUAKE-T here because there is only one 219 edit-to-edit contamination when all 1,868 cases from MQUAKE-T are edited (case\_id:424).

Table 2: Error statistics of MQUAKE-CF-3K (Zhong et al., 2023) in terms edited cases contaminating each others. k-edited means k cases out of the total 3,000 cases are edited.



Figure 2: Example of contamination between two edited cases

#### 3.3 **CONFLICTING EDITS**

249 The two types of contamination introduced in §3.1 and §3.2 are indeed subtle and hard to detect. 250 However, MQUAKE-CF-3K also includes some straightforward edit conflicts, such as for the subquestion "Which company is Ford Mustang produced by?" we have the following edits: 252

- case\_id:2566 (edited): Ford Moter Company Nintendo.
- case\_id:231/2707 (edited): Ford Moter Company Fiat S.p.A.

255 This is going to cause a direct conflict when case\_id:2566 and any of the 256 case\_id:231/2707 are both selected as edited cases, as they shall confuse any knowl-257 edge edited LLM for having two answers to the same questions. Fortunately, such types of errors 258 are rather minuscule in MQUAKE-CF-3K, with the abovementioned Ford Mustang question and 259 three cases being the only affected data samples.

261

220

222

224

225

230

231

233

241

242

244 245

246 247

248

251

253

254

260

262

#### MISSING INFORMATION IN MULTI-HOP QUESTION INSTRUCTIONS 3.4

As mentioned in §2, the MQUAKE dataset is built upon a severely filtered Wikidata: RDF knowl-263 edge graph (Vrandečić & Krötzsch, 2014). Specifically, the triples of a certain  $\{2, 3, 4\}$ -hop walk 264 on this subgraph are then fed into a GPT-3.5-turbo model to generate three multi-hop question 265 instructions in a natural language format. During evaluation, an LLM is considered right should it 266 correctly answer against any three of the multi-hop question instructions (Zhong et al., 2023). 267

However, while repeating generation three times definitely reduces the chances of having incompre-268 hensible question instructions, we noticed some of such instructions in MQUAKE are still incom-269 plete. We take the following triple set and its generated 3-questions as an example:

case\_id:546 (unedited): We have a 2-hop question with "Albert Mohler" as the subject and (employer, religion or worldview) as the relation chain. MQUAKE-CF-3K provides the following generated multi-hop questions:

- ♦ Generation #1: What religion is Albert Mohler associated with?
- ♦ Generation #2: Which religion does Albert Mohler follow?
- ♦ Generation #3: With which religious faith does Albert Mohler identify?

All three generated questions omit the part mentioning which company/institution Albert Mohler is
employed by and essentially reduce themselves to single-hop questions, where a correct generation
should read like "*What religion is Albert Mohler's employer associated with*?" Without the complete question, suppose there is an edit on Albert Mohler's employer (which there indeed is one),
the final answer would likely change. However, with this omission of information, even the best
knowledge-edited LLM cannot answer the question correctly with a faithful approach.

As a general analysis, we find the natural language question instructions of 672 cases in MQUAKE-CF-3K are missing information in comparison to their raw triplet chain. This number is counted in the sense that one or more pieces of information present in the triple chain are missing from all three variants of the generated natural language instruction. Similarly, there are 2,830 and 233 cases of erroneous instructions in MQUAKE-CF and MQUAKE-T, respectively.

288 3.5 DUPLICATED CASES

The last kind of error we discovered in MQUAKE is simply unintended duplication — i.e., two or
more cases sharing the same start subjects, edited facts, chain of triples, and final answer — i.e.,
they are the carbon copy of each other, yet simultaneously exist in the dataset. We discovered 47, 4,
and 4 cases of duplication, respectively, in MQUAKE-CF, MQUAKE-CF-3K, and MQUAKE-T.

294

287

289

273

274

275

#### 4 REMASTERING MQUAKE

295 296 297

298

299

300

In this section, we illustrate how we modified and improved the MQUAKE dataset to MQUAKE-REMASTERED with various fixes on the data samples themselves, as well as providing utility modules to facilitate how one interacts with such datasets. We further provide audit correctness analysis in Appendix C. Furthermore, we demonstrate the impact of our improvements through ablation studies that analyze the types of errors addressed, as discussed in Appendix D.

301 302 303

#### 4.1 HARD CORRECTIONS

304 Three types of error existing in MQUAKE can be fixed once and for all with some careful hard 305 corrections, they are namely Conflicting Edits (§3.3), Missing Information in Multi-hop Question 306 Instructions (§3.4), and Duplicated Cases (§3.5). For Conflicting Edits and Duplicated Cases, since 307 there are only a few such errors (<50 per type per dataset), we employ some manual corrections to address these errors: in the former case, we flip the minority edits to align with the majority edits 308 (and adjust their answers to their subsequence subquestions, should there be any); in the latter case, 309 we simply remove such duplicated cases (except for MQUAKE-CF-3K, which we manually select 310 4 more cases from MQUAKE-CF to keep the dataset having 3,000 cases in total and a 1,000 cases 311 for  $\{2,3,4\}$ -hops). For consistency, we rewrite the natural language question instructions for all 312 questions in the datasets using meta-llama/Llama-3.1-405B (Dubey et al., 2024). 313

314

 4.2 DYNAMIC MASKING FOR MAXIMUM COVERAGE: MQUAKE-REMASTERED-CF, MQUAKE-REMASTERED-CF-3K, AND MQUAKE-REMASTERED-T

Due to the contamination count of Intra Edited-to-Unedited Contamination (§3.1) and Inner Editedto-Edited Contamination (§3.2) tend to grow in the opposite direction as shown in Table 1 and 2, it is impossible to find a fix within the current MQUAKE that can address both issues without significantly decreasing the dataset size. As an alternative, we develop an API that will take a case\_id and an edited\_flag as input, indicating the evaluating case-in-question and whether this case is considered edited; our API shall then return a set of triples that are contamination free by dynamically masking out the conflicting edits from other cases. After such, the user may build up an editing knowledge bank upon such triplets and conduct evaluations for any memory-based knowledge editing methods without losing any cases caused by contaminations. Due to the nature of the N-hop question, at most N edited facts would be removed for each case, marginal compared to the number of edited facts in Table 12.

Specifically, once case\_id-of-interest is given, our API would loop through all of its subquestions and identify if any is considered edited under another case. If there is a hit, the triple for such edited subquestions is removed from the bank of edited triples in constant time. This dynamic masking mechanism would ensure all cases within the original MQUAKE be usable against memory-based knowledge editing methods. **However, the drawback of masking is it won't support parameterbased knowledge editing methods**, where weight update is required. We additionally provide a MQUAKE-REMASTERED-CF-6334 to address the need for such methods (Appendix E.1).

334 335

328

329

330

331

332

333

336 337

338 339 340

## 5 MAKING SAFE AND FAITHFUL APPROACH TO MQUAKE AND MQUAKE-REMASTERED

In addition to our dataset audit, fix, and the benchmark results we'd show below, it is our observa-341 tion that many multi-hop knowledge editing methods with decent accuracy reports on MQUAKE 342 or MQUAKE-REMASTERED are utilizing designs that leverage dataset idiosyncrasies unique to 343 MQUAKE. For example, methods like GLAME (Mengqi et al., 2024) utilize Wikidata (Vrandečić 344 & Krötzsch, 2014) as the external knowledge graph to better detect the edit-induced conflicts, which 345 happen to be the source of MQUAKE as discussed in §2.1. While these methods might have decent 346 performance on MQUAKE, the cost of maintaining a positive knowledge graph on the correct — but 347 not just edited — knowledge facts is undoubtedly a non-trivial operation cost. Yet, whether sourc-348 ing the same Wikidata knowledge graph as MQUAKE might bring them data-specific advantages 349 remains unanswered. Similarly, PokeMQA (Gu et al., 2024) utilizes the 6,218 cases included in MQUAKE-CF but not in MQUAKE-CF-3K as the train set to train its auxiliary components. Given 350 MQUAKE is a dataset with relatively low diversity (e.g., it only includes 37 types of relations), 351 whether having a heavily overlapped train and test set will result in data-specific advantages unique 352 to MQUAKE and its variants, again remains unanswered. 353

- 354
- 355 356

357

#### 5.1 A MINIMALLY INVASIVE BUT PERFORMANT APPROACH: GWALK

Here, we provide a brief walkthrough of a simple method we designed, namely <u>GraphWalk</u>. GWalk
 does not leverage any data-specific property unique to MQUAKE or MQUAKE-REMASTERED.
 Yet, it still presents SOTA performance surpassing many, if not all, established baselines. We illus trate this pilot method as concrete guidance and potential inspiration to our future multi-hop
 knowledge editing scholars.

The design of GWalk hinges on the fundamental pipeline of memory-based knowledge editing methods: where the pool of source only contains *edited facts*. This school of editing methods has proven to be successful, mainly because it can leverage the power of retrieval-argument generation (RAG) combined with the in-context learning (ICL) capability of LLMs. Further, it is common sense that edited knowledge facts will be much less than unedited knowledge facts, making maintaining a knowledge pool exclusively containing edited facts a viable option — like done so in MeLLo.

369 Different from MeLLo, where all edited facts are converted from triples to natural language (NL) 370 descriptions in its edited bank, GWalk preserves the edited facts in their original triples fashion and 371 leverages the graph topology. This makes maintaining this edited bank much easier — as one can 372 easily adjust the entity or relation on a knowledge graph without rewriting every natural language 373 description of every related edited fact. It also brings more precise retrieval mapping when a ques-374 tion of a certain edited fact is asked. Methods like MeLLo rely on RAG from a pool of edited facts 375 in NL format. This can lead to unintended retrievals, where irrelevant facts with similar embeddings 376 are retrieved, potentially causing hallucinations. However, if we simply query the entity and relations implied in a question against a knowledge graph, there is less chance of retrieving unintended 377 materials. We share the detailed pseudocode of GWalk in Algorithm 1.

378	Algorithm 1: General Procedure of GWalk on One Multi-hop Question
379	Input:
380	M, the Question Answering Language Model;
381	T, a Text-embedding model;
382	<i>Q</i> , a Multi-hop Question;
202	<i>E</i> , a bank of edited facts as a knowledge graph.
303	Output:
384	op, meanswer to g.
385	i = 1, the subquestion counter:
386	$o_p =$ None, the answer from the previous subquestion.
387	$s \leftarrow Extracted subject from Q;$
2001	$2 rels \leftarrow Prompt M$ to breakdown Q into a sequence of relations.
300	/* If $Q$ is `What is the official language of the country where Karl Alvarez
389	holds citizenship?', then s would be 'Karl Alvarez' and a possible rels is
390	[ CICIZENSIEP, OFFICIAL LANGUAGE ] */
391	4   Ouerv $\langle s, r, ? \rangle$ against E using T, namely we do $T(s)$ first to determine if there is a retrievable $s \in E$ , then
392	inspect if the $s \in E$ has an relation edge retrievable by $T(r)$ .
202	/* We set a threshold on embedding similarity for $T$ to determine whether an
393	item is retrievable or not. */
394	<sup>5</sup> Prompt M to generate subquestion $q_i$ with s and r.
395	6 $o_p \leftarrow$ the <i>M</i> -generated answer to $q_i$ .
396	7 If $T(s, r)$ has a valid reflection $s, r, 0 >$ then
397	
001	/* The answer to this subquestion will be the start subject of the next
398	subquestion. */
399	$\begin{array}{c} 1 \\ 1 \\ 1 \\ 0 \end{array} \stackrel{()}{\leftarrow} i + 1; \end{array}$
400	11 <b>Return</b> $o_p$ ;
401	

#### 6 BENCHMARK AND DISCUSSION

 Given almost all proposed multi-hop knowledge editing methods are evaluated on the original, errorcontained, MQUAKE datasets. Here, we provide a re-benchmark of those methods against post-fix
 MQUAKE-REMASTERED datasets for a more reliable reporting of each method's performance.
 Please refer to Appendix G for reproducibility.

409 6.1 EXPERIMENT COVERAGE

402

403

408

410 **Compared Methods** In this work, we aim to cover most if not all, open-sourced knowledge 411 editing methods specifically evaluated on the original MQUAKE. This includes MeLLo (Zhong 412 et al., 2023), PokeMQA (Gu et al., 2024), RAE (Shi et al., 2024), and DeepEdit (Wang et al., 2024) as methods specifically proposed to target this multi-hop knowledge editing problem and 413 evaluated on MQUAKE. We additionally include ICE (Cohen et al., 2023) and IKE (Zheng et al., 414 2023a) as these are also methods purposed for the (single-edit) multi-hop knowledge editing task, 415 though not specifically evaluated on MOUAKE in their original publications. We note that we 416 are aware methods like GLAME (Mengqi et al., 2024), StableKE (Wei et al., 2024), Temple-MQA 417 (Cheng et al., 2024), and GMeLLo (Anonymous, 2024) are also evaluated on MQUAKE, but they 418 are purposely omitted from our re-benchmark coverage due to lack of open-sourced implementation, 419 likely because most of these works are still in submission. 420

421 Covered Models We opt to use lmsys/vicuna-7b-v1.5 (Zheng et al., 2023b), mistralai/Mistral422 7B-Instruct-v0.2 (Jiang et al., 2023), and meta-llama/Meta-Llama-3-8B-Instruct (AI@Meta, 2024)
423 as the choice of question-answering models, both for alignment with existing works (Zhong et al., 2023; Shi et al., 2024; Gu et al., 2024) as well as providing coverage the most recent language
425 models. For methods that require a text-embedding model as a retriever, we use facebook/contriever426 msmarco (Izacard et al., 2022) for alignment with MeLLo (Zhong et al., 2023).

426 msmarco (Izacard et al., 2022) for alignment with MeLLo (Zhong et al., 2023).
427
428
428
429
429
429
430
430
430
430
431
430
431
432
433
434
434
435
435
436
436
437
438
439
439
430
430
430
430
431
431
432
433
434
434
435
435
436
436
437
438
439
439
430
430
430
430
431
431
432
432
433
434
434
435
436
436
437
438
438
439
439
430
430
430
430
431
431
432
432
433
434
434
434
435
435
436
436
437
438
438
439
439
430
430
430
431
431
431
432
432
433
434
434
434
435
434
436
436
437
438
438
439
439
439
430
430
430
431
431
431
432
432
433
434
434
434
435
434
436
436
436
437
438
438
438
439
439
439
430
430
430
431
431
432
432
433
434
434
434
434
434
435
436
436
436
437
437
438
438

#### 6.2 **RESULTS AND DISCUSSION**

Table 3: Performance Comparison of Original MQUAKE and our MQUAKE-REMASTERED datasets. The original MQUAKE cannot faithfully reflect the true capacities of the methods due to errors specified in §3, especially if the method-in-question is performant. 

Mathad	MQuA	KE-CF-3k	MQuAKE-T		
Wethou	Original	Remastered	Original	Remastered	
MeLLo (Zhong et al., 2023)	6.7	6.77	30.84	44.37	
GWalk (Ours)	36.23	66.33	46.41	54.88	

Given that our MQUAKE-REMASTERED is mostly provided as a fix to MQUAKE, we would like to highlight the drastic results difference when the same method is evaluated on these two datasets. Table 3 shows our fixing can indeed result in drastically different experiment reports. Where such difference is especially significant for stronger methods, suggesting all previous reporting on MQUAKE has room for reliability improvements, which we filled here with MQUAKE-REMASTERED.

Table 4: Experiments on MQUAKE-REMASTERED-CF-6334 with numbers of edited cases and Total Accuracy

methods. Results are reported in the format: (Test Edited Accuracy, Train Edited Accuracy, Unedited Accuracy)

Method	100-edit	MQUAKE-REMA 1000-edit	ASTERED-CF-6334 3000-edit	6344-edit
	vicuna-7b	-v1.5 (Zheng et al., 2023	3b)	
MeLLo (Zhong et al., 2023)	19.16 (0, 10.99, 19.37)	19.27 (5.1, 9.58, 24.53)	11.17 (4.31, 8.55, 23.3)	6.83 (4.58, 7.72, 19.05)
ICE (Cohen et al., 2023) IKE (Zheng et al., 2023a)	OOM OOM	OOM OOM	OOM OOM	OOM OOM
PokeMQA (Gu et al., 2024)	-	-	-	21.77 (3.25, 30.82, 1.59)
DeepEdit Wang et al. (2024)	<1	<1	<1	<1
GWalk (Ours)	<b>57.55</b> (22.22, 64.84, 57.48)	<b>61.79</b> (29.08, 66.17, 63.23)	<b>59.1</b> (39.3, 63.74, 64.33)	<b>56.62</b> (44.64, 62.11, 68.25
	Mistral-7B-In	struct-v0.2 (Jiang et al.,	2023)	
MeLLo (Zhong et al., 2023)	27.5 (<1, 23.08, 27.65)	27.54 (12.76, 24, 30.4)	24.37 (11.88, 25.51, 32.06)	21.26 (13.29, 24.9, 30.16)
ICE (Cohen et al., 2023)	OOM	OOM	OOM	OOM
IKE (Zheng et al., 2023a)	8.82 (11.11,6.59,8.86)	OOM	OOM	OOM
PokeMQA (Gu et al., 2024)	-	-	-	20.38 (3.99, 27.41, 69.84)
DeepEdit Wang et al. (2024)	<1	<1	<1	<1
GWalk (Ours)	<b>56.25</b> (33.33, 57.14, 56.28)	<b>58.9</b> (34.69, 60.57, 60.6)	<b>56.03</b> (42.69, 59.04, 59.85)	<b>54.43</b> (47.49, 57.74, 52.38
	Meta-Llama-3	-8B-Instruct (AI@Meta	, 2024)	
MeLLo (Zhong et al., 2023)	<1	<1	1.12 (1.17, 1.48, 0.22)	1.27 (<1, 1.4, 1.59)
ICE (Cohen et al., 2023)	OOM	OOM	OOM	OOM
IKE (Zheng et al., 2023a)	<1	OOM	OOM	OOM
PokeMQA (Gu et al., 2024)	-	-	-	20.38 (1.08, 28.41, 76.19)
DeepEdit (Wang et al., 2024)	24.13 (11.1, 19.78, 24.29)	24.35 (8.16, 20.52, 26.27)	21.01 (7.57, 19.65, 25.38)	18.90 (7.48, 18.81,28.57)
RAE (Shi et al., 2024)	29.33 (22.22, 12.09, 29.74)	25.65 (33.67, 11.67, 32.49)	15.59 (23.11, 10.12, 33.48)	11.58 (18.75, 11.39, 28.57)
GWalk (Ours)	<b>67.01</b> (33.33, 74.73, 66.92)	<b>71.89</b> (47.45, 80.94, 70.65)	<b>73.76</b> (54.05, 81.6, 71.12)	<b>74.22</b> (61.02, 80.47, 73.02)

In Table 4, we present benchmark results on MQUAKE-REMASTERED-CF-6334. GWalk consis-tently outperforms other methods in terms of models and edit numbers. The "OOM" in ICE and IKE are due to memory overload from concatenating all edited facts in the in-context learning prompt. Whereas, the "<1" results likely stem from the LLM's failure to recognize the few-shot examples, often generating irrelevant tokens or failing to follow the few-shot format. This issue was ob-served with MeLLo using Meta-Llama-3-8B-Instruct, and with DeepEdit using vicuna-7b-v1.5 and Mistral-7B-Instruct-v0.2. Due to page limitation, we refer our readers to Appendix H for bench-marks of MQUAKE-REMASTERED-CF, MQUAKE-REMASTERED-CF-3K, and MQUAKE-**REMASTERED-T.** We present MQUAKE-REMASTERED-CF-6334 in main text solely because it can feature the most methods.

# 486 7 RELATED WORKS

488

489

490 491

504

**Audit and Fix of MQUAKE** To the best of our knowledge, no work has conducted a comprehensive audit to MQUAKE as we do, but two prior arts have touched on the errors existing in MQUAKE: GMeLLo (Anonymous, 2024) and DeepEdit (Wang et al., 2024).

Table 5: Comparison of error analysis/quantification/fix of MQUAKE provided in different works.

Ref.	Error Types Found	Error Quantified	Error Scopes Fixed	Cost of Fixing
GMeLLo (Anonymous, 2024) DeepEdit (Wang et al., 2024)	Missing Instruction	No CF-3K in <b>3000</b> -edit	No CF-3K in <b>3000</b> - edit	N/A 998 out of 300 cases remove from CF-3K
Ours	Intra Contamination, Inner Contamination, Conflicting Edits, Missing Instructions, Duplicated Cases	CF-3K in {1,100,1000,3000}-edit, T in {1,100,500,all}-edit, CF-9K in {1,100,1000,3000, 6000,all}-edit	CF-3K, T, CF-9K in <b>any</b> -edit Remastered- CF-6334 in <b>any</b> -edit	No case re moved from CF-3K, CF-9F or T.

Specifically, GMeLLo briefly discusses the inconsistency between the triple chain and the corresponding generated instructions in its §4.5.1, which is the same type of error we discussed in §3.4.
However, GMeLLo merely presents two examples of such an error without providing any quantitative error analysis or fix; we did both in §3.4 and §4.1.

509 DeepEdit (Wang et al., 2024) discovered the same inner contamination error (edited-to-edited) as we 510 discussed in §3.2, but limited to one dataset (MQUAKE-CF-3K) under one setting (when all 3000 511 cases are considered edited). Further, DeepEdit removed all 998 inner contaminated cases from the 512 MQUAKE-CF-3K dataset — which is (supposedly) the same 998 cases we detect in Table 2 under the 3000-edit column - and named it MQUAKE-2002. While this fix is, of course, helpful, we 513 argue our Remastered fixes are much more comprehensive and effective since they patched 514 many more errors revealed in §3 (the other four types of errors still exist in MQUAKE-2002), 515 and most importantly, done so without scarifying almost 1/3 of the capacity of the original 516 dataset thanks to masking utility we proposed in §4.2. We further demonstrate the quantifiable 517 difference between our work, GMeLLo, and DeepEdit in Table 5. 518

519 Multi-hop Knowledge Editing Datasets RippleEdit (Cohen et al., 2023) is the only other publicly 520 available multi-hop knowledge editing dataset. However, it is actually a single-edit dataset, meaning 521 only one edited fact is addressed at a time. We consider this an oversimplification of real-world 522 scenarios, where systems must handle multiple edits simultaneously. This design also inherently 523 avoids contamination. For additional exercise, we convert RippleEdit to a multi-edit setup to 1) make it more challenging, 2) show that our audit can also "fix" issues within a different dataset, and 524 3) demonstrate our proposed GWalk is indeed faithful and doesn't depend on MQuAKE-specific 525 data. More in Appendix F. 526

- Benchmark and Guidance Our work re-benchmarks nearly all open-sourced knowledge editing methods on MQUAKE and guides on safely and faithfully approaching such datasets. To the best of our knowledge, no other work offers this level of benchmarking or touches on the same issues. Notably, we are likely the only work to evaluate on MQUAKE-CF/MQUAKE-CF-9K, the largest dataset that even the original MQUAKE paper did not assess due to resource constraints. Table 13 illustrates the significant difference in evaluation coverage between our work and previous efforts.
- 533 534

535

## 8 CONCLUSION

In this work, we present a comprehensive audit, a fix of the MQUAKE dataset, and the only avail able evaluation conducted on the error-free datasets: MQUAKE-REMASTERED. We further
 re-benchmarked all open-sourced knowledge editing methods evaluated on MQUAKE with our
 MQUAKE-REMASTERED datasets and provided guidance and examples on how to faithfully approach these datasets with our GWalk.

# 540 REFERENCES

- 542 AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/ llama3/blob/main/MODEL\_CARD.md.
  - Anonymous. Graph memory-based editing for large language models. *Submission to ACL ARR* 2024 February, 2024.
- Keyuan Cheng, Gang Lin, Haoyang Fei, Yuxuan zhai, Lu Yu, Muhammad Asif Ali, Lijie Hu, and Di Wang. Multi-hop question answering under temporal knowledge editing. *arXiv*, 2024.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects
   of knowledge editing in language models. *Transactions of the Association for Computational Linguistics*, 2023.
- 552 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 553 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony 554 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 558 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, 559 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-561 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 562 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 Ma-564 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy 565 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, 566 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-567 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, 568 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, 569 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-570 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, 571 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, 572 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur 573 Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-574 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Oing He, Oingxiao Dong, 575 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, 576 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-577 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, 578 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, 579 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, 581 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, 582 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-583 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, 584 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, 585 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 588 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda 589 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 592 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

594 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, 595 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 596 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana 597 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 598 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 600 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory 601 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, 602 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-603 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, 604 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer 605 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 606 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie 607 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun 608 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, 609 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 610 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, 611 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-612 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 613 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-614 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-615 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, 616 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, 617 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, 618 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, 619 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, 620 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, Sa-621 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-622 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang 623 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 624 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, 625 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, 626 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-627 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, 628 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 629 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-630 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, 631 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 632 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. 633 URL https://arxiv.org/abs/2407.21783. 634

Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ninghao Liu, Ruobing Wang, and Xin Wang. Pokemqa: Programmable knowledge editing for multi-hop question answering. *arXiv*, 2024.

635

636

- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
  Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large
  language models: Principles, taxonomy, challenges, and open questions. *arXiv*, 2023.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. *Transactions on Machine Learning Research*, 2022.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *arXiv*, 2023.

- 648 Tianjie Ju, Yijin Chen, Xinwei Yuan, Zhuosheng Zhang, Wei Du, Yubin Zheng, and Gongshen Liu. 649 Investigating multi-hop factual shortcuts in knowledge editing of large language models, 2024. 650 URL https://arxiv.org/abs/2402.11900. 651 Yi Lu, Jing Nathan Yan, Songlin Yang, Justin T Chiu, Siyu Ren, Fei Yuan, Wenting Zhao, Zhiyong 652 Wu, and Alexander M Rush. A controlled study on long context extension and generalization in 653 llms. arXiv preprint arXiv:2409.12181, 2024. 654 655 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. Advances in Neural Information Processing Systems, 35, 2022. 656 657 Zhang Mengqi, Ye Xiaotian, Liu Qiang, Ren Pengjie, Wu Shu, and Chen Zhumin. Knowledge graph 658 enhanced large language model editing. arXiv, 2024. 659 Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model 660 editing at scale. In International Conference on Learning Representations, 2022. 661 662 Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 663 Retrieval-enhanced knowledge editing for multi-hop question answering in language models. 664 arXiv, 2024. 665 Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitry Pyrkin, Sergei Popov, and Artem Babenko. Editable 666 neural networks. In International Conference on Learning Representations, 2020. 667 668 Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. Communi-669 cations of the ACM, 2014. 670 Yiwei Wang, Muhao Chen, Nanyun Peng, and Kai-Wei Chang. Deepedit: Knowledge editing as 671 decoding with constraints. arXiv, 2024. 672 673 Zihao Wei, Liang Pang, Hanxing Ding, Jingcheng Deng, Huawei Shen, and Xueqi Cheng. Stable knowledge editing in large language models. arXiv, 2024. 674 675 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 676 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, 677 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jin-678 gren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin 679 Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, 680 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 681 Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, 682 Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024a. URL 683 https://arxiv.org/abs/2407.10671. 684 685 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Shaochen 686 Zhong, Bing Yin, and Xia Hu. Harnessing the power of llms in practice: A survey on chatgpt and 687 beyond. ACM Trans. Knowl. Discov. Data, 2024b. 688 Jiayi Yuan, Hongyi Liu, Shaochen Zhong, Yu-Neng Chuang, Songchen Li, Guanchu Wang, Duy Le, 689 Hongye Jin, Vipin Chaudhary, Zhaozhuo Xu, Zirui Liu, and Xia Hu. KV cache compression, but 690 what must we give in return? a comprehensive benchmark of long context capable approaches. 691 In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Findings of the Association 692 for Computational Linguistics: EMNLP 2024, pp. 4623–4648, Miami, Florida, USA, Novem-693 ber 2024. Association for Computational Linguistics. URL https://aclanthology.org/ 694 2024.findings-emnlp.266. Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, 696 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 697 Siren's song in the ai ocean: A survey on hallucination in large language models. arXiv, 2023. 698 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, 699
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Hanyi Tang, Xiaolei Wang, Tupeng Hou, Tingqian Win, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *arXiv*, 2023.

- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can we edit factual knowledge by in-context learning? *arXiv*, 2023a.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
  Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
  Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv*, 2023b.

# Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. MQuAKE: Assessing knowledge editing in language models via multi-hop questions. In *The* 2023 Conference on Empirical Methods in Natural Language Processing, 2023.

7	1	4
7	1	5
7	1	6
7	1	7
7	1	8
7	1	9
7	2	0
7	2	1
7	2	2
7	2	3
7	2	4
7	2	5
7	2	6
7	2	7
7	2	8
7	2	9
7	3	0
7	3	1
7	3	2
7	3	3
7	3	4
7	3	5
7	3	6
7	3	7
7	3	8
7	3	9
7	4	0
7	4	1
7	4	2
7	4	3
7	4	4
7	4	5
7	4	6
7	4	7
7	4	8
7	4	9
7	5	0
7	5	1
7	5	2

753 754 755

## A EXTENDED PRELIMINARY

## A.1 DEMO REPORT OF MQUAKE

Table 6: Standard reporting format of MQUAKE-CF-3K, and MQUAKE-T demoed with MeLLo on Vicuna-7B (Zheng et al., 2023b); k-edited means k cases out of the total cases are edited. Abbreviated table courtesy of Zhong et al. (2023) (Table 3).

Madal	Mathad		MQUA	KE-CF-3K			MQ	UAKE-T	
Widdei	Method	1-edit	100-edit	1000-edit	3000-edit	1-edit	100-edit	500-edit	1868-edit
Vicuna-7B	MeLLo (Zhong et al., 2023)	20.3	11.9	11.0	10.2	84.4	56.3	52.6	51.3

#### A.2 DATASET STATISTICS

Table 7: Dataset Statistics of MQUAKE. Numbers are in terms of cases (a case in MQUAKE is a chain consisting of multiple subquestions).

Dataset	# of 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
	1	2,454	855	446	3,755
	2	2,425	853	467	3,745
MQUAKE-CF	3	-	827	455	1,282
	4	-	-	436	436
	All	4,879	2,535	1,804	9,218
MQUAKE-T	1 (All)	1,421	445	2	1,868

Table 8: Dataset Statistics of MQUAKE-REMASTERED. Numbers are in terms of cases (a case in MQUAKE is a chain consisting of multiple subquestions).

Dataset	# of Edits	2-hop	3-hop	4-hop	Tota
	1	513	356	224	1,093
	2	487	334	246	1,06
MQUAKE-REMASTERED-CF-3K	3	-	310	262	572
	4	-	-	268	268
	All	1,000	1,000	1,000	3,00
	1	2,446	850	441	3,73
MQUAKE-REMASTERED-CF	2	2,415	852	463	3,73
	3	-	823	451	1,27
	4	-	-	430	430
	All	4,861	2,525	1,785	9,17
MQUAKE-REMASTERED-T	1 (All)	1,421	441	2	1,86
	1	1,971	77	0	2,04
	2	2,415	476	14	2,90
MQUAKE-REMASTERED-CF-6334	3	-	823	128	951
	4	-	-	430	430
	All	4,386	1,376	572	6,33

# 810 B EXTENDED AUDITING

# B.1 EXAMPLE OF INTRA CONTAMINATION BETWEEN AN EDITED CASE TO AN UNEDITED CASE (§3.1)

For a concrete example, consider the following two multi-hop questions from MQUAKE-CF-3K (we also additionally provide the subquestion breakdown and intermediate answers of the two questions for better presentation, we note that such auxiliary information is not part of the instruction visible to the question-answering LLM):

- case\_id:126 (unedited): What is the continent of the country where Kamal Haasan holds citizenship?
  - ♦ What is the country of citizenship of Kamal Haasan? India.
  - $\diamond$  What is the continent of India? Asia.
- case\_id:482 (unedited): What is the capital of the country where Kamal Haasan holds citizenship?
  - ♦ What is the country of citizenship of Kamal Haasan? India.
  - ♦ What is capital of India? New Delhi.

The correct pre-edited answer should be "*Asia*" and "*New Delhi*" respectively. As Kamal Haasan is an Indian citizen, India is located in Asia and is the capital of New Delhi. However, suppose case\_id:482 is sampled as an edited case while case\_id:126 remains unedited, we will be provided with the additional triple containing the knowledge of "*The official language of United States of America is Arabic.*"

Since the unedited case\_id:126 and the edited case\_id:482 share the same subquestion of 833 "What is the country of citizenship of Kamal Haasan?" The answer of case\_id:482 will be 834 rightfully updated to "USA" per the new knowledge. However, the unedited case\_id:126 still 835 considers the original answer "India" to be correct, and is therefore contaminated by the edited 836 case case\_id: 482 in an unintended fashion. This is problematic because a successful knowledge 837 editing method should be able to retrieve the edited knowledge — "Kamal Haasan is a citizen of 838 USA?" — upon the relevant questions (in this case the shared one), and thus answering "North 839 America" to case\_id:126. This is technically correct, but in conflict with MQUAKE-CF-3K's 840 label, causing inaccurate experiment readings.

841 842

843 844

845

846 847

848

849

850

851

852

853

854

855

856

812

813

814 815

816

817

818 819

820

821

822

823

824

825

826

#### B.2 EXAMPLE OF INNER CONTAMINATION BETWEEN DIFFERENT EDITED CASES (§3.2)

Again, we walk through two cases from MQUAKE-CF-3K as a concrete example. First, we show them in their unedited format (again, subquestion breakdowns and intermediate answers are here for demonstration purposes and are not visible to the question-answering LLM during evaluation):

- case\_id:1570 (unedited): Who was the creator of the official language used in the work location of Matti Vanhanen?
  - ♦ Which city did Matti Vanhanen work in? Helsinki.
  - ♦ What is the official language of Helsinki? Finnish.
  - ♦ Who was Finnish created by? Mikael Agricola.
- case\_id:1968 (unedited): Who created the official language of Housemarque's headquarters location?
  - ♦ Which city is the headquarter of Housemarque located in? Helsinki.
  - ♦ What is the official language of Helsinki? Finnish.
  - ♦ Who was Finnish created by? Mikael Agricola.

Suppose case\_id:1570 and case\_id:1968 are both selected as editing cases, two triples containing the following knowledge will be available: *"The official language of Helsinki is Black Speech"* (intended for case\_id:1570), and *"Finnish was created by William Shakespeare"* (intended for case\_id:1968), leading to the following edited breakdown.

- case\_id:1570 (edited): Who was the creator of the official language used in the work location of Matti Vanhanen?
  - ◊ Which city did Matti Vanhanen work in? Helsinki.

- ♦ What is the official language of Helsinki? Finnish Black Speech.
  - ♦ Who was Finnish Black Speech created by? J. R. R. Tolkien.

case\_id:1968 (edited): Who created the official language of Housemarque's headquarters location?

- ♦ Which city is the headquarter of Housemarque located in? Helsinki.
- ♦ What is the official language of Helsinki? Finnish.
- ♦ Who was Finnish created by? Mikael Agricola William Shakespeare.

Much like the previous conflict between unedited and edited cases, these two edited cases share a common subquestion: "What is the official language of Helsinki?" However, such subquestion is edited in case\_id:1570 while unedited in case\_id:1968, causing unintended contamination.

B.3 ERROR ANALYSIS OF MQUAKE-CF

Table 9: Error statistics of MQUAKE-CF (Zhong et al., 2023) in terms of edited cases contaminating unedited cases §3.1. *k*-edited means *k* cases are edited out of the total 9218 cases.

# of Contominated			MQUA	KE-CF-3K			
# of Contaminated	1-edit	100-edit	1000-edit	2000-edit	3000-edit	5000-edit	9218-edit
Cases	62	3307	5275	5110	4578	3346	0
Subquestions	62	4525	8751	8989	8326	6364	0

Table 10: Error statistics of MQUAKE-CF (Zhong et al., 2023) in terms edited cases contaminating each others §3.2. *k*-edited means *k* cases are edited out of the total 9218 cases.

# of Contaminated	1-edit	100-edit	1000-edit	2000-edit	3000-edit	5000-edit	9218-edit
Cases	0	8	192	441	732	1397	2873
Subquestions	0	12	270	606	1027	1986	4250

#### 918 C ERROR DETECTION PROCEDURE AND POST-AUDIT CHECKING 919

In this section, we discuss how exactly we carry out our audit and fixes and how we conduct our post-audit checking to ensure our audited datasets are error-free to the best of our ability.

C.1 INTRA AND INNER CONTAMINATION

As discussed in §3.1 and §3.2, we observed that some edited facts were retrieved for subquestions that were not intended to involve an edit. We categorized this issue as contamination, where edited facts inadvertently influence the correct reasoning path. To carry out the audit, we made the following observation: regardless of whether a case is edited or unedited, a valid reasoning path must always exist from the initial subquestion to the last subquestion. Thus, suppose any unedited subquestion on this reasoning path shares the same subject and relation with a triple reflecting an edited fact; then this unedited subquestion is contaminated and therefore flagged.

We programmed the abovementioned filtering mechanism and identified the contaminated edit facts
against different subquestions/cases. We then employed the API described in §4.2 to dynamically
mask out contaminated cases. Last, we confirmed that there is no contamination remaining by reexecuting our filtering program upon the dynamically masked dataset.

937 C.2 CONFLICTING EDITS

As illustrated in §3.3, we noticed some edits within the editing knowledge bank are self-contradicting, where edits with the same subject and relation led to different tail entities. Again, we follow the intended reasoning path as introduced above and check if there are multiple edit-reflecting triples that share the source and relation with an edited subquestion. If so, this suggests there are conflicting edits. We flagged all those edited triples, put ones with shared sources and relations into the same group, then flipped the minority edit to the majority edit and updated their subquestions accordingly. We then reran the program to ensure no more flagged triples.

945 946

947

936

920

921

922 923

924

#### C.3 MULTI-HOP QUESTION INSTRUCTION REWRITE

As highlighted in §3.4, we identified some questions lacked a complete set of relations in their instructions, thus essentially omitting necessary information for a model to provide the correct answer.
We collect a list of synonyms of all relations of an editing path, then evaluate if a certain instruction is not using any of the corresponding synonyms when its reasoning path indicates it should reflect a certain relation.

Subsequently, we prompted the original meta-llama/Llama-3.1-405B to regenerate all instructions with the few-shot demonstration prompt demoed in Appendix E.2 and reran the detection procedure. This process resulted in only a small number of instructions that still didn't meet our predefined rules due to the fact that our lists of synonyms per each relation cannot be exhaustive by design. We then manually inspect and, in a few occupations, manually fix those flagged cases.

958 959 C.4 DUPLICATE CASES

<sup>960</sup> Upon investigating conflicting edits, we accidentally discovered that there exist cases with identical reasoning paths to each other, as illustrated in §3.5. We simply opt to retain only one of such cases and remove the duplicated rests. We then keep track of a set of reasoning paths from all cases and see if the cardinality of the set is equivalent to the number of cases.

- 964
- 965
- 966
- 967 968
- 969
- 970
- 971

# D ERROR TYPE ABLATION STUDY

In this section, we provide ablation studies demonstrating the benefits of addressing errors in the MQuAKE dataset, aligning with the observed error patterns in Tables 2 and 1. Using our proposed GWalk and the Llama-3.1-8B-Instruct model, we evaluate datasets corrected for major error types, including Inner Contamination, Intra Contamination, and Missing Information in Multi-hop Question Instructions. These error types significantly affect the performance of both edited and unedited

accuracies. We opt to exclude minor errors, such as duplicate questions and conflict edits, which are automatically addressed across all settings due to their limited prevalence in the original datasets.

The observed impact is consistent with our analysis: The Inner Contamination fix has the most impact when editing intensity is high (e.g., 3000-edit). Yet, the Intra Contamination fix has the most impact with lower editing intensity (e.g., 100-edit). The Missing Instruction fix consistently improves performance across all editing intensities.

 

 Table 11: Performance comparison across dataset variants of MQUAKE-CF-3K on meta-Total Accuracy

 Ilama/Llama-3.1-8B-Instruct. Results are reported as (Test Edited Accuracy, Unedited Accuracy).

Toma of Famous Firmal	MQUAKE-CF-3K						
Type of Errors Fixed	100-edit	1000-edit	3000-edit				
Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)							
None	45.47	42.73	39.57				
INOILE	(38, 45.72)	(41.3, 43.45)	(39.57, -)				
Innan Contomination	46.83	53.3	71.36				
Inner Contamination	(70, 46.03)	(73.9, 43)	(71.36, -)				
Intro Contomination	71.17	61.73	39.87				
Intra Contamination	(37, 72.35)	(40.9, 72.15)	(39.87, -)				
Missing Instruction	49.33	47.2	45.1				
ivitssing msuluction	(41, 49.62)	(45.6, 48)	(45.1, -)				
All (Our proposed)	76.83	75.03	71.53				
All (Our proposed)	(69, 77.1)	(74.6, 75.25)	(71.53, -)				

## 1004 E EXTENDED REMASTERING

<sup>06</sup> E.1 Contamination Free Subset: MQUAKE-Remastered-cf-6334

While MQUAKE-REMASTERED-MASKED with masking operation can well support memorybased knowledge editing methods, it will not be compatible with parameter-based methods. This is
because, for parameter-based methods, the set of edited facts used for training and evaluation needs
to be constant yet consistent with each other at all times; whereas dynamic masking cannot suffice
as it is essentially adjusting the dataset on the fly during inference time.

To effectively evaluate parameter-based knowledge editing methods, we present MQUAKE-REMASTERED-CF-6334. MQUAKE-REMASTERED-CF-6334 is a dataset extracted from MQUAKE-CF, where all 6,334 cases are edited cases; and they are completely contamination-free from each other. This dataset is suitable for LLM editing with parameter-based approaches, as one can make careful splits among the 6,334 cases of MQUAKE-REMASTERED-CF-6334 to serve as train, validation, and evaluation sets.

1020Table 12: The number of unique edited facts for a varied number of edited cases in MQUAKE-1021REMASTERED-CF

1022						
1023	Number of Edited Cases	100	1000	3000	6000	All (9171)
1024	Number of Unique Edited Facts	150	1171	2991	5137	7252
1005						

1026 Table 13: Experiment coverage comparison among our and other works. For brevity and better rele-1027 vance, "Method Coverage" only includes open-sourced methods specifically designed for multi-hop editing, as adopted single-hop editors are often too weak to deliver usable results. "Separate Metrics?" means that 1028 both the accuracy of edited cases and unedited cases are reported. We consider the inclusion of both metrics 1029 paramount, as editing is often a double-edged sword, causing potential hallucinations under unedited scenarios. 1030 Prior work often only tests on the former but ignores the latter. We did both in our work. 1031

Ref.	Dataset Coverage	Method Coverage	Separate Metrics?	Error Fix?	
MQuAKE (Zhong et al., 2023)	CF-3K {1, 100, 1000, all}-edit; T {1, 100, 500, all}-edit	MeLLo	No	No	
Temple-MQA (Cheng et al., 2024)	CF-3K {1, 100, all}-edit; T {1, all}- edit	MeLLo, PokeMQA	No	No	
Ju et al. (2024) PoleMQA (Gu et al., 2024)	CF-3K {all}-edit CF-3K {1, 100, all}-edit; T {1, 100, all}-edit	N/A MeLLo, PokeMQA	No No	No No	
Ours	CF-3K {1, 100, 1000, all}-edit; T {1, 100, 500, all}-edit; CF-9K {1, 1000, 3000, 6000, all}-edit; CF- 6334 {100, 1000, 3000, all}-edit	MeLLo, ICE, IKE, PokeMQA, GWalk, RAE, DeepEdit	Yes	Yes	

#### **E.2** PROMPT FOR REWRITING INSTRUCTIONS

#### Few-shot Prompt

**Instruction:** Given a chain of relations, generate 3 multi-hop questions that comprehensively include the semantics of the relations.

#### Example 1:

```
Relation Chain:
XXX \rightarrow 'The author of is' \rightarrow ' is a citizen of' \rightarrow ?
Generated Questions:
```

- 1. What is the country of citizenship of the author of XXX?
- 2. What country is the author of XXX a citizen of?
- 3. What is the nationality of the author of XXX?

#### **Example 2:**

```
Relation Chain:
```

```
XXX -> ' was developed by' -> 'The chairperson of is' -> '
is a citizen of ' -> ' is located in the continent of ' -> ?
Generated Questions:
```

- 1. What continent is the country located in, where the chairperson of the developer of XXX is a citizen?
- 2. On which continent is the country located, whose citizen is the chairperson of the company that developed XXX?
- 3. Which continent houses the country of the chairperson of the developer of XXX?

## Example 3:

```
Relation Chain:
<The relational chain we want the generated questions to be
based on>
```

#### 1074 1075

1046

1047 1048

1049

1050

1051 1052

1053

1054

1055 1056

1057

1058

1059

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1077

1078

#### F RIPPLEEDIT

We consider MQuAKE's task design and setup to be more reflective of real-world editing tasks, 1079 as naturally, there will always be more than one edited fact stored for any system with reasonable complexity. That being said, we are happy to report our proposed pilot method, GWalk, performs decently on RippleEdit. Here are some snapshot results on Llama-2-7b-chat:

Table 14: Single-edit result of RippleEdit-Popular/Recent/Random. C1/2 means the edit is happening at the 1st or the 2nd hop (RippleEdit cases only have 2 hops).

Method	Popular C1 Acc.	Popular C2 Acc.	Recent C1 Acc.	Recent C2 Acc.	Random C1 Acc.	Random C2 Acc.
ROME	37.4	16.2	47.8	50.0	35.5	49.5
ICE	85.1	67.6	74.8	85.0	73.8	80.3
MeLLo	45.1	77.1	50.2	80.0	40.2	68.3
GWalk (ou	rs) 85.7	81.8	80.9	87.6	76.1	82.9

We additionally convert RippleEdit to a multi-edit setup — i.e., there are multiple edited facts within the editing knowledge bank at the same time — to a) make it more challenging and, b) show that our audit can also "fix" issues within a different dataset. Note we put the "fix" in quotes as RippleEdit is not designed with multi-edit in mind, so the things we fixed are not necessarily errors but just some adjustments required for making a proper multi-edit dataset. In any case, here are the snapshot results on Llama-3-8b-Instruct:

1099Table 15: Multi-edit result of RippleEdit-Popular/Recent/Random. In this case, we fixed 21/2/01100conflict edits and 3/120/1 case-to-case contamination within RippleEdit-Ropular/Recent/Random1101datasets, respectively.

Method	Popular C1 Acc.	Popular C2 Acc.	Recent C1 Acc.	Recent C2 Acc.	Random C1 Acc.	Random C2 Acc.
MeLLo	35.1	40.3	41.1	42.4	49.5	50.0
GWalk (ours)	79.0	66.9	79.2	63.9	72.9	60.0

# G REPRODUCIBILITY

All experiments are conducted with one or more 80G NVIDIA A100 GPUs from the DGX A100 server. Please refer to https://anonymous.4open.science/r/MQuAKE-Remastered-118E for assets.

#### 1134 Η ADDITIONAL EXPERIMENT RESULTS

1135

#### 1136 One observation we made in §6.2 is in-context learning-based methods — like ICE (Cohen et al., 1137 2023) and IKE (Zheng et al., 2023a) tend to "OOM" when facing a larger amount of edited facts. 1138 This is because these two methods — originally designed for single-edit tasks — essentially dump

1139 all edited facts as a long concatenated prompt and expect the model to figure out the corresponding 1140 editings naturally. They face OOM issues because when the number of editing facts grows, the prompt becomes extremely long and, therefore, introduces a large amount of KV cache and poses 1141 significant memory footprint issues. 1142

1143 While efficiently and effectively handling long input is out-of-scope of our work, as general guid-1144 ance, we refer interested readers to efficient long context-handing survey/benchmark works like 1145 Yuan et al. (2024), which cover the schools and performance of several popular long context-1146 handling methods. Other than the system challenges, another necessary aspect is to improve LLM long context performance, as most LLMs are pretrained on limited context length and thus cannot 1147 effectively handle long input even if the system challenge is addressed. In this regard, we again 1148 recommend survey/benchmark works like Lu et al. (2024) for insights. Further, one can certainly 1149 convert this long context scenario to leverage the power of the RAG pipeline, much like the majority 1150 of multi-hop knowledge editing methods featured in this work. 1151

- 1152
- 1153

1154

Table 16: This is the benchmark result for MQUAKE-REMASTERED-CF-3K reported in the Total Accuracy

format of: (Edited Accuracy, Unedited Accuracy)

$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	1155		,		27					
Method         1-edit         100-edit         1000-edit         3000-edit           1155         vicuna-7b-v1.5 (Zheng et al., 2023b) $(16.54)$ 18 $14.63$ $6.77$ 1160         ICE (Cohen et al., 2023) $(16.54)$ 18 $14.63$ $6.77$ 1161         IKE (Zheng et al., 2023a) $<1$ $<1$ OOM         OOM           1162         DeepEdit Wang et al. (2024) $<1$ $<1$ $<1$ $<1$ $<1$ 1163         GWalk (Ours) $54.89$ $60.9$ $57.37$ $66.33$ N/A)           1166         MeLLo (Zhong et al., 2023) $19.73$ $18.6$ $16.33$ $15.93$ 1166         MeLLo (Zhong et al., 2023) $11.73$ $18.6$ $16.33$ $15.93$ 1167         ICE (Cohen et al., 2023) $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$	1156	Mathad	M	QUAKE-REN	MASTERED-CI	F-3K				
vicuna-7b-v1.5 (Zheng et al., 2023b)         MeLLo (Zhong et al., 2023)       16.54       18       14.63       6.77         MeLLo (Zhong et al., 2023)       (100, 16.51)       (9.0, 18.31)       (8.0, 17.95)       (6.77, N/A)         MeLLo (Zhong et al., 2023)                 Mistal       General (2024)                    Mistal-7B-Instruct-v0.2 (Jiang et al., 2023)       MeLLo (Zhong et al., 2023)       19.73       18.6       16.33       15.93       N/A)         Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)       19.73       18.50       116.33       15.93       N/A)         MeLLo (Zhong et al., 2023)       19.73       18.50       16.33       15.93       N/A)         MelLo (Zhong et al., 2023)        1        4.43       OOM       OOM       OOM         Methol (Murs)       Methol (3.55)       (47, 62.45)       (51.5, 60.0)       (51.0, N/A)       (51.0, N/A)         Mittal       Methol (2.00001       1       (2.1       1.03       2.3       (3.0, <1)       (2.3, N/A)       (3.0, <1)       (2.3, N/A)	1157	Method	1-edit	100-edit	1000-edit	3000-edit				
1159 1160 1161 1161 1162 1163 1163 1164 1163 1164 1164 1164 1165 1166	1158	vicun	a-7b-v1.5 (Zh	eng et al., 202	23b)					
1160ICE (Cohen et al., 2023) IKE (Zheng et al., 2023a) DeepEdit Wang et al. (2024) GWalk (Ours)(100, 16.51)(90, 18.31)(8.0, 17.95)(6.77, NA) OOM OOM OOM OOM OOM (41 $<1$ $<0$ 1163 1164GWalk (Ours) $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ <td>1159</td> <td>MeLLo (Zhong et al., 2023)</td> <td>16.54</td> <td>18</td> <td>14.63</td> <td>6.77</td>	1159	MeLLo (Zhong et al., 2023)	16.54	18	14.63	6.77				
1161ICE (Conclust et al., 2023) DeepEdit Wang et al. (2024) GWalk (Ours) $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$ $<<1$	1160	ICE (Cohen et al. $2023$ )	(100, 16.51)	(9.0, 18.31)	(8.0, 17.95)	(6.77, N/A)				
1162Interpletit Wang et al. (2024) DeepEdit Wang et al. (2024) $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ 1163 1164GWalk (Ours)Mistral-7B-Instruct-v0.2 (Jiang et al., 2023) $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ 1165MetLo (Zhong et al., 2023)19.7318.616.3315.9315.931166MetLo (Zhong et al., 2023)19.7318.616.3315.931167ICE (Cohen et al., 2023)19.71 $< 21$ , 18.52 $(17.8, 15.6)$ $(15.93, N/A)$ 1169IKE (Zheng et al., 2023a) $< 1$ $< 4$ .43OOMOOM1170DeepEdit Wang et al. (2024) $< 1$ $< 1$ $< 1$ $< 1$ 1171GWalk (Ours) $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ 1172Meta-Llama-3-8B-Instruct (AI@Meta, 2024) $< 1$ $< 1$ $< 1$ $< 1$ 1173MetLo (Zhong et al., 2023) $< 1$ $< 1$ $< 1$ $< 0$ $< 0$ 1174MetLo (Zhong et al., 2023) $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ 1174MetLo (Zhong et al., 2023) $< 1$ $< 1$ $< 1$ $< 0$ $< 0$ $< 0$ 1175ICE (Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $0$ $0$ $< 0$ 1176IKE (Zheng et al., 2023a) $< 1$ $< 1$ $< 1$ $0$ $< 0$ $< 0$ 1177DeepEdit Wang et al. (2024) $< 1$ $< 1$ $< 0$ $0$ $< 0$ $< 0$	1161	IKE (Concil et al., $2023$ )		OOM	OOM	OOM				
54.89 (100, 54.87)60.9 (54, 61.14)57.37 (54, 45, 58.85)66.33 (66.33, N/A)1164 1165Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)1166 1166Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)1167 1168ICE (Cohen et al., 2023)1169 1169ICE (Cohen et al., 2023)1160 1170ICE (Cohen et al., 2023)1169 1170ICE (Cohen et al., 2023)1171 1171GWalk (Ours)1172 11721173 11741174 1175 11751175 1176 1176 11761176 1177 1176 11761177 1178 11761178 1179 11761179 1176 1186 11881174 1176 11801175 1176 11801176 1180 11801177 1178 11801178 1180 11801179 11801180 1181 11801181 1181 11811181 11821181 11831182 1184 11841184 11851184 11851184 1186	1162	DeepEdit Wang et al. (2024)	<1	<1	<1	<1				
(100, 94.37)(34, 61.14)(34, 4, 36.35)(06.35, 10/A)1165Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)1166MeLLo (Zhong et al., 2023)1167ICE (Cohen et al., 2023)1168IKE (Zheng et al., 2023a)1170DeepEdit Wang et al. (2024)1171GWalk (Ours)1172MetLo (Zhong et al., 2023)1174MetLo (Zhong et al., 2023)1175ICE (Cohen et al., 2023)1176ICE (Cohen et al., 2023)1177MetLo (Zhong et al., 2023)1178MetLo (Zhong et al., 2023)1174ICE (Cohen et al., 2023)1175ICE (Cohen et al., 2023)1176IKE (Zheng et al., 2023)1177ICE (Cohen et al., 2023)1178GWalk (Ours)1179GWalk (Ours)1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours)1182GWalk (Ours)1183Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024)1184GWalk (Ours)11851186	1163	GWalk (Ours)	<b>54.89</b>	<b>60.9</b>	<b>57.37</b>	66.33				
Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)MeLLo (Zhong et al., 2023)19.7318.616.3315.93MeLLo (Zhong et al., 2023)19.7318.616.3315.93Mice(Ce (Cohen et al., 2023) $< 1$ $< 1$ OOMOOMMike(Zheng et al., 2023a) $< 1$ $< 4.43$ OOMOOMDeepEdit Wang et al. (2024) $< 1$ $< 1$ $< 1$ $< 1$ $< 1$ GWalk (Ours)Meta-Llama-3-8B-Instruct (AI@Meta, 2024)Meta-Llama-3-8B-Instruct (AI@Meta, 2024)Meta-Llama-3-8B-Instruct (AI@Meta, 2024) $< 1$ $< 1$ $< 2.3$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $< 3$ Meta-Llama-3-8B-Instruct (AI@Meta, 2024) $< 1$ $< 2.3$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 0$ $< 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 0$ $< 2.3$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $< 0$ $< 2.3$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $< 0$ $> 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $> 0$ $> 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $> 0$ $> 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $> 0$ $> 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $> 0$ $> 0$ Meta-Cohen et al., 2023) $< 1$ $< 1$ $< 1$ $> 0$ $> 0$ Meta-Cohen et al., 2024) <td>1164</td> <td></td> <td>(100, 34.87)</td> <td>(34, 01.14)</td> <td>(34.4, 38.83)</td> <td>(00.55, IVA)</td>	1164		(100, 34.87)	(34, 01.14)	(34.4, 38.83)	(00.55, IVA)				
1166 1167 1168 1168 1169MeLLo (Zhong et al., 2023) ICE (Cohen et al., 2023) IKE (Zheng et al., 2023a)19.73 19.73 (100, 19.71)18.6 (21, 18.52) (17.8, 15.6) (17.8, 15.6) (15.93, N/A) (17.8, 15.6) (15.93, N/A) (15.93, N/A) (15.93, N/A) (100, 19.71)1169 1169 1170 1170 1170 1171 1171 1172 1172 1172 1173 1174 1175 1174 1175 1175 1176 1176 1176 1177 1176 1177 1178 1176 1176 1176 1177 1176 1177 1178 1176 1176 1176 1177 1177 1178 1176 1176 1177 1178 1176 1177 1178 1176 1177 1178 1179 1180 1181 1180 1181 1181 1181 1181 118219.73 1182 1182 1181 1182 1181 1182 1183 1184 1184 1184 118419.73 1184 1184 1184 118419.73 1184 1184 118519.73 1184 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118419.73 1184 1184 118419.73 1184 1184 118419.73 1184 1184 118419.73 1184 1184 118419.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118419.73 1184 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 1184 1184 118419.73 1184 1184 118619.73 1184 1184 118619.73 1184 1184 118619.73 1184 118619.73 1187 1181 1181 1184 118619.73 1187 11	1165	Mistral-7	B-Instruct-v0	.2 (Jiang et al.	., 2023)					
1167 1168 1169ICE (Cohen et al., 2023) IKE (Zheng et al., 2023a) $<1$ $<1$ $<1$ $OOM$ $OOM$ 1170 1170DeepEdit Wang et al. (2024) GWalk (Ours) $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$	1166	MeLLo (Zhong et al., 2023)	19.73	18.6	16.33	15.93 (15.93 N/A)				
1168INC (Zheng et al., 2023a)<1 $4.43$ ( $4,4.49$ )OOMOOM1170DeepEdit Wang et al. (2024)<1	1167	ICE (Cohen et al., 2023)	<1	<1	OOM	OOM				
1170 1171 1172DeepEdit Wang et al. (2024) GWalk (Ours) $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ $<1$ 1172 1173Meta-Llama-3-8B-Instruct (AI@Meta, 2024)Meta-Llama-3-8B-Instruct (AI@Meta, 2024)Meta-Llama-3-8B-Instruct (AI@Meta, 2024)1173 1174MeLLo (Zhong et al., 2023) ICE (Cohen et al., 2023) $<1$ $<1$ $1.03$ $2.3$ 1175 1176 1176 1177 1177 1178 1178 1179Meta-Llama et al., 2023) IKE (Zheng et al., 2023a) $<1$ $<1$ $0OM$ $0OM$ 1176 1177 1178 1178 1178 1179GWalk(Ours) $<1$ $<1$ $0OM$ $0OM$ 1179 1180 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024) $(10, 68.99)$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1181 1182 1183 1184 1185GWalk (Ours) <b>73.3</b> (100, 73.3) <b>76.83</b> (69, 77.1) <b>75.03</b> (74.6, 75.25) <b>71.53</b> (71.53, N/A)1184 1185GWalk (Ours) <b>65.33</b> (100, 65.35) <b>65.27</b> (65, 65.28) <b>65.07</b> (68.4, 63.4) <b>66.74</b> (66.74, N/A)	1168	IKE (Zheng et al., 2023a)	<1	4.43	OOM	OOM				
1170DeepLant wing et al. (2024) (100, 56.55)56.761.9357.1751.01171GWalk (Ours)Meta-Llama-3-8B-Instruct (AI@Meta, 2024)1173MetLo (Zhong et al., 2023) (E (Cohen et al., 2023))<1	1109	DeenEdit Wang et al. (2024)	<1	(4,4.49)	<1	<1				
1171GWalk (Ours) $(100, 56.55)$ $(47, 62.45)$ $(51.5, 60.0)$ $(51.0, N/A)$ 1172Meta-Llama-3-8B-Instruct (AI@Meta, 2024)1173MeLLo (Zhong et al., 2023) $<1$ $<1$ $1.03$ $2.3$ 1174MeLLo (Zhong et al., 2023) $<1$ $<1$ $0.00M$ $0.00M$ 1175ICE (Cohen et al., 2023) $<1$ $<1$ $0.00M$ $0.00M$ 1176IKE (Zheng et al., 2023a) $<1$ $<1$ $0.00M$ $0.00M$ 1177DeepEdit Wang et al. (2024) $22.93$ $17.27$ $15.03$ $12.63$ 1178GWalk(Ours) $0.00, 68.99$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours) $73.3$ $76.83$ $75.03$ $71.53$ 1182GWalk (Ours) $65.33$ $65.27$ $65.07$ $66.74$ 1184GWalk (Ours) $65.33$ $65.28$ $(68.4, 63.4)$ $(66.74, N/A)$ 1186	1170		56.57	61.93	57.17	51.0				
1172 1173Meta-Llama-3-8B-Instruct (AI@Meta, 2024)1173MeLLo (Zhong et al., 2023)<1	1171	Gwalk (Ours)	(100, 56.55)	(47, 62.45)	(51.5, 60.0)	(51.0, N/A)				
1173 1174MeLLo (Zhong et al., 2023) ICE (Cohen et al., 2023)<1<1 $1.03$ $2.3$ 1175 1176ICE (Cohen et al., 2023) IKE (Zheng et al., 2023a)<1	1172	Meta-Lla	Meta-Llama-3-8B-Instruct (AI@Meta, 2024)							
1174McLEO (Zhong et al., 2023)C1 $(2.0, <1)$ $(3.0, <1)$ $(2.3, N/A)$ 1175ICE (Cohen et al., 2023) $<1$ $<1$ OOMOOM1176IKE (Zheng et al., 2023a) $<1$ $<1$ OOMOOM1177DeepEdit Wang et al. (2024) $22.93$ $17.27$ $15.03$ $12.63$ 1178GWalk(Ours) <b>69.076.7375.4770.6</b> 1179GWalk(Ours) $(100, 68.99)$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours) <b>73.376.8375.0371.53</b> 1182GWalk (Ours) <b>65.3365.2765.0766.74</b> 1184GWalk (Ours) <b>65.3365.2765.0766.74</b> 1185GWalk (Ours) <b>163.365.2765.0766.74</b> 11861186 $(100, 65.35)$ $(65, 65.28)$ $(68.4, 63.4)$ $(66.74, N/A)$	1173	MeLLo (Zhong et al. 2023)		<1	1.03	2.3				
1175ICE (Cohen et al., 2023)<1<1OOMOOM1176IKE (Zheng et al., 2023a)<1	1174			(2.0, <1)	(3.0, <1)	(2.3, N/A)				
1176IKE (Zheng et al., 2023a) $< 1$ $< 1$ $< 1$ $< 00M$ $00M$ 1177DeepEdit Wang et al. (2024) $22.93$ $17.27$ $15.03$ $12.63$ $(0, 22.94)$ $(11, 17.48)$ $(15.1, 15.0)$ $(12.63, N/A)$ $GWalk(Ours)$ $69.0$ $76.73$ $75.47$ $70.6$ $(100, 68.99)$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ $1180$ Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024) $1181$ GWalk (Ours) $73.3$ $76.83$ $75.03$ $71.53$ $(100, 73.3)$ $(69, 77.1)$ $(74.6, 75.25)$ $(71.53, N/A)$ $1183$ Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a) $1184$ GWalk (Ours) $65.33$ $65.27$ $65.07$ $66.74$ $1186$ GWalk (Ours) $65.33$ $65.28$ $(68.4, 63.4)$ $(66.74, N/A)$	1175	ICE (Cohen et al., 2023) IKE ( $7h$ and $a$ at al. 2022a)	<1	<1	OOM	OOM				
1177DeepEdit Wang et al. (2024) $11727$ $11727$ $11503$ $112.03$ 1178GWalk(Ours) $(0, 22.94)$ $(11, 17.48)$ $(15.1, 15.0)$ $(12.63, N/A)$ 1179GWalk(Ours) $(100, 68.99)$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours) $73.3$ $76.83$ $75.03$ $71.53$ 1182GWalk (Ours) $(100, 73.3)$ $(69, 77.1)$ $(74.6, 75.25)$ $(71.53, N/A)$ 1183Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a)1184GWalk (Ours) $65.33$ $65.27$ $65.07$ $66.74$ 11851186 $(100, 65.35)$ $(65, 65.28)$ $(68.4, 63.4)$ $(66.74, N/A)$	1176	IKE (Zheng et al., 2023a)	22.03	<1 17.27	15.03	12.63				
69.076.7375.4770.61179GWalk(Ours) $(100, 68.99)$ $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours) <b>73.376.8375.0371.53</b> 1182GWalk (Ours) $(100, 73.3)$ $(69, 77.1)$ $(74.6, 75.25)$ $(71.53, N/A)$ 1183Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a)1184GWalk (Ours) <b>65.3365.2765.0766.74</b> 1185GWalk (Ours)100, 65.35) $(65, 65.28)$ $(68.4, 63.4)$ $(66.74, N/A)$	1177	DeepEdit Wang et al. (2024)	(0, 22.94)	(11, 17,48)	(15.1, 15.0)	(12.63, N/A)				
1179(Twark(Ours))(100, 68.99) $(67, 77.07)$ $(74.2, 76.1)$ $(70.6, N/A)$ 1180Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)1181GWalk (Ours) <b>73.376.8375.0371.53</b> 1182(100, 73.3)(69, 77.1)(74.6, 75.25)(71.53, N/A)1183Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a)1184GWalk (Ours) <b>65.3365.2765.0766.74</b> 1185(100, 65.35)(65, 65.28)(68.4, 63.4)(66.74, N/A)	1178	CWalls(Ours)	69.0	76.73	75.47	70.6				
Meta-Llama/Llama-3.1-8B-Instruct (Dubey et al., 2024)         I181       GWalk (Ours)       73.3       76.83       75.03       71.53         1182       (100, 73.3)       (69, 77.1)       (74.6, 75.25)       (71.53, N/A)         1183       Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a)         1184       GWalk (Ours)       65.33       65.27       65.07       66.74         1185       (100, 65.35)       (65, 65.28)       (68.4, 63.4)       (66.74, N/A)	1179	Gwark(Ours)	(100, 68.99)	(67, 77.07)	(74.2, 76.1)	(70.6, N/A)				
1181 1182 1182 1183GWalk (Ours)73.3 $(100, 73.3)$ 76.83 $(69, 77.1)$ 77.03 $(74.6, 75.25)$ 71.53 	1180	Meta-Llama/Ll	ama-3.1-8B-I	nstruct (Dube	y et al., 2024)					
1182       Gwaik (Ours)       (100, 73.3)       (69, 77.1)       (74.6, 75.25)       (71.53, N/A)         1183       Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a)         1184       GWalk (Ours) <b>65.33 65.27 65.07 66.74</b> 1185       (100, 65.35)       (65, 65.28)       (68.4, 63.4)       (66.74, N/A)         1186	1181	CWalls (Ours)	73.3	76.83	75.03	71.53				
GWalk (Ours)         65.33         65.27         65.07         66.74           (100, 65.35)         (65, 65.28)         (68.4, 63.4)         (66.74, N/A)	1182	Owark (Ours)	(100, 73.3)	(69, 77.1)	(74.6, 75.25)	(71.53, N/A)				
I184         GWalk (Ours)         65.33         65.27         65.07         66.74           1185         (100, 65.35)         (65, 65.28)         (68.4, 63.4)         (66.74, N/A)           1186	1183	Qwen/Qwe	en2.5-7B-Instr	uct (Yang et a	l., 2024a)					
1185 (100, 65.35) (65, 65.28) (68.4, 63.4) (66.74, N/A) 1186	1184	GWalk (Ours)	65.33	65.27	65.07	66.74				
1186	1185		(100, 65.35)	(65, 65.28)	(68.4, 63.4)	(66.74, N/A)				
	1186									

# Table 17: This is the benchmark results of MQUAKE-REMASTERED-T. The reported format is: Total Accuracy (Edited Accuracy, Unedited Accuracy)

Mathed MQUAKE-REMASTERED-T							
Method	1-edit	100-edit	500-edit	1864-edit			
vicu	na-7b-v1.5 (Zh	eng et al., 202	3b)				
MeLLo (Zhong et al., 2023)	19.31	18.88 (45.0, 17.4)	22.16 (40.4, 15.47)	44.37 (44.37 N/A)			
ICE (Cohen et al., 2023)	<1	<1	<1	OOM			
IKE (Zheng et al., 2023a)	<1	<1	<1	OOM			
DeepEdit Wang et al. (2024)	<1	<1	<1	<1			
	35.52	46.51	48.93	54.88			
Gwalk (Ours)	(100, 35.48)	(49.0, 46.37)	(56.0, 46.33)	(54.88, N/A)			
Mistral-	7B-Instruct-v0	.2 (Jiang et al.,	, 2023)				
M. L.L. (71 ( 1. 2022)	10.3	10.25	18.78	47.75			
MeLLo (Zhong et al., 2023)	(0, 10.31)	(59.0, 7.48)	(48.4, 7.92)	(47.75, N/A)			
ICE (Cohen et al., 2023)	<1	<1	<1	OOM			
IKE (Zheng et al., 2023a)	<1	<1	<1	OOM			
DeepEdit Wang et al. (2024)	<1	<1	<1	<1			
CWalls (Ours)	34.07	45.76	46.78	50.7			
Gwaik (Ours)	(0, 34.08)	(47, 45.69)	(51.2, 45.16)	(50.7, N/A)			
Meta-Llama-3-8B-Instruct (AI@Meta, 2024)							
Mal La (Zhang at al. 2023)	~1	1.13	4.72	16.58			
Mello (Zhong et al., 2023)	<1 <1	(17, <1)	(17.4, <1)	(16.58, N/A			
ICE (Cohen et al., 2023)	<1	<1	<1	OOM			
IKE (Zheng et al., 2023a)	<1	<1	<1	OOM			
DeepEdit Wang et al. (2024)	6.49	8.48	14.74	34.71			
DeepLuit Wang et al. (2024)	(0, 6.49)	(36.0, 6.92)	(36.20, 6.89)	(34.71, N/A			
GWalk (Ours)	70.12	73.28	76.61	84.01			
G Wark (Gurs)	(100, 70.1)	(84.0, 72.68)	(87, 72.8)	(84.01, N/A)			
Meta-Llama/L	lama-3.1-8B-I	nstruct (Dubey	r et al., 2024)				
GWalk (Ours)	74.68	76.34	77.74	83.32			
	(100, 74.66)	(85, 75.85)	(85.4, 74.91)	(83.32, N/A)			
Qwen/Qw	en2.5-7B-Instr	uct (Yang et al	., 2024a)				
	44.23	46.03	55.1	86.32			
(Walk (Ours)							

 1253
 1254 Table 18: Experiments on MQUAKE-REMASTERED-CF with numbers of edited cases and methods. *Total Accuracy* Results are reported in the format: (*Edited Accuracy*, Unedited Accuracy)

Mathad	Mathed MQUAKE-REMASTERED-CF						
Method	1-edit	1000-edit	3000-edit	6000-edit	9171-edit		
	vicuna-7b	-v1.5 (Zheng et	al., 2023b)				
MeI I o (Zhong et al. 2023)	22.55	21.54	17.79	12.62	6.95		
	(100, 22.54)	(8, 23.2)	(7.43, 22.83)	(7.28, 22.58)	(6.95, N/A)		
ICE (Cohen et al., 2023)	<1	OOM	OOM	OOM	OOM		
IKE (Zheng et al., 2023a)	<1	OOM	OOM	OOM	OOM		
DeepEdit Wang et al. (2024)	<1	<1	<1	<1	<1		
GWalk (Ours)	61.89	56.98	56.37	54.93	54.15		
G Walk (Guis)	(100, 61.89)	(56.2, 57.07)	(53.97, 57.54)	(53.27, 58.06)	(54.15, N/A)		
	Mistral-7B-In	struct-v0.2 (Jiai	ng et al., 2023)				
MaLL o (Zhong et al. 2022)	19.83	19.08	18.9	18.27	18.09		
Mello (Zhong et al., 2023)	(<1, 19.84)	(20.6, 18.9)	(19.47, 18.62)	(19.02, 16.87)	(18.09, N/A)		
ICE (Cohen et al., 2023)	<1	OOM	OOM	OOM	OOM		
IKE (Zheng et al., 2023a)	<1	OOM	OOM	OOM	OOM		
DeepEdit Wang et al. (2024)	<1	<1	<1	<1	<1		
CW-II- (O	61.42	57.79	56.35	53.73	51.53		
Gwalk (Ours)	(100, 61.42)	(51.8, 58.52)	(52.3, 58.32)	(50.93, 59.04)	(51.53, N/A)		
Meta-Llama-3-8B-Instruct (AI@Meta, 2024)							
MeLLo (Zhong et al., 2023)	<1	<1	<1	<1	<1		
ICE (Cohen et al., 2023)	<1	OOM	OOM	OOM	OOM		
IKE (Zheng et al., 2023a)	<1	OOM	OOM	OOM	OOM		
Deer Edit War a st sl (2024)	22.16	19.26	21.09	23.04	24.25		
DeepEdit wang et al. (2024)	(100, 22.15)	(21.29, 19.01)	(24.48, 19.44)	(23.77, 21.67)	(24.25, N/A)		
	74.09	73.67	72.4	71.62	70.08		
Gwalk (Ours)	(100, 74.09)	(71.1, 73.98)	(70.9, 73.13)	(70.33, 74.05)	(70.08, N/A)		
Meta-	Llama/Llama-	3.1-8B-Instruct	(Dubey et al., 2	2024)			
CWalls (Ques)	76.27	73.48	72.86	72.03	70.94		
Gwalk (Ours)	(1, 76.27)	(73.1, 73.53)	(71.98, 73.29)	(70.96, 74.08)	(70.94, N/A)		
Q	wen/Qwen2.5	-7B-Instruct (Ya	ang et al., 2024a	)			
CW-11- (O-m)	64.4	62.61	63.35	64.93	66.79		
Gwaik (Ours)	(0, 64.41)	(66.6, 62.12)	(66.34, 61.9)	(66.28, 62.4)	(66.79, N/A)		

1310Table 19: Additional experiments on meta-llama/Llama-3.1-8B-Instruct (Dubey et al., 2024) and1311Qwen/Qwen2.5-7B-Instruct (Yang et al., 2024a) on MQUAKE-REMASTERED-CF-6334. Results1312Total Accuracy

 Total Accuracy

 are reported in the format: (Test Edited Accuracy, Train Edited Accuracy, Unedited Accuracy).

Method				
Wiethou	100-edit	1000-edit	3000-edit	6344-edit
	Vicuna-7B-v	1.5 (Zheng et al., 2023	ib)	
ROME (Meng et al., 2022)	<1	<1	<1	<1
MEND (Mitchell et al., 2022)	12.75	10.36	9.56	7.24
	(11.11, 11, 13.25)	(7.33, 9.6, 13.64)	(6.1, 7.2, 11.9) <b>50 1</b>	(6.38, 6.49, 10.3)
GWalk (Ours)	(22.22, 64.84, 57.48)	(29.08, 66.17, 63.23)	(39.3, 63.74, 64.33)	(44.64, 62.11, 68.25
	Mistral-7B-Instr	uct-v0.2 (Jiang et al.,	2023)	
ROME (Meng et al., 2022)	<1	<1	<1	<1
MEND (Mitchell et al. 2022)	11.84	11.57	8.39	6.82
WILLIND (WITCHICH et al., 2022)	(11.11, 9, 12.36)	(6.95, 8.7, 12.12)	(3.41, 6.6, 10.1)	(2.33, 6.4, 8.4)
GWalk (Ours)	56.25	58.9	56.03	54.43
	(33.33, 37.14, 30.28)	(34.09, 60.37, 60.6)	(42.09, 59.04, 59.85)	(47.49, 57.74, 52.38
	Meta-Llama-3-8	B-Instruct (AI@Meta,	2024)	
ROME (Meng et al., 2022)	<1	<1	<1	<1
MEND (Mitchell et al., 2022)	13.04	13.3	9.81	7.42
((((((((((((((((((((((((((((((((((((((	(11.11, 10, 13.47)	(5.33, 8.4, 14.33)	(4.21, 8.63, 11.1)	(5.12, 7.45, 7.3)
GWalk (Ours)	<b>07.01</b> (33.33,74.73,66.92)	/1.89 (47.45,80.94,70.65)	/3./0 (54.05_81.6_71.12)	7 <b>4.</b> 22 (61.02 80.47 73.02
	Meta-Llama/Llama-3.	1-8B-Instruct (Dubey)	et al., 2024)	(01.02, 00.47, 73.02
	66.79	73.66	72.09	73.3
GWalk (Ours)	(33.33, 72, 66.66)	(49.47, 73.68, 73.02)	(51.23, 75.1, 70.6)	(55.39, 73.84, 71.55
	Qwen/Qwen2.5-7	B-Instruct (Yang et al.,	, 2024a)	
GWalk (Ours)	60.59	65.42	68.75	70.49
G traik (Ours)	(33.33, 62, 60.56)	(30.13, 68.6, 63.83)	(43.65, 69.9, 64.99)	(59.12, 70.51, 68.25