

Can We Edit LLMs for Long-Tail Biomedical Knowledge?

Anonymous ACL submission

Abstract

Knowledge editing has emerged as an effective approach for updating large language models (LLMs) by modifying their internal knowledge. However, their application to the biomedical domain faces unique challenges due to the long-tailed distribution of biomedical knowledge, where rare and infrequent information is prevalent. In this paper, we conduct the first comprehensive study to investigate the effectiveness of knowledge editing methods for editing *long-tail* biomedical knowledge. Our results indicate that, while existing editing methods can enhance LLMs’ performance on *long-tail* biomedical knowledge, their performance on long-tail knowledge remains inferior to that on high-frequency popular knowledge, even after editing. Our further analysis reveals that long-tail biomedical knowledge contains a significant amount of one-to-many knowledge, where one subject and relation link to multiple objects. This high prevalence of one-to-many knowledge limits the effectiveness of knowledge editing in improving LLMs’ understanding of long-tail biomedical knowledge, highlighting the need for tailored strategies to bridge this performance gap¹.

1 Introduction

Recently, knowledge editing (Meng et al., 2022a; Yao et al., 2023) has emerged as a promising approach to efficiently update large language models (LLMs) by injecting new knowledge into their internal knowledge (Touvron et al., 2023; Achiam et al., 2023). These methods have shown remarkable performance in enhancing LLMs’ performance across several general-domain tasks, such as question answering (QA) (Huang et al., 2023), knowledge injection (Li et al., 2024), and knowledge reasoning (Wang et al., 2024a).

While knowledge editing methods have proven effective in general-domain tasks, their application

¹Code: https://anonymous.4open.science/r/edit_bio_long_tail-951A/

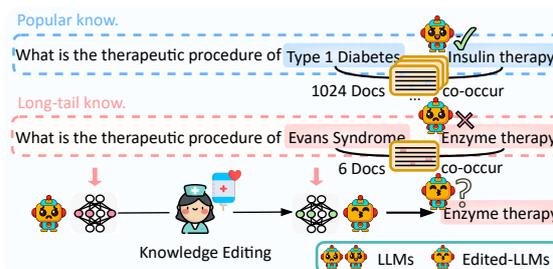


Figure 1: LLMs often struggle with long-tail biomedical knowledge, where entities co-occur in a few documents. Knowledge editing offers a potential solution by injecting this rare information into LLMs, improving their ability to handle such long-tail knowledge.

to the biomedical domain presents unique challenges (Wu et al., 2024b). Specifically, real-world biomedical data often exhibit a long-tailed distribution, with a small amount of popular knowledge and a large amount of long-tail knowledge that appears rarely or only once (Wu et al., 2024b; Delile et al., 2024). For example, the common disease “Type 1 Diabetes” is mentioned in over 106,138 papers in PubMed (Roberts, 2001), while a rare disease like “Evans Syndrome” appears in only about 23 papers (Wei et al., 2013). Recent studies indicate that the low frequency of knowledge in the pre-training corpus can hinder LLMs’ understanding of this knowledge (Kandpal et al., 2023; Wu et al., 2024b). Figure 1 illustrates an example where LLMs struggle with low-frequency biomedical knowledge. This is particularly problematic as LLMs are increasingly being used by healthcare professionals, including doctors, to assist in diagnosis and treatment recommendations (Tian et al., 2024). As LLMs become more integrated into clinical practice, their ability to accurately handle rare but critical biomedical knowledge becomes essential. This raises a critical question for knowledge editing in the biomedical domain:

Can knowledge editing methods effectively edit large language models to incorporate long-tail biomedical knowledge?

In this work, we present the first comprehensive study to investigate the effectiveness of knowledge editing for long-tail biomedical knowledge. We focus on biomedical knowledge represented as knowledge triples and leverage knowledge probing (Alghanmi et al., 2021) to evaluate whether LLMs have effectively acquired this knowledge. Specifically, knowledge probing is a technique that queries LLMs to assess their internal factual knowledge (Meng et al., 2022b). As illustrated in Figure 1, we query LLMs with questions generated from biomedical knowledge triples to determine whether they can correctly recall the target knowledge. By comparing the knowledge probing results of LLMs before and after editing, we can evaluate how effectively knowledge editing enhances LLMs’ ability to handle long-tail biomedical knowledge. Our key findings can be summarised as follows:

- LLMs struggle to capture long-tail biomedical knowledge through pre-training;
- Knowledge editing can improve LLMs’ performance on long-tail biomedical knowledge, but the post-edit performance still lags behind that of popular knowledge;
- Edited LLMs can memorise the form of long-tail knowledge, but their ability to generalise such knowledge is limited.
- The prevalence of one-to-many knowledge in long-tail biomedical knowledge is a key factor contributing to LLMs’ poor performance in capturing such long-tail knowledge;
- Effectively handling one-to-many knowledge is critical for improving LLMs’ performance on long-tail biomedical knowledge through knowledge editing.

2 Background and Definitions

This section defines long-tail biomedical knowledge and briefly introduces the knowledge probing and editing techniques used in our experiments.

2.1 Long-Tail Biomedical Knowledge

We present biomedical knowledge using knowledge triple $\langle s, r, o \rangle$, where s is the subject, r is the relation, and o is the object. Let \mathcal{D} be the set of documents in the pre-training corpus, and $\mathcal{D}(s, o)$ be the subset of documents where both s and o co-occur. We define the *co-occurrence number* of the knowledge triple as $|\mathcal{D}(s, o)|$, which represents the

frequency of knowledge $\langle s, r, o \rangle$ within the document set \mathcal{D} (Kandpal et al., 2023). In this paper, following Mallen et al. (2023) and Kandpal et al. (2023), we define *long-tail knowledge* as:

$$\mathcal{K}_1 = \{\langle s, r, o \rangle \mid |\mathcal{D}(s, o)| < \alpha\}, \quad (1)$$

where \mathcal{K}_1 denotes the set of long-tail knowledge and α represents a predefined threshold.

2.2 Knowledge Probing

Knowledge probing aims to evaluate LLMs’ ability to capture factual knowledge (Meng et al., 2022b), and can serve as an evaluation method to assess the effectiveness of knowledge editing (Hernandez et al., 2023). Specifically, given a subject s and a relation r in a triple $\langle s, r, o \rangle$, we use a manually designed template $\mathcal{T}(s, r)$ to generate a natural language question, which is then fed into an LLM f_θ to generate the object o as the answer. Following the work of Meng et al. (2022a) and Kassner et al. (2021), accuracy (ACC) is used to evaluate the performance of LLM in recalling the correct target entity o , which is formulated as:

$$\mathbb{E}_{\langle s, r, o \rangle \sim \mathcal{P}} \mathbb{I} \left\{ \arg \max_y f_\theta(y \mid \mathcal{T}(s, r)) = o \right\}, \quad (2)$$

where $\mathbb{E}_{\langle s, r, o \rangle \sim \mathcal{P}}$ denotes the expectation over a set of knowledge triples \mathcal{P} , y indicates the predicted answer and $\mathbb{I}\{\cdot\}$ is the indicator function. In this paper, we compare the knowledge probing results of LLMs before and after knowledge editing to investigate the effectiveness of editing methods in handling long-tail biomedical knowledge.

2.3 Knowledge Editing

Knowledge editing (Yao et al., 2023) aims to inject a new knowledge $\langle s, r, o \rangle$ into an LLM through a specific edit descriptor (x_e, y_e) (Yao et al., 2023). Given a knowledge $\langle s, r, o \rangle$ for editing, x_e can be formulated as $\langle s, r \rangle$, and $y_e = o$. The ultimate target of knowledge editing is to obtain an edited model f_{θ_e} , which effectively integrates the intended modifications within the editing scope, while preserving the model’s performance for out-of-scope unrelated facts:

$$f_{\theta_e}(x) = \begin{cases} y_e & \text{if } x \in I(x_e, y_e) \\ f_\theta(x) & \text{if } x \in O(x_e, y_e) \end{cases} \quad (3)$$

Here, the *in-scope* set $I(x_e, y_e)$ includes x_e and its equivalence neighborhood $N(x_e, y_e)$, which includes related input/output pairs. In contrast, the

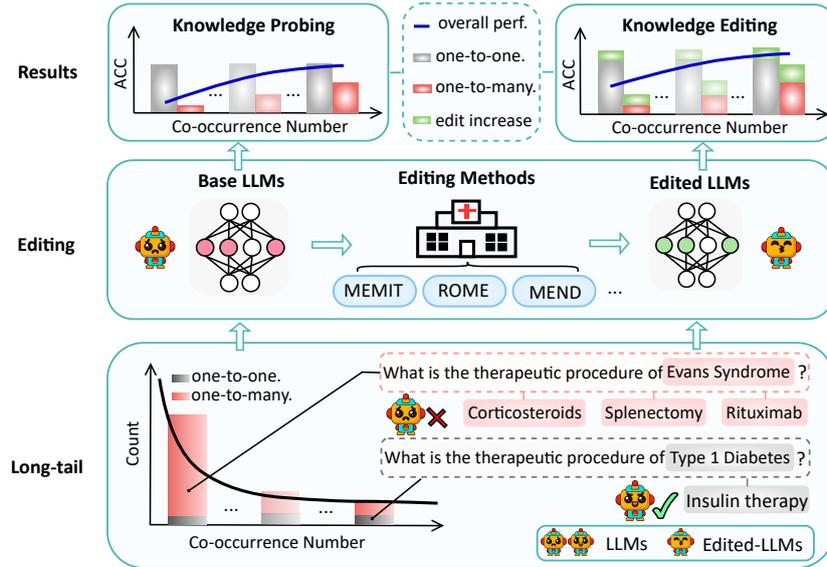


Figure 2: An overview of probing and editing for biomedical knowledge. These knowledge triples are classified into different groups based on co-occurrence number and further divided into one-to-one and one-to-many categories based on the number of correct answers (see § 4.4). The increasing performance with the number of co-occurrence number indicates that LLMs struggle to effectively capture long-tail biomedical knowledge before and after editing.

out-of-scope $O(x_e, y_e)$ contains inputs that are unrelated to the edit descriptor (x_e, y_e) .

3 Identifying Long-Tail Biomedical Knowledge

Due to the lack of biomedical datasets specifically designed to evaluate long-tail knowledge, we develop a pipeline to extract such knowledge. In this section, we outline the procedures for extracting long-tail biomedical knowledge, with further details provided in Appendix A.

Specifically, we focus on biomedical knowledge represented as knowledge triples. We extract triples from SNOMED CT (Donnelly et al., 2006), which is a comprehensive biomedical knowledge graph comprising over 1.4 million clinical triples (Benson and Grieve, 2021), and widely used for assessing LLMs’ understanding of biomedical knowledge (Meng et al., 2022b). To identify the long-tail knowledge within these triples, we use an entity linking pipeline to compute the co-occurrence number of each triple in the PubMed corpus², which is a widely used biomedical corpus for pre-training. In the entity linking pipeline, we first use PubTator (Wei et al., 2013) to annotate entities in the PubMed corpus and then use SapBERT (Liu et al., 2021) to link knowledge triple entities to PubMed entities.

²<https://pubmed.ncbi.nlm.nih.gov/>

Subsequently, we calculate the co-occurrence number for each triple. Long-tail knowledge is defined as triples with a co-occurrence number less than 10 (Kandpal et al., 2023). To evaluate LLMs’ ability to capture these triples, we generate question-answer pairs following Meng et al. (2022a). For each triple, we construct a question using the subject and relation, with the object serving as the answer. For example, for the triple $\langle \text{Diabetes}, \text{treated_by}, \text{Insulin} \rangle$, the corresponding QA pair is: *What is Diabetes treated by? Answer: Insulin.* The statistics of our extracted data are presented in Table 1 and the template for constructing questions is provided in Table 3. We refer to our dataset as CliKT (Clinical Knowledge Triples). Details of the construction process can be found in Appendix A and Figure 7.

4 Knowledge Editing for Long-Tail Biomedical Knowledge

In this section, we investigate the effectiveness of knowledge editing methods in enhancing LLMs’ ability to handle long-tail biomedical knowledge. Since some editing methods like MEND (Mitchell et al., 2022) and IKE (Zheng et al., 2023a) require additional training data, we follow Meng et al. (2022a) to divide our CliKT dataset into training, validation, and test sets (See Table 1), and report the results on the test set. Specifically, we detail

Item	Train	Valid	Test
# Triples	59,705	14,087	28,375
$ \mathcal{D}(s, o) < 10^1$	52,297	11,476	22,952
$ \mathcal{D}(s, o) \in [10^1, 10^2)$	5,363	2,055	4,110
$ \mathcal{D}(s, o) \in [10^2, 10^3)$	1,659	551	1,103
$ \mathcal{D}(s, o) \geq 10^3$	386	105	210
# Relations	21	21	21
# Subjects	39,654	12,267	21,872
# Objects	7,867	3,526	4,706

Table 1: The statistics of CliKT dataset. $|\mathcal{D}(s, o)|$ represents the oc-occurrence number of knowledge triple.

the experimental setup in § 4.1, and introduce the results of LLMs before and after editing in § 4.2 and § 4.3, respectively.

4.1 Experimental Setup

LLMs. In our experiments, we employ two widely used biomedical LLMs primarily pre-trained on the PubMed corpus: **BioGPT-Large** (Luo et al., 2022) and **BioMedLM** (Bolton et al., 2024). Additionally, we include two general-domain LLMs: **Llama2** (Touvron et al., 2023) and **GPT-J** (Wang and Komatsuzaki, 2021) to evaluate whether our findings generalise to models that are not specifically trained on biomedical data. Details of these LLMs are provided in Appendix B.1.

Knowledge Editing Methods. For knowledge editing, we employ the following methods, which have demonstrated strong effectiveness in knowledge injection tasks (Wang et al., 2025):

- **ROME** (Meng et al., 2022a): ROME updates an MLP layer to encode new information by treating the MLP module as a key-value memory. It relies on causal mediation analysis to precisely identify the location for editing.
- **MEMIT** (Meng et al., 2023): it employs the localisation strategies from ROME and applies explicit parameter adjustments to inject new knowledge across multiple layers.
- **MEND** (Mitchell et al., 2022): MEND enables efficient, targeted updates to LLMs by leveraging low-rank gradient transformations. It enables quick, localised modifications in model behaviour using only a single input-output example, while preventing overfitting.
- **IKE** (Zheng et al., 2023a): IKE modifies factual knowledge in LLMs through in-context learning without updating parameters. It corrects specific knowledge using demonstration

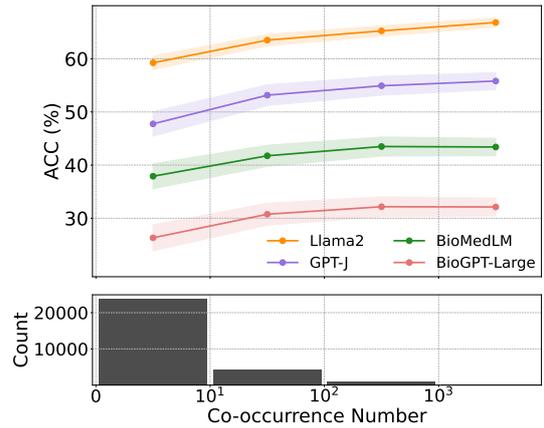


Figure 3: The overall performance of pre-edit probing on Llama2, GPT-J, BioMedLM and BioGPT-Large. The shaded areas indicate the standard deviation and Count denotes the number of triples within each group.

contexts, reducing over-editing and preserving previously stored knowledge.

- **FT** (Yao et al., 2023): FT updates model parameters using gradient descent on a single MLP layer identified by ROME. We employ the FT implementation within the EasyEdit framework (Wang et al., 2023b).

Evaluation Metrics. We use knowledge probing to evaluate whether LLMs have successfully acquired biomedical knowledge within the CliKT dataset. Specifically, we focus on the zero-shot QA performance of LLMs in answering questions from the CliKT dataset. The questions are used as inputs, and the accuracy (ACC) metric is employed to evaluate the correctness of the generated answers, as described in § 2.2.

In addition to knowledge probing, we follow previous works (Meng et al., 2022a; Yao et al., 2023) and use the following metrics to evaluate the comprehensive effectiveness of knowledge editing: (1) **Reliability**: This metric measures the mean accuracy on a specific collection of pre-defined input-output pairs (x_e, y_e) ; (2) **Generalisation**: Considering that paraphrased sentences should be modified accordingly by editing, this metric measures the average accuracy on equivalent neighbours $R(x_e, y_e)$; (3) **Locality**: This metric quantifies how often the predictions of the post-edit model remain unchanged for out-of-scope neighbours $O(x_e, y_e)$. Detailed definitions of these metrics are provided in Appendix B.2.

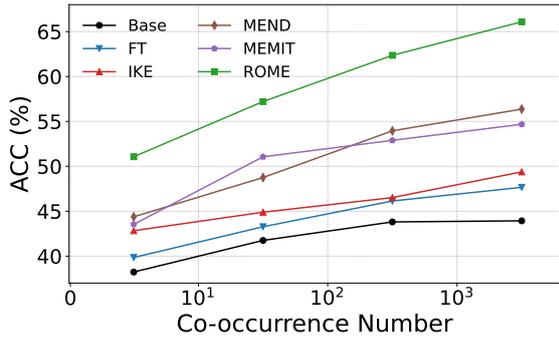


Figure 4: The performance of knowledge probing after editing with different editing methods on BioMedLM, where “Base” denotes LLM without editing.

4.2 Pre-Edit Results on Long-Tail Biomedical Knowledge

Finding 1: *LLMs struggle to capture long-tail biomedical knowledge through pre-training.*

To investigate whether LLMs face challenges in capturing long-tail biomedical knowledge during pre-training, we categorise biomedical knowledge triples in CliKT into different groups based on their co-occurrence number $|\mathcal{D}(s, o)|$ and evaluate the probing results of LLMs across these groups.

The bottom portion of Figure 3 shows the distribution of triples across different group, which highlights the long-tail nature of biomedical knowledge, where long-tail knowledge accounts for the majority of the data. The results for biomedical LLMs and general-domain LLMs are illustrated in the top portion of Figure 3. Specifically, Figure 3 shows that the performance of LLMs declines as the co-occurrence number decreases. In particular, the performance of BioMedLM on long-tail knowledge ($|\mathcal{D}(s, o)| < 10$) is 22.86% lower relative to its performance on popular knowledge ($|\mathcal{D}(s, o)| \geq 10^3$). This trend is also evident in general-domain LLMs. For example, Llama2 experiences an accuracy drop of 16.86% when handling long-tail biomedical knowledge compared with popular knowledge. These results indicate that LLMs struggle with long-tail biomedical knowledge, highlighting the challenge of accurately capturing long-tail knowledge during pre-training. Furthermore, Figure 3 shows that as the co-occurrence number decreases, the standard deviation of ACC increases. This observation implies that LLMs exhibit greater confidence when processing popular biomedical knowledge than long-tail biomedical knowledge.

Based on the above analysis, we conclude that

Group	Edit	Reliability [↑]	Gen. [↑]	Locality [↑]
<10 ¹	ROME	98.02	68.42	83.70
	MEMIT	86.21	<u>47.36</u>	98.10
	MEND	<u>91.32</u>	46.75	89.60
	IKE	83.87	43.70	<u>97.81</u>
	FT	32.52	40.36	96.80
[10 ¹ , 10 ²)	ROME	98.11	70.10	84.60
	MEMIT	89.21	48.21	97.30
	MEND	88.90	47.80	89.83
	IKE	84.52	45.12	96.80
	FT	33.35	40.78	97.90
[10 ² , 10 ³)	ROME	98.63	72.50	84.62
	MEMIT	89.01	<u>51.47</u>	97.90
	MEND	88.94	48.83	91.40
	IKE	85.89	46.74	<u>96.85</u>
	FT	33.89	44.62	96.66
≥ 10 ³	ROME	98.66	72.54	84.45
	MEMIT	89.87	<u>50.00</u>	<u>97.43</u>
	MEND	<u>90.96</u>	49.86	90.92
	IKE	85.91	48.76	96.87
	FT	34.84	44.62	97.57

Table 2: Performance of knowledge editing methods on the CliKT dataset across different co-occurrence number groups. The best performance per group is marked in boldface, while the second-best performance is underlined. [↑] indicates that higher values reflect better performance, and “Gen.” stands for Generalisation.

LLMs indeed struggle to capture long-tail biomedical knowledge. As long-tail knowledge constitutes the majority of biomedical data, it is crucial to explore methods that can effectively improve LLMs’ performance on long-tail biomedical knowledge.

4.3 Post-Edit Results for Long-Tail Biomedical Knowledge

Finding 2: *Knowledge editing can improve LLMs’ performance on long-tail biomedical knowledge, but the post-edit performance still lags behind that of popular knowledge.*

Subsequently, we investigate the effectiveness of knowledge editing for long-tail biomedical knowledge. We apply existing knowledge editing methods to inject biomedical knowledge from the CliKT dataset into LLMs and then follow the procedures in the pre-edit experiments for evaluation.

The post-edit probing results for BioMedLM are presented in Figure 4, while the results for other LLMs can be found in Figure 8. These results yield the following findings: (1) Knowledge editing methods, especially ROME, can enhance LLM’s ability in handling long-tail biomedical knowledge. For example, Figure 4 shows that BioMedLM edited with ROME achieves an improvement of approximately 52.08% in ACC on

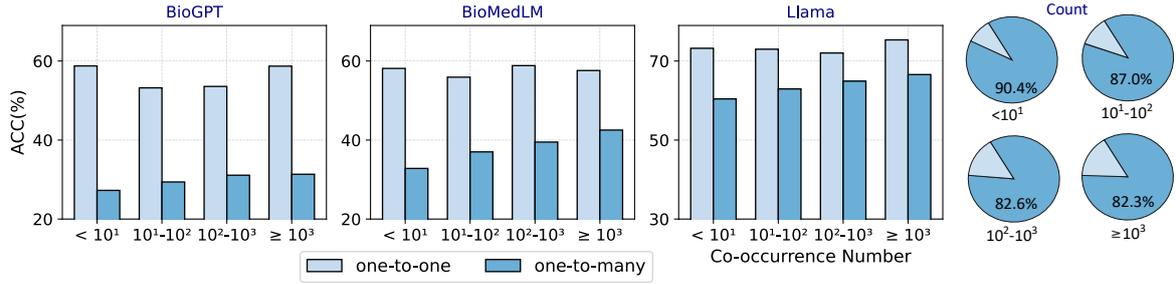


Figure 5: The comparison of knowledge probing performance between one-to-one and one-to-many settings across different co-occurrence Number, with the pie chart on the far right illustrating the data distribution.

347 long-tail knowledge ($|\mathcal{D}(s, o)| < 10$) compared to the base model before editing; (2) Despite the
 348 improvements from knowledge editing, Figure 4
 349 also reveals that ACC of post-edit LLMs consistently drops as the co-occurrence number decreases
 350 across all the editing methods. Specifically, for
 351 ROME, the ACC on long-tail knowledge is still
 352 16.15% relatively lower than on popular knowl-
 353 edge ($|\mathcal{D}(s, o)| \geq 10^3$). This indicates that even
 354 after editing, the edited LLMs continue to suffer
 355 from long-tail biomedical knowledge.
 357

358 **Finding 3:** Edited LLMs can memorise the form of
 359 long-tail knowledge, but their ability to generalise
 360 such knowledge is limited.

361 In addition to the post-edit probing results, we
 362 also calculate the other editing metrics outlined in
 363 §4.1 to comprehensively evaluate the effectiveness
 364 of the editing methods. Specifically, we calculate
 365 the Reliability, Generalisation and Locality metrics
 366 of edited models across different groups of biomed-
 367 ical knowledge. From the results in Table 2, we
 368 observe that ROME’s Reliability remains above
 369 98% across all groups, with no significant varia-
 370 tion. Similarly, the Reliability of MEMIT, MEND,
 371 and IKE is largely unaffected by the co-occurrence
 372 number, indicating that the edited LLMs’ ability
 373 to memorise the form of inserted knowledge is
 374 not influenced by long-tail knowledge. However,
 375 the generalisation performance declines as the co-
 376 occurrence number decreases, which aligns with
 377 the observed reduction in post-edit ACC for edited-
 378 LLMs as the co-occurrence number decreases.
 379 This observation suggests that, although edited
 380 LLMs can memorize the form of long-tail knowl-
 381 edge itself after knowledge editing, their ability
 382 to generalise this long-tail knowledge, especially
 383 in reasoning and responding to related questions,
 384 remains influenced by low co-occurrence numbers.

385 Furthermore, we observe that, though all the

386 editing methods exhibit relatively strong perform-
 387 ance in terms of locality across groups, ROME
 388 is affected more than the other methods. This in-
 389 dicates that while ROME achieves the best reli-
 390 ability and generalisation, it may slightly affect
 391 unrelated knowledge, consistent with the observa-
 392 tions of Wang et al. (Wang et al., 2024b).

393 4.4 In-depth Analysis

394 In this section, to further investigate the cause of
 395 the performance gap between long-tail and popu-
 396 lar biomedical knowledge before and after edit-
 397 ing, we further subdivide the data of long-tail
 398 and popular knowledge into *one-to-one* and *one-*
 399 *to-many* knowledge categories. The *one-to-one*
 400 knowledge means the subject is linked to a sin-
 401 gle object through the same relation, and *one-to-*
 402 *many* knowledge means the subject is linked to
 403 multiple objects through the same relation. For
 404 example, the triple $\langle \text{Type 1 diabetes, therapeutic}$
 405 $\text{procedure, insulin therapy} \rangle$ represents a one-to-one
 406 knowledge, where “Type 1 diabetes” is associated
 407 with a single object, “insulin therapy”. In contrast,
 408 $\langle \text{hypertension, associated with, heart disease} \rangle$ ex-
 409 emplifies a one-to-many knowledge, where “hyper-
 410 tension” can be linked to multiple objects, such as
 411 “stroke” or “kidney disease”.

412 4.4.1 Pre-Edit Probing of Different Types of 413 Knowledge

414 **Finding 4:** The prevalence of one-to-many knowl-
 415 edge in long-tail biomedical knowledge is a key
 416 factor contributing to LLMs’ poor performance in
 417 capturing such long-tail knowledge.

418 Figure 5 presents the pre-edit probing results
 419 of one-to-one and one-to-many knowledge across
 420 different co-occurrence number groups. We found
 421 that one-to-one knowledge is almost unaffected
 422 by co-occurrence numbers and consistently outper-
 423 forms one-to-many knowledge in all groups. For

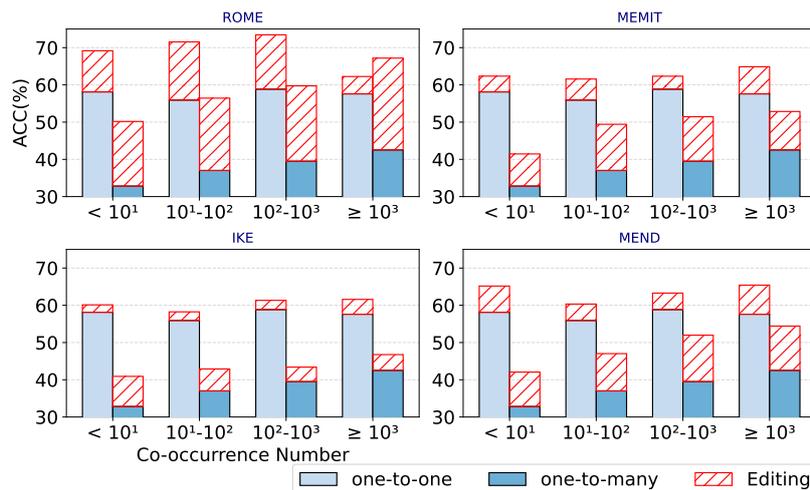


Figure 6: The knowledge probing performance of BioMedLM on both one-to-one knowledge and one-to-many knowledge before and after editing.

instance, BioGPT achieves an ACC that is approximately 115.56% higher on one-to-one knowledge compared to one-to-many knowledge. In contrast, for one-to-many knowledge, results from BioGPT, BioMedLM, and Llama2 all show a steady increase in ACC as the co-occurrence number increases. This suggests that co-occurrence number, or knowledge frequency, has a significant impact on LLMs’ ability to accurately comprehend one-to-many knowledge. We further analysed the distribution of one-to-one and one-to-many knowledge. Figure 5 shows that as the co-occurrence number increases, the proportion of one-to-many knowledge decreases while one-to-one knowledge increases. In the long-tail knowledge group ($|\mathcal{D}(s, o)| < 10$), 90.4% of the knowledge is one-to-many. This analysis reveals that LLMs’ difficulty with long-tail biomedical knowledge before editing is primarily due to the large proportion of one-to-many knowledge, which is challenging for LLMs to comprehend, as it increases the probability that the correct answers will not align with the model’s output.

4.4.2 Knowledge Editing for Different Types of Knowledge

Finding 5: *Effectively handling one-to-many knowledge is critical for improving LLMs’ performance on long-tail biomedical knowledge through knowledge editing.*

Next, we apply editing methods to both one-to-one and one-to-many knowledge. The results for BioMedLM are provided in Figure 6, while the results for other LLMs can be found in Figure 9. As

shown in Figure 6, while editing methods enhance performance on one-to-many knowledge, the improvement remains limited. For instance, in the ROME-edited BioMedLM for the long-tail knowledge ($|\mathcal{D}(s, o)| < 10$), the ACC for one-to-one knowledge was initially 42.19% higher than that for one-to-many knowledge. After applying the editing, this gap decreased to 16.43%. However, the persistent gap also highlights that even after editing, the model’s performance on one-to-many knowledge, which constitutes the majority of long-tail knowledge, remains constrained. This finding suggests that *despite knowledge editing can enhance LLMs’ capability in handling one-to-many knowledge, there remains a challenge in bridging the performance gap between one-to-one and one-to-many knowledge.* This limitation is critical given that one-to-many knowledge constitutes the majority of long-tail knowledge.

5 Related Work

5.1 LLMs for the Biomedical Domain

LLMs have made significant success in the biomedical domain, with an increasing variety of models contributing to advancements across different tasks (Tian et al., 2024). In the initial stages of their application, BERT (Vaswani et al., 2017) and its variants, such as BioBERT (Lee et al., 2020) and ClinicalBERT (Huang et al., 2019), demonstrated notable improvements in named entity recognition and relation extraction when applied to large datasets such as PubMed and clinical notes (Perera et al., 2020; Sun et al., 2021). GPT-based models, including GPT-J (Wang and

Komatsuzaki, 2021), BioGPT (Luo et al., 2022) and BioMedLM (Bolton et al., 2024), further enhanced biomedical text generation and question answering (Tian et al., 2024). Recent LLMs like Llama (Touvron et al., 2023), Falcon (Almazrouei et al., 2023), and Palm (Chowdhery et al., 2023) have scaled transformer architectures to address more complex tasks, such as biomedical knowledge reasoning (Wu et al., 2024a; Watanabe et al., 2024) and assisting in clinical decision-making (Sandmann et al., 2024). This work explores LLMs’ performance on long-tail biomedical knowledge. We present the first study to investigate how long-tail knowledge impacts LLMs in knowledge editing, offering new insights into improving LLMs’ handling of rare biomedical information through knowledge editing techniques.

5.2 Knowledge Editing

Knowledge editing methods can be broadly classified into three distinct categories (Yao et al., 2023): memory-based (Zheng et al., 2023b), meta learning (Mitchell et al., 2022), and locate-then-edit (Meng et al., 2022a). Memory-based methods, like IKE (Zheng et al., 2023b), enhance LLMs with external memory modules to update knowledge without changing the model’s parameters. Meta-learning approaches, such as KE (Cao et al., 2021), train a hyper-network to generate updated weights. MEND (Mitchell et al., 2022) improves on this by using low-rank gradient updates for more efficient model edits. However, meta-learning methods still require substantial computational resources and may unintentionally affect unrelated knowledge.

Locate-then-edit approaches aim for more targeted knowledge editing. Methods like KN (Dai et al., 2022) use knowledge attribution to locate relevant neurons but struggle with precise weight updates. ROME (Meng et al., 2022a) advances this by using causal tracing to locate and edit the Feed Forward Network (FFN) layers, which act as key-value memories (Geva et al., 2021, 2023). MEMIT (Meng et al., 2023) further expands this technique for batch editing. To the best of our knowledge, this work is the first to investigate the effectiveness of knowledge editing on long-tail biomedical knowledge.

5.3 Long-Tail Knowledge within LLMs

Existing studies have explored how long-tail knowledge, affects LLMs’ performance (Shin et al.,

2022; Han and Tsvetkov, 2022; Elazar et al., 2022; Mallen et al., 2023; Kandpal et al., 2023). Mallen et al. (2023) find that commonsense QA accuracy is strongly correlated with the frequency of entity popularity in the pre-training data from Wikipedia (Milne and Witten, 2008). Similarly, Elazar et al. (2022) employ causal inference to investigate how pre-training data statistics affect commonsense QA, highlighting how models rely on co-occurrence patterns between subjects, objects, and text to answer questions. More recently, Kandpal et al. (2023) explore the connection between the knowledge LLMs acquire for general-domain QA tasks and its frequency in the pre-training corpus, introducing comparative experiments involving model retraining and scaling.

Despite these findings, prior work has focused on general-domain QA, with the long-tail biomedical domain remaining largely unexplored (Wu et al., 2024b). This is especially concerning as LLMs are increasingly being used by healthcare professionals. Our research fills this gap by investigating the influence of long-tail biomedical knowledge on LLMs through knowledge probing and examining its impact on the effectiveness of knowledge editing. This is particularly problematic as LLMs are increasingly being used by healthcare professionals, including doctors, to assist in diagnosis and treatment recommendations.

6 Conclusion

In this paper, we investigated the effectiveness of knowledge editing methods for addressing the challenges of long-tail biomedical knowledge in LLMs. Our findings show that while existing techniques enhance performance on long-tail knowledge, their performance remains inferior to that on high-frequency popular knowledge. This problem is primarily attributed to the high presence of one-to-many knowledge in the biomedical domain, which complicates the models’ ability to effectively comprehend such knowledge. To address these challenges, we recommend the development of advanced editing techniques specifically tailored to long-tail knowledge. These techniques should prioritise strategies for effectively handling the intricacies of one-to-many knowledge scenarios, which are particularly common in the biomedical domain and remain a significant obstacle for current methods.

588 Limitations

589 We identify the following limitations of our
590 work: (1) First, our approach to extracting long-
591 tail knowledge is based on document-level co-
592 occurrence frequency (Kandpal et al., 2023), which
593 captures general patterns of occurrence but lacks
594 refinement at the sentence level. This limitation
595 may cause our analysis to miss finer patterns in
596 knowledge distribution, especially in instances
597 where sentence-level context provides essential nu-
598 ances. Future work could enhance the long-tail
599 knowledge extraction pipeline by investigating co-
600 occurrence on the sentence-level to improve the
601 granularity of knowledge editing. (2) Second, our
602 experimental framework is limited to the collection
603 of over 100,000 biomedical knowledge extracted
604 from PubMed, an extensive repository of biomed-
605 ical literature. While we believe the scale of this
606 collection offers a robust foundation for evaluat-
607 ing our methods, our future research should focus
608 on extracting long-tail knowledge from a broader
609 range of domains to further validate the generalis-
610 ability of our findings. (3) Finally, we concentrate
611 on analysing limitations without proposing spec-
612 ific solutions, prioritising the establishment of a
613 comprehensive understanding. Future work will fo-
614 cus on developing methods to improve knowledge
615 editing performance on long-tail knowledge.

616 References

617 Mohd Hafizul Afifi Abdullah, Norshakirah Aziz,
618 Said Jadid Abdulkadir, Hitham Seddig Alhassan Al-
619 hussian, and Noureen Talpur. 2023. Systematic liter-
620 ature review of information extraction from textual
621 data: recent methods, applications, trends, and chal-
622 lenges. *IEEE Access*, 11:10535–10562.

623 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
624 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
625 Diogo Almeida, Janko Altschmidt, Sam Altman,
626 Shyamal Anadkat, et al. 2023. GPT-4 technical re-
627 port. *arXiv preprint arXiv:2303.08774*.

628 Israa Alghanmi, Luis Espinosa Anke, and Steven
629 Schockaert. 2021. Probing pre-trained language
630 models for disease knowledge. In *Findings of the
631 Association for Computational Linguistics*, volume
632 ACL/IJCNLP 2021 of *Findings of ACL*, pages 3023–
633 3033.

634 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
635 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,
636 Mérouane Debbah, Étienne Goffinet, Daniel Hess-
637 low, Julien Launay, Quentin Malartic, et al. 2023.

The falcon series of open language models. *arXiv
preprint arXiv:2311.16867*. 638 639

Tim Benson and Grahame Grieve. 2021. *SNOMED CT*,
pages 293–324. Springer International Publishing,
Cham. 640 641 642

Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga,
David Hall, Betty Xiong, Tony Lee, Roxana
Daneshjou, Jonathan Frankle, Percy Liang, Michael
Carbin, et al. 2024. Biomedlm: A 2.7 b parameter
language model trained on biomedical text. *arXiv
preprint arXiv:2403.18421*. 643 644 645 646 647 648

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Edit-
ing factual knowledge in language models. In *Pro-
ceedings of the 2021 Conference on Empirical Meth-
ods in Natural Language Processing*, pages 6491–
6506. 649 650 651 652 653

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
Maarten Bosma, Gaurav Mishra, Adam Roberts,
Paul Barham, Hyung Won Chung, Charles Sutton,
Sebastian Gehrmann, et al. 2023. Palm: Scaling lan-
guage modeling with pathways. *Journal of Machine
Learning Research*, 24(240):1–113. 654 655 656 657 658 659

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao
Chang, and Furu Wei. 2022. Knowledge neurons in
pretrained transformers. In *Proceedings of the 60th
Annual Meeting of the Association for Computational
Linguistics*, pages 8493–8502. 660 661 662 663 664

Julien Delile, Srayanta Mukherjee, Anton Van Pamel,
and Leonid Zhukov. 2024. Graph-based retriever
captures the long tail of biomedical knowledge.
arXiv preprint arXiv:2402.12352. 665 666 667 668

Kevin Donnelly et al. 2006. Snomed-ct: The advanced
terminology and coding system for ehealth. *Studies
in health technology and informatics*, 121:279. 669 670 671

Yanai Elazar, Nora Kassner, Shauli Ravfogel, Amir
Feder, Abhilasha Ravichander, Marius Mosbach,
Yonatan Belinkov, Hinrich Schütze, and Yoav Gold-
berg. 2022. Measuring causal effects of data statis-
tics on language model’sfactual’ predictions. *arXiv
preprint arXiv:2207.14251*. 672 673 674 675 676 677

Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci,
Christophe Gravier, Jonathon Hare, Frederique
Laforest, and Elena Simperl. 2018. T-rex: A large
scale alignment of natural language with knowledge
base triples. In *Proceedings of the Eleventh Inter-
national Conference on Language Resources and
Evaluation (LREC 2018)*. 678 679 680 681 682 683 684

Hao Fei, Yafeng Ren, Yue Zhang, Donghong Ji, and
Xiaohui Liang. 2021. Enriching contextualized lan-
guage model from knowledge graph for biomedical
information extraction. *Briefings in bioinformatics*,
22(3):bbaa110. 685 686 687 688 689

690	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12216–12235.	745
691		746
692		747
693		748
694		749
695		
696	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495.	750
697		751
698		752
699		753
700		754
701	Xiaochuang Han and Yulia Tsvetkov. 2022. Orca: Interpreting prompted language models via locating supporting data evidence in the ocean of pretraining data. <i>arXiv preprint arXiv:2205.12600</i> .	755
702		9802–9822.
703		757
704		758
705	Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models. <i>arXiv preprint arXiv:2304.00740</i> .	759
706		760
707		761
708		762
709	Kexin Huang, Jaan Altonaar, and Rajesh Ranganath. 2019. Clinicalbert: Modeling clinical notes and predicting hospital readmission. <i>arXiv preprint arXiv:1904.05342</i> .	763
710		764
711		765
712		
713	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron. In <i>The Eleventh International Conference on Learning Representations</i> .	766
714		767
715		768
716		769
717		770
718		771
719		772
720		773
721		774
722		775
723	Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. Multilingual LAMA: investigating knowledge in multilingual pretrained language models. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3250–3258.	776
724		777
725		778
726		779
727		780
728		
729	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. <i>Bioinformatics</i> , 36(4):1234–1240.	781
730		782
731		783
732		784
733		
734	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2024. Pmet: Precise model editing in a transformer. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18564–18572.	785
735		786
736		787
737		788
738		
739	Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. 2021. Self-alignment pretraining for biomedical entity representations. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 4228–4238.	789
740		790
741		791
742		792
743		793
744		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

799	Seongjin Shin, Sang-Woo Lee, Hwijeen Ahn, Sungdong Kim, HyoungSeok Kim, Boseop Kim, Kyunghyun Cho, Gichang Lee, Woo-Myoung Park, Jung-Woo Ha, and Nako Sung. 2022. On the effect of pre-training corpora on in-context learning by a large-scale language model. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 5168–5186.	854
800		855
801		856
802		857
803		
804		858
805		859
806		860
807		
808	Cong Sun, Zhihao Yang, Lei Wang, Yin Zhang, Hongfei Lin, and Jian Wang. 2021. Biomedical named entity recognition using bert in the machine reading comprehension framework. <i>Journal of Biomedical Informatics</i> , 118:103799.	
809		
810		
811		
812		
813	Shubo Tian, Qiao Jin, Lana Yeganova, Po-Ting Lai, Qingqing Zhu, Xiuying Chen, Yifan Yang, Qingyu Chen, Won Kim, Donald C Comeau, et al. 2024. Opportunities and challenges for chatgpt and large language models in biomedicine and health. <i>Briefings in Bioinformatics</i> , 25(1):bbad493.	861
814		862
815		863
816		864
817		
818		
819	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	865
820		866
821		867
822		868
823		
824		
825	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	869
826		870
827		871
828		872
829		873
830	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 billion parameter autoregressive language model.	874
831		875
832	Benyou Wang, Qianqian Xie, Jiahuan Pei, Zhihong Chen, Prayag Tiwari, Zhao Li, and Jie Fu. 2023a. Pre-trained language models in biomedical domain: A systematic survey. <i>ACM Computing Surveys</i> , 56(3):1–52.	876
833		877
834		878
835		
836		
837	Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao, Jiarong Xu, and Fandong Meng. 2024a. Cross-lingual knowledge editing in large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics</i> , pages 11676–11686.	879
838		880
839		881
840		882
841		883
842		884
843	Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-jun Chen. 2024b. Wise: Rethinking the knowledge memory for lifelong model editing of large language models. <i>arXiv preprint arXiv:2405.14768</i> .	885
844		886
845		887
846		888
847		889
848	Peng Wang, Ningyu Zhang, Bozhong Tian, Zekun Xi, Yunzhi Yao, Ziwen Xu, Mengru Wang, Shengyu Mao, Xiaohan Wang, Siyuan Cheng, et al. 2023b. Easyedit: An easy-to-use knowledge editing framework for large language models. <i>arXiv preprint arXiv:2308.07269</i> .	891
849		892
850		893
851		894
852		895
853		896
		897
	Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2025. Knowledge editing for large language models: A survey. <i>ACM Comput. Surv.</i> , 57(3):59:1–59:37.	
	Natsumi Watanabe, Kudoro Kinaseka, and Akira Nakamura. 2024. Empower llama 2 for advanced logical reasoning in natural language understanding.	
	Chih-Hsuan Wei, Alexis Allot, Robert Leaman, and Zhiyong Lu. 2019. Pubtator central: automated concept annotation for biomedical full text articles. <i>Nucleic acids research</i> , 47(W1):W587–W593.	
	Chih-Hsuan Wei, Hung-Yu Kao, and Zhiyong Lu. 2013. Pubtator: a web-based text mining tool for assisting biocuration. <i>Nucleic acids research</i> , 41(W1):W518–W522.	
	Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024a. Pmc-llama: toward building open-source language models for medicine. <i>Journal of the American Medical Informatics Association</i> , page ocae045.	
	Zheng Wu, Kehua Guo, Entao Luo, Tian Wang, Shoujin Wang, Yi Yang, Xiangyuan Zhu, and Rui Ding. 2024b. Medical long-tailed learning for imbalanced data: bibliometric analysis. <i>Computer Methods and Programs in Biomedicine</i> , page 108106.	
	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10222–10240.	
	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023a. Can we edit factual knowledge by in-context learning? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4862–4876.	
	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023b. Can we edit factual knowledge by in-context learning? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4862–4876, Singapore. Association for Computational Linguistics.	

Appendix

In the Appendix, we introduce more details along with dataset construction, additional experimental results, discussions, and related works:

- **Appendix A:** CliKT Construction (cf. Section 3).
- **Appendix B:** Experimental Details (cf. Section 2 and 3).
- **Appendix C:** Additional Results (cf. Section 3).

A CliKT Construction

Due to the lack of datasets dedicated for evaluating long-tail biomedical knowledge, we propose CliKT, a new benchmark specifically designed to evaluate LLMs’ performance on long-tail biomedical knowledge. Notably, given that PubMed is a widely used biomedical corpus for pre-training LLMs (Wang et al., 2023a), which contains over 37 million abstracts of biomedical papers (Wei et al., 2013), we mainly focus on PubMed data to extract long-tail biomedical knowledge. Specifically, we first extract knowledge triples from SNOMED CT (Donnelly et al., 2006) (§A.1) to obtain a comprehensive set of biomedical concepts and their relationships. Next, we employ an entity linking pipeline to map these triples back to their corresponding documents in the PubMed (Roberts, 2001) corpus (§A.2), enabling us to identify whether a triple represents long-tail knowledge based its occurrence in the corpus. Finally, we generate question-answer (QA) pairs based on the knowledge triples to evaluate the ability of LLMs to capture the factual knowledge, and conduct a human evaluation to show that our entity linking pipeline accurately identifies relevant documents for the majority of the QA pairs.

A.1 Extracting Biomedical Knowledge Triples

We focus on the long-tail biomedical knowledge from the PubMed corpus. However, directly extracting such knowledge from the entire corpus is a challenging task (Shetty and Ramprasad, 2021; Nguyen et al., 2021; Abdullah et al., 2023). Therefore, following previous work (Alghanmi et al., 2021; Fei et al., 2021), we leverage information from existing biomedical knowledge graphs to facilitate more efficient extraction. Specifically, we extract all the knowledge triples from SNOMED CT (Donnelly et al., 2006), which is a comprehensive biomedical knowledge graph comprising over 200K triples and widely used for assessing LLMs’ understanding of biomedical knowledge (Meng et al., 2022b). Each triple is denoted as (head entity, relation, tail entity), representing the relationship between two entities, e.g., (Type 1 Diabetes, Therapeutic Procedure, Insulin therapy).

A.2 Mapping Knowledge Triples to PubMed Documents

We then develop an entity linking pipeline to map the extracted knowledge triples back to documents in Pubmed (Roberts, 2001) to identify long-tail knowledge. The detailed procedure is as follows:

Entity Annotation. To facilitate the mapping of knowledge triples to specific PubMed documents, we first need to annotate the entities within the PubMed corpus. To this end, we use PubTator (Wei et al., 2013), a robust web-based text-mining tool that provides automatic annotations of biomedical concepts in PubMed. Following the work of Wei et al. (2019), we obtain entity annotations within 37 million PubMed abstracts³.

Entity Linking. After obtaining annotated entities, the next step is to map the knowledge triples to their corresponding PubMed documents. Previous studies (Elsahar et al., 2018; Kandpal et al., 2023) suggest that when the head entity and the tail entity of a knowledge triple co-occur within a document, it is likely that the knowledge represented by the triple is expressed in that document. Based on this observation, we define documents where both the head and tail entities of a knowledge triple co-occur as its *related documents*, and the count of such documents as the *co-occurrence number*.

³The annotated data is available at <https://ftp.ncbi.nlm.nih.gov/pub/lu/PubTatorCentral/>

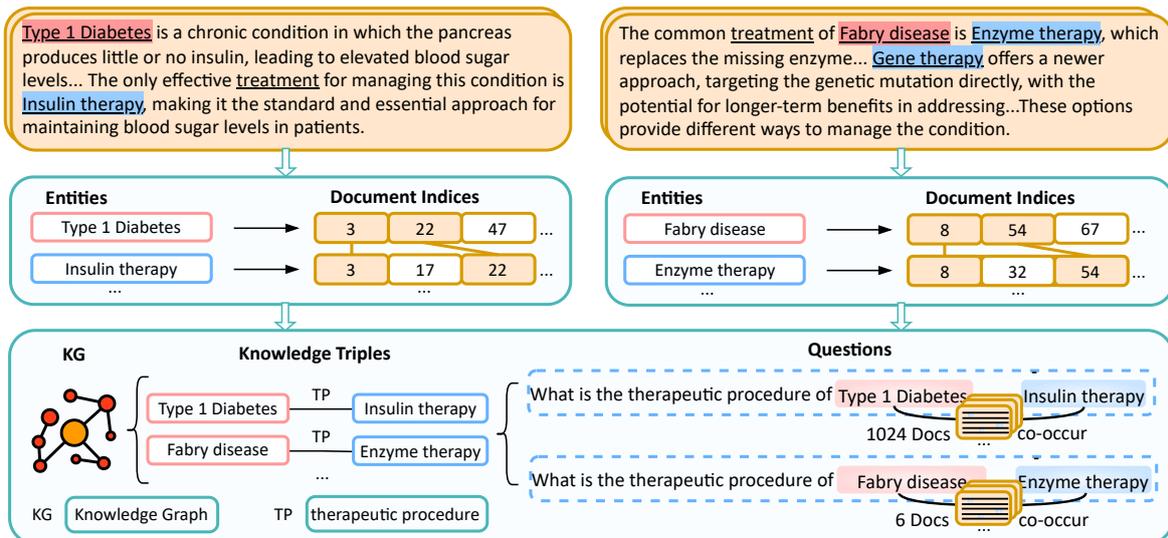


Figure 7: The pipeline for identifying long-tail biomedical knowledge consists of a systematic process encompassing documents collection, entity linking, knowledge graph traversal, and question generation.

To determine whether both the head and tail entities of a triple co-occur in a document, we use SapBERT (Liu et al., 2021), an effective biomedical entity linking model, to match these entities to those present in the document. For instance, given the triple (Hypertension, causes, heart disease) from SNOMED CT, SapBERT can link “Hypertension” to its equivalent term “high blood pressure” in PubMed, ensuring an accurate match with related documents. We iterate through the entire corpus to calculate the co-occurrence number for each triple. We define triples with a low co-occurrence number as long-tail biomedical knowledge.

Question Generation. Finally, we generate QA pairs based on the resulting triples to assess the LLMs’ ability to capture these knowledge triples. Following Meng et al. (2022a), we manually design templates to generate questions using the head entity and the relation, while considering the tail entity as the answer. For example, given a triple (Diabetes, treated_by, Insulin), the corresponding QA pair would be: *Question: What is Diabetes treated by? Answer: Insulin.*

B Experimental Details

B.1 Details of Large Language Models

We employ two biomedical LLMs and two general-domain LLMs in our experiments:

- **BioGPT-Large (Luo et al., 2022):** A 1.5 billion parameter model from Microsoft, primarily pre-trained on PubMed, excelling in drug discovery and medical record analysis.
- **BioMedLM (Bolton et al., 2024):** A Stanford-developed model optimised for biomedical tasks, pretrained on PubMed with 2.7 billion parameters, ideal for literature retrieval and information extraction.
- **Llama2 (Touvron et al., 2023):** A Meta-developed model with 7 billion parameters, designed for general-purpose language tasks. It has been leveraging large-scale pretraining on diverse datasets, including biomedical corpora.
- **GPT-J (Wang and Komatsuzaki, 2021):** A 6 billion parameter open-source model by EleutherAI, trained on the Pile dataset, which includes a significant portion of biomedical texts from PubMed.

B.2 Details of Evaluation Metrics

(1) **Reliability:** This metric measures the average accuracy over a predefined set of input-output pairs (x_e, y_e) . It is aimed to evaluate the ability of memorising the form of edit Prompt after knowledge editing.

$$\mathbb{E}_{x'_e, y'_e \sim \{(x_e, y_e)\}} \mathbf{1} \left\{ \underset{y}{\operatorname{argmax}} f_{\theta_e}(y | x'_e) = y'_e \right\} \quad (4)$$

Relation	Template
Finding site	Edit Prompt: “The finding site of [SUBJECT] is.” Question: “What is the finding site of [SUBJECT]?” Rephrase: “Where is [SUBJECT] typically found?”
Associated morphology	Edit Prompt: “The associated morphology of [SUBJECT] is.” Question: “What is the associated morphology of [SUBJECT]?” Rephrase: “Can you describe the morphology associated with [SUBJECT]?”
Causative agent	Edit Prompt: “The causative agent of [SUBJECT] is” Question: “What is the causative agent of [SUBJECT]?” Rephrase: “Which pathogen causes [SUBJECT]?”
Interprets	Edit Prompt: “[SUBJECT] interprets.” Question: “What does [SUBJECT] interprets?” Rephrase: “What is interpreted by [SUBJECT]?”
Procedure site	Edit Prompt: “The procedure site of [SUBJECT] is” Question: “What is the indirect procedure site of [SUBJECT]?” Rephrase: “Where is the procedure site for [SUBJECT]?”
Pathological process	Edit Prompt: “The pathological process of [SUBJECT] involves.” Question: “What is the pathological process of [SUBJECT]?” Rephrase: “Which pathological process does [SUBJECT] involve?”
Due to	Edit Prompt: “[SUBJECT] is due to.” Question: “What is the [SUBJECT] due to?” Rephrase: “What is the cause of [SUBJECT]?”
Has active ingredient	Edit Prompt: “The active ingredient of [SUBJECT] is.” Question: “What is the active ingredient of [SUBJECT]?” Rephrase: “What active ingredient does [SUBJECT] have?”
Part of	Edit Prompt: “[SUBJECT] is a part of.” Question: “What is the [SUBJECT] a part of?” Rephrase: “To what is [SUBJECT] a part?”
Has definitional manifestation	Edit Prompt: “The definitional manifestation of [SUBJECT] is.” Question: “What is the definitional manifestation of [SUBJECT]?” Rephrase: “How is [SUBJECT] manifested definitionally?”
Component	Edit Prompt: “The component of [SUBJECT] is.” Question: “What is the component of [SUBJECT]?” Rephrase: “What components does [SUBJECT] consist of?”

Table 3: Examples of relation templates demonstrate how each relation is transformed into input prompts, which can be categorized into three parts: Edit Prompt, Question, and Rephrase. The “Edit Prompt” is used for knowledge editing and reliability evaluation, the “Question” is designed for knowledge probing, and the “Rephrase” is used to assess generalisation metrics. The complete template for all the relations can be found in our github repository.

(2) **Generalisation:** Considering that paraphrased sentences are modified accordingly through editing, this metric measures the average accuracy on equivalent neighbours $R(x_e, y_e)$, where equivalent neighbours are rephrased questions based on the edited knowledge.

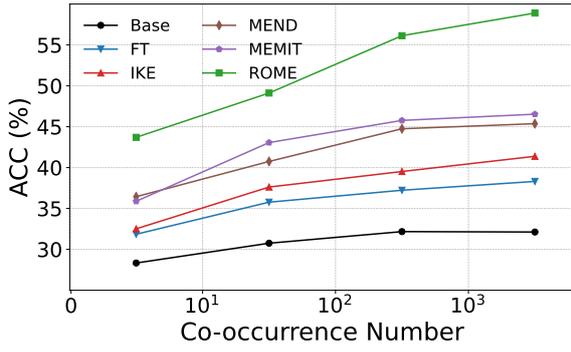
$$\mathbb{E}_{x'_e, y'_e \sim R(x_e, y_e)} \mathbf{1} \left\{ \underset{y}{\operatorname{argmax}} f_{\theta_e}(y | x'_e) = y'_e \right\} \quad (5)$$

(3) **Locality:** This metric measures the frequency with which the predictions of the post-edit model remain consistent for out-of-scope neighbors $O(x_e, y_e)$.

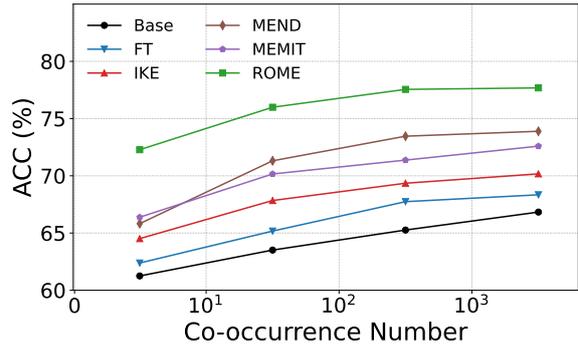
$$\mathbb{E}_{x'_e, y'_e \sim O(x_e, y_e)} \mathbf{1} \left\{ f_{\theta_e}(y | x'_e) = f_{\theta}(y | x'_e) \right\} \quad (6)$$

C Additional Results

We present the performance of knowledge editing on the other base LLMs in this section. Specifically, the performance of knowledge probing after editing with different editing methods on BioGPT and Llama2



(a) The performance on BioGPT.



(b) The performance on Llama2.

Figure 8: The performance of knowledge probing after editing with different editing methods on BioGPT and Llama2, where “Base” denotes LLM without editing.

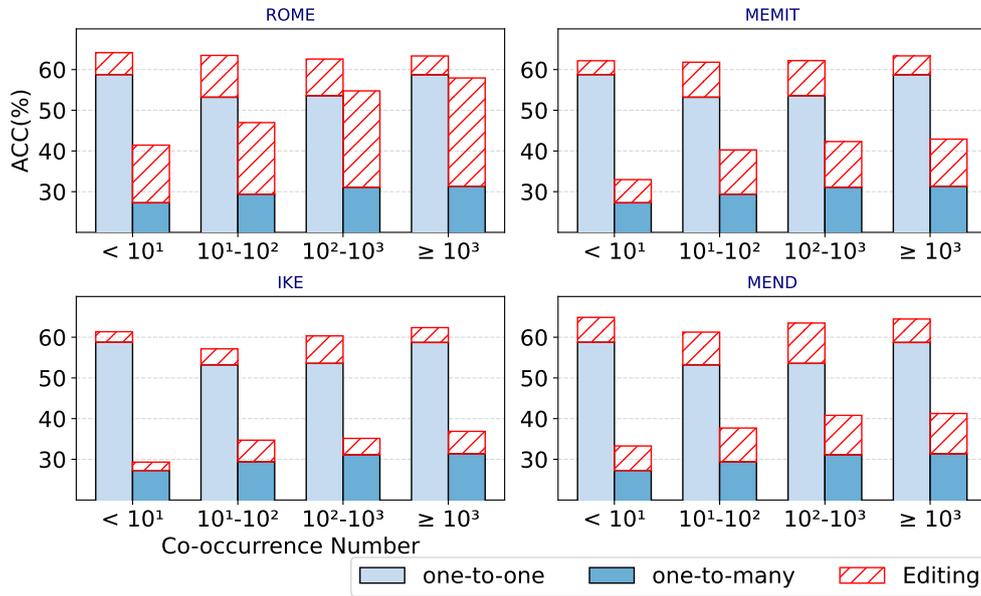


Figure 9: The knowledge probing performance of BioGPT on both one-to-one knowledge and one-to-many knowledge before and after editing.

can be seen in figure 8(a) and figure 8(b). We have also conducted the further analysis on BioGPT and Llama2, which can be seen in figure 9 and figure 10.

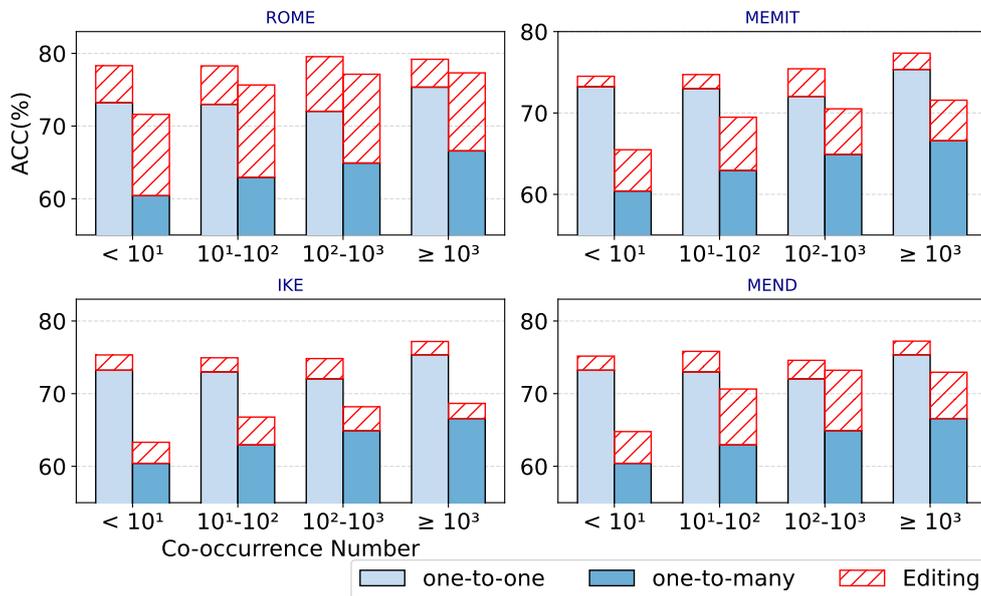


Figure 10: The knowledge probing performance of Llama2 on both one-to-one knowledge and one-to-many knowledge before and after editing.