

# Comparing Knowledge Injection Methods for LLMs in a Low-Resource Regime

Hugo Abonizio<sup>1,4</sup>, Thales Almeida<sup>2,4</sup>, Roberto Lotufo<sup>1,3</sup>, Rodrigo Nogueira<sup>4</sup>

<sup>1</sup>Faculdade de Engenharia Elétrica e de Computação (FEEC), University of Campinas (Unicamp)

<sup>2</sup>Instituto de Computação (IC), University of Campinas (Unicamp)

<sup>3</sup>NeuralMind

<sup>4</sup>Maritaca AI, Campinas, SP – Brazil

**Abstract**—Large language models (LLMs) often require vast amounts of text to effectively acquire new knowledge. While continuing pre-training on large corpora or employing retrieval-augmented generation (RAG) has proven successful, updating an LLM with only a few thousand or million tokens remains challenging. In this work, we investigate the task of injecting small, unstructured information into LLMs and its relation to the catastrophic forgetting phenomenon. We use a dataset of recent news – ensuring no overlap with the model’s pre-training data – to evaluate the knowledge acquisition by probing the model with question-answer pairs related the learned information. Starting from a continued pre-training baseline, we explored different augmentation algorithms to generate synthetic data to improve the knowledge acquisition capabilities. Our experiments show that simply continuing pre-training on limited data yields modest improvements, whereas exposing the model to diverse textual variations significantly improves the learning of new facts – particularly with methods that induce greater variability through diverse prompting. Furthermore, we shed light on the forgetting phenomenon in small-data regimes, illustrating the delicate balance between learning new content and retaining existing capabilities. We also confirm the sensitivity of RAG-based approaches for knowledge injection, which often lead to greater degradation on control datasets compared to parametric methods. Finally, we demonstrate that models can generate effective synthetic training data themselves, suggesting a pathway toward self-improving model updates. All code and generated data used in our experiments are publicly available, providing a resource for studying efficient knowledge injection in LLMs with limited data at <https://github.com/hugoabonizio/knowledge-injection-methods>.

**Index Terms**—Large language models, Knowledge injection, Data augmentation, Synthetic data

## I. INTRODUCTION

Previous work has shown that large language models (LLM) trained with a self-supervised objective learn large amounts of information, effectively acting as knowledge bases [1]–[5]. Likewise, updating the knowledge of a model through continued pre-training [6], [7], with the goal of specializing an existing LLM in a domain such as mathematics [8], [9], medicine [10], [11] or code [12], has also shown to be fruitful. In these scenarios, the amount of training data spans from billions to trillions of tokens.

However, perhaps surprisingly, incorporating relatively small amounts of information (e.g., thousands or millions of tokens) has proven to be more challenging, often resulting in performance degradation or only marginal gains when done

naively [13]–[17], and potentially even increasing *hallucinations* [18]. An illustrative case is that, while one might expect fine-tuning on new data using self-supervision to enable the LLM to internalize additional information into its parametric knowledge, it has not been shown to seamlessly learn certain relational inferences (e.g., from “A is B” to “B is A”) [19].

An additional complication is the forgetting problem, in which new information can often be successfully injected at the expense of forgetting previously learned knowledge [17], [20]–[24], a phenomenon commonly referred to as catastrophic forgetting [25]–[27].

To address these problems, recent literature on knowledge injection often concentrates on two extremes. In the first, also referred to as model editing, methods that require the information to be learned can be expressed as well-defined entities and relations [19], [28]–[30]. However, the techniques proposed in these works cannot be easily applied to real documents without modifications, as it is challenging to convey the complex information as a knowledge graph with a finite number of possible relationships. For example, converting the following headline into a set of discrete triples presents a non-trivial challenge:

*“Advocates for Ukraine’s surprise incursion into the Russian territory say it will provide Kyiv with vital leverage for any future peace talks (December 2nd, 2024)”*

On the other end of the spectrum, works on domain adaptation methods often builds upon datasets on the order of billions of tokens to learn new information [12], [16], [31], [32]. These approaches require substantial computational resources and are, therefore, only feasible for organizations with large-scale infrastructure and access to vast amounts of data. In many private applications or niche domains, however, data can be scarce, and compute resources are limited, making such large-scale methods impractical in some real-world settings.

In parallel, retrieval-augmented generation (RAG) [33] methods aim to inject knowledge through in-context learning, rather than by updating the model’s parametric knowledge. Some studies have compared these two approaches, highlighting different advantages in each case [34], [35]. However, it is important to note that these approaches are orthogonal: parametric knowledge injection and in-context learning methods can be combined and are not mutually exclusive.

In this work, we focus on the middle ground between the

two extremes. Our goal is to study the learning dynamics of small, unstructured information in an efficient manner, without forgetting previously acquired knowledge. We investigate the learning–forgetting tradeoff using a small corpus and evaluate different continual pre-training techniques including augmentation algorithms that leverage synthetic variations of the original data, aiming to overcome the challenges of learning in a small-data regime.

To measure the learning effectiveness, we need a dataset containing information that is both new to the model – i.e., it was not seen during its pre-training – and complex, to make the results impactful for real-world applications. Thus, we chose the TiEBE dataset [36], which contains news articles with the required recency, and question-answer pairs to measure knowledge and understanding.

Our results indicate that directly continuing pre-training on a set of documents using self-supervised learning [6] (i.e., next-token prediction) has limited effectiveness. Additionally, augmentation methods suggest that training on multiple versions of the same source document is necessary. Our findings also shed light on why LLMs require vast amounts of training data: intuitively, learning a small piece of information should not require 20 variations – one or two should suffice. Addressing this inefficiency could make training these models – which currently cost millions of dollars – significantly more affordable.

In summary, our main contributions are as follows:

- We conduct an extensive analysis of different knowledge injection approaches, including RAG, continued pre-training and augmentation techniques. We compare different techniques and their variations, providing evidence that models benefit from exposure to diverse data variations to effectively learn new information.
- We propose an evaluation methodology to assess the effectiveness of knowledge injection using small and unstructured datasets.
- We provide insights into the challenges of continued pre-training in small-scale data regimes and propose strategies to address training instabilities in this context.
- We release a set of synthetically augmented corpora along with the code to reproduce and expand them, supporting future research.

## II. RELATED WORK

Recent work on injecting new knowledge into LLMs can be divided into three categories: (1) model editing techniques that modify entity relationships within the model parameters, (2) knowledge injection through in-context learning, and (3) knowledge injection via some form of continual training. In this paper, we turn our attention to the latter two, examining the dichotomy between retrieval-augmented generation and fine-tuning.

The majority of model editing works focus on entity-centric tasks [37]. For instance Zhang et al. [38] propose KnowEdit, a benchmark consisting mostly of entity-centric datasets, such as Wikidata [39], ZsRE [40], and WikiBio [41].

Eva-KELLM [42] is another example of benchmark, in which the source of information to be learned comes from a raw document instead of a triple. However, in their evaluation methodology they still use triples to evaluate whether the knowledge was successfully learned. In our work, we depart from this constraint and evaluate the effectiveness of knowledge injection techniques by letting the LLM generate free-form answer for a given prompt that probes whether the model knows a particular piece of information.

A common approach to knowledge ingestion involves coupling the LLM with a retriever that has access to a database containing the relevant information [43], [44], a method generally referred to as RAG [33]. Because it does not require model training, it is relatively cost-effective for real-world scenarios. However, later in this work, we will demonstrate that this approach has some drawbacks. Some of the most frequently discussed in the literature include its dependence on retriever quality [45], the chunking strategy [46], and the model’s ability to handle the provided context [47].

Finally, in the last category of related work are methods that leveraged continued training, either by continual pre-training or fine-tuning. Here, we refer to continual pre-training as the self-supervised training using next-token prediction, while fine-tuning methods are the ones trained using instruction-tuning [48].

Ovadia et al. [15] compared RAG with continual pre-training on synthetic paraphrases, finding that exposing the new information in diverse ways through paraphrases play an important role. Their results found that RAG knowledge injection to be more effective than self-supervised finetuning. The work also generated questions to probe the effectiveness of knowledge injection.

Cheng et al. [16], [31] proposed methods for synthetically augmenting the pre-training corpus by transforming the original examples using different tasks. Their results show advantages over vanilla pre-training in the original documents, highlighting the importance of variations for effective learning.

Balaguer et al. [34] and Mecklenburg et al. [35] compared RAG and fine-tuning methods by training on synthetic pairs of questions and answers. Both work showed a significant increase in performance after fine-tuning and Balaguer et al. showed that RAG and fine-tuning can be combined synergistically.

Wu et al. [49] also investigated the knowledge injection through fine-tuning on question-answer (QA) pairs and used recent news as one of their evaluation datasets. Their results showed that learning from this fine-tuning is limited.

Yang et al. [50] introduced the EntiGraph, a knowledge graph-based augmentation technique for generating diverse synthetic text. They continued pretraining on small, domain-specific corpora and showed that their approach outperformed standard continued pretraining and simple paraphrasing using question-answering accuracy. Additionally, they found that their approach is complementary to RAG, improving downstream performance even when the original documents were available at inference time. While their work is closely related

to ours, we place a stronger emphasis on mitigating data contamination risks and explore knowledge injection with an even smaller training corpus.

### III. METHODS

In this section, we describe our evaluation methodology, including the dataset used and the knowledge injection techniques evaluated.

#### A. Dataset

We chose to inject knowledge related to recent news articles for two reasons. First, news articles carry complex forms of knowledge expression, which we argue are more aligned with real-world challenges researchers and practitioners will face when keeping an LLM continuously up-to-date. This contrasts with simpler forms of knowledge, such as facts encoded as knowledge triples (e.g., “John Smith works at ACME Corporation”). A news article might be incomplete (e.g., it mentions key people participating in an ongoing event but does not define their roles, assuming the reader has been following the event for a while), or it might contradict other documents (e.g., “investment X is no longer recommended due to fraud scandals”).

Second, using the news domain we can ensure that the model has not been exposed to that specific information previously. Given the recency of the news, we can mitigate the contamination problem using news about events that occurred after the model’s training cutoff.

To address these two criteria, we used the TiEBE dataset [36]<sup>1</sup>, a dataset of news articles spanning from 2015 to 2024, along with a corresponding set of question-answer pairs, to evaluate models on their knowledge of specific events. The questions and answers were generated using GPT-4o-2024-08-06, by prompting it to produce four pairs per article in a single generation. More details on the creation of the data set can be found in Almeida et al. [36].

To run the knowledge injection experiments, we chose the Llama-2 model [51], which has an old enough knowledge cutoff – September 2022 – while having strong performance in the question-answering task due to being instruction-finetuned. This mitigates the risk of data contamination because the model was not exposed to the specific events covered in the documents during the pre-training.

However, since Llama-2 has a limit of 4k tokens in its context, we selected a subset of the *World* category filtering articles with up to 3,500 tokens to fit the article, the instruction template, and the generated answer within the model’s context length. Additionally, we selected only the recent documents, from 2023 to 2024, as the training corpus. The final dataset comprises 117 documents, each paired with four QA pairs – 468 QA pairs, in total.

For the automatic evaluation of the models’ answers, we followed the methodology described in Almeida et al., applying the process commonly known as LLM-as-a-judge [52],

[53]. This approach leverages the expected answers provided in the dataset, prompting an LLM to assess the correctness of candidate answers based on these true answers. This more sophisticated way of evaluating the answers is required to check whether the model has learned the complex facts contained in the documents, whereas strict approaches, such as exact matching, may fail to capture these nuances. The prompt and evaluation code we used are available in the released repository.

#### B. Control Datasets

One of the challenges in updating a model’s knowledge is ensuring it retains previously learned capabilities. Large models can easily memorize (i.e., fully reconstruct) a small set of documents, which is a key goal of the knowledge injection task. However, this often comes at the cost of significant performance drops on unrelated tasks where the original model previously excelled.

To quantify this forgetting gap, we use the average accuracy across the following seven datasets, collectively referred to as the *Control datasets*: OpenBookQA [54], ARC-Easy and ARC-Challenge [55], WinoGrande [56], HellaSwag [57], PIQA [58], and BoolQ [59]. These datasets are implemented in the Language Model Evaluation Harness [60], which has been used in prior work for similar evaluations [61]–[63].

#### C. Knowledge Injection Techniques

We investigate the following knowledge injection techniques:

**Retrieval-Augmented Generation (RAG):** One of the simplest and most straightforward methods for injecting knowledge into an LLM is to use a retrieval mechanism to locate relevant information within a corpus and include this information in the prompt provided to the LLM. Specifically, we employ BM25 [64], [65] to retrieve the top-N documents from the corpus of recent news, which consists of 117 documents.

To evaluate different configurations, we tested a document retrieval approach, where the best-matching document is prepended to the prompt, followed by the test question. Additionally, we evaluated a chunking-based approach, where documents are divided into chunks of 512 tokens with an overlap of 64 tokens. In this setup, we retrieved the top-5 chunks, which were then prepended to the question to allow the LLM to generate an answer.

**Continual Pre-training (CPT):** This approach involves continuing the pre-training of the LLM directly on the target document or article using the causal language modeling objective (i.e., next-token prediction) applied to the unmodified corpus [6].

**Rephrasing the Web (RTW):** Following the four prompts proposed by Maini et al. [66], we generate rephrased versions of a training example. Three of these prompts instruct the LLM to rephrase the input document in styles with varying levels of complexity: easy (simplified, suitable for a toddler), medium (clear and high-quality, similar to Wikipedia), and hard (terse

<sup>1</sup><https://huggingface.co/datasets/TimelyEventsBenchmark/TiEBE>

and abstruse). The fourth prompt asks the LLM to generate QA pairs that the document is likely to address.

We experiment with two models for generating these rephrases: (1) GPT-4o,<sup>2</sup> which provides a high-quality upper bound for rephrases but may represent an unrealistic setup since it could have been exposed to the content during its pre-training;<sup>3</sup> and (2) the model itself (i.e., Llama-2-7B), representing a more practical scenario in which the model does not have prior knowledge of the document being rephrased.

We report the results using all the four styles of prompts, using only the first three, and using only the QA-style prompt. This way we can keep only the rephrasing prompts and isolate the possibility of cross-contamination by generating similar questions as the ones used during the test phase.

**Instruction Pre-training (IPT):** We applied the instruction generation method proposed by Cheng et al. [31] to our training corpus. Using their instruction generation model,<sup>4</sup> we generated synthetic instructions for each training document in a 1-shot setting. A 1-shot approach was necessary, as using more examples would result in truncation due to exceeding the model’s context length.

**Paraphrasing (Para):** To evaluate the effect of simple paraphrasing on training examples, we adapted the prompt proposed by Ovadia et al. [15]. The prompt instructed the LLM to rephrase the content while maintaining factual accuracy and maintaining the original text length. We generated multiple paraphrases by applying token sampling with a specified temperature. For this process, we used two models: GPT-4o, providing high-quality paraphrases as an upper bound, and Llama-2-7B, representing a more realistic, self-contained approach.

#### D. Training Setup

The knowledge injection methods relying on continued pre-training were implemented by mixing synthetic augmented examples with the original documents. Therefore, all reported results, along with the corresponding variations, refer to the original 117 documents supplemented with an additional N synthetic variations in the training set.

Starting from the Llama-2-7B-chat checkpoint<sup>5</sup>, we further trained the model using the causal language modeling objective with the traditional cross-entropy loss. Training was conducted on batches of 8 examples, with a learning rate of  $5e-5$ , and the AdamW [67] optimizer. Given the task of injecting small amounts of knowledge, our training runs were intentionally short, often involving fewer than 15 training steps per epoch. To ensure stable training and reduce variance in results, we carefully tuned the hyperparameters in preliminary experiments, determining that training for two epochs with a relatively large learning rate warmup yielded the best performance.

<sup>2</sup>more specifically, GPT-4o-2024-08-06

<sup>3</sup><https://platform.openai.com/docs/models>

<sup>4</sup><https://huggingface.co/instruction-pretrain/instruction-synthesizer>

<sup>5</sup><https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

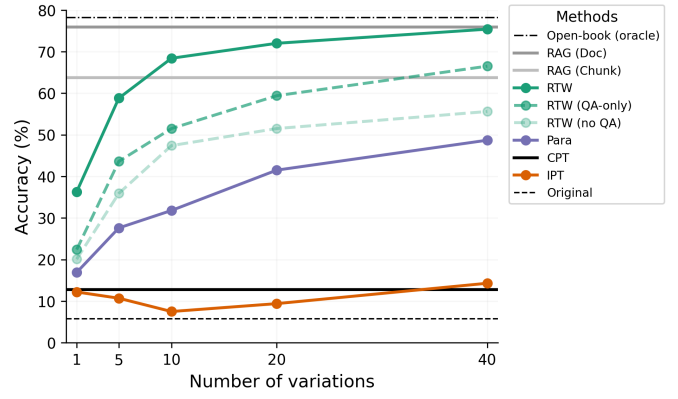


Fig. 1. Comparison of different knowledge injection methods, including parametric and non-parametric techniques. The upper bound is represented by open-book answering with access to the source document (oracle), while the lower bound corresponds to the model’s original performance in closed-book answering. Colored lines represent the knowledge injection methods evaluated in this study. The y-axis represents accuracy on the TIEBe dataset, considering only events from 2023 and 2024. The x-axis indicates the number of variations used for each augmentation method.

As described in [32], we applied a re-warmup and re-decay strategy to the learning rate, utilizing linear warmup and cosine decay. Specifically, the warmup phase was applied throughout the first epoch, while the decay phase occurred during the second and final epoch. This approach allowed us to re-warm the learning rate during the first half of training and re-decay it during the second half, optimizing the stability and effectiveness of the training process.

## IV. RESULTS

In this section, we describe and analyze the results of our experiments and conduct ablations comparing variations of the studied methods.

### A. Which method is the best?

Fig. 1 summarizes the comparison of the methods evaluated in this work. The dashed black lines show the lower and upper bounds, which correspond to the closed-book answering performance (when the model relies solely on its parametric knowledge) and the open-book answering performance (when the model has access to the context that answers the given question) using the original model.

Next, still using open-book answering but without providing the exact correct context (oracle), we evaluated RAG approaches. This is a more realistic scenario because in real-world applications the pairing of a question and its relevant context is not known a priori. For RAG evaluation, we used the top-1 most similar document according to the BM25 similarity score, as well as the top-5 most similar chunks. The results show very similar performance between the upper-bound oracle and RAG using the top-1 document. However, a significant drop is observed when using chunks, highlighting a known caveat of RAG systems and their sensitivity to chunking strategies.

The solid black line represents the CPT performance, serving as the baseline method for injecting knowledge from unstructured text. The other colored lines correspond to each augmentation method under comparison. For all methods, we repeated the synthetic generation  $N$  times (shown on the x-axis), using a sampling temperature of 1.0 to introduce variation in the generated examples.

For the RTW method, we report separate results for different prompt configurations: (i) using all four proposed prompts (easy, medium, hard, and QA-style), and (ii) two variations – one excluding the QA-style prompt and another using only the QA-style prompt [66]. This separation was implemented to measure the impact of the generated QA pairs and to mitigate the risk of indirect contamination, where the model might generate QA pairs similar to those encountered during testing, thus making the task artificially easier. Although we ensured that the QA pairs generated by RTW do not overlap with those in the test set, we report these results separately as a precaution.

The results indicate a monotonic increase in performance when using the RTW and Para methods, with RTW achieving a performance level comparable to document-level RAG, only a few points below the upper bound. The RTW variant that excludes the QA-style prompt and the Para method both yield similar performance, approaching that of chunk-based RAG. However, the superior performance of RTW suggests that leveraging multiple prompts to generate textual variations is more effective than relying on a single paraphrasing prompt.

The RTW variation that uses only QA-style prompts outperforms the configuration using only the other three prompts. This result indicates that synthetically generated QA pairs play an important role in increasing the model’s ability to answer questions based on learned information. However, concerns about indirect contamination remain, suggesting that the no-QA variant may be better suited for real-world applications.

The IPT method scored lower than expected, hovering around the CPT baseline. This might be due to using a one-shot approach instead of the few-shot approach that worked best in the original study – a strategy that was not feasible here because of the longer texts used, which would exceed the model’s context-length limit.

### B. Learning-forgetting tradeoff

To measure the possible catastrophic forgetting, we evaluated each checkpoint on the control dataset and averaged the accuracy across seven different tasks. This gives us a measure of how the model performs on tasks it previously knew.

Fig. 2 compared the control dataset accuracy with the total training tokens of each method on each number of generated variations (1, 5, 10, 20, and 40). The CPT baseline is indicated by a star, since its amount of training tokens is fixed, and we compared with the original performance and the RAG variants.

To evaluate the RAG performance, we prepended the retrieved context on each evaluated prompt, following the same prompt used on non-RAG evaluations. Thus, their evaluations are exactly the same as the other methods, the only difference

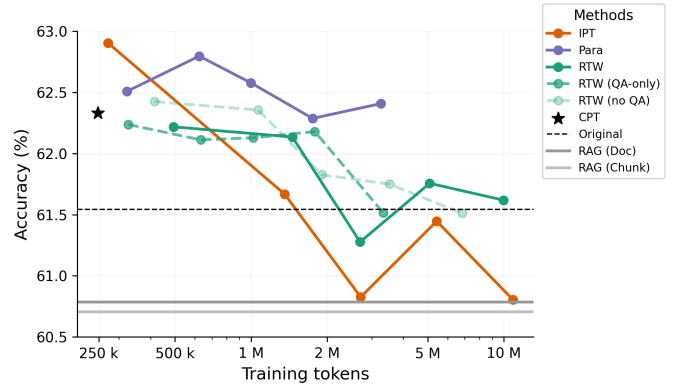


Fig. 2. Comparison of average accuracy across seven control datasets against the number of training tokens. Each point on the lines represents a different variation level (1, 5, 10, 20, 40). The results indicate a trend of accuracy degradation in the control sets as more variations are introduced during training, suggesting the onset of catastrophic forgetting of previously learned capabilities. This effect is also observed in RAG variants, which show performance degradation when exposed to in-domain data and evaluated on out-of-domain data. Note that the x-axis is on a logarithmic scale.

is that the retrieval result is included in the context. It is noteworthy that both RAG approaches lead to the highest degradation compared to the other methods, with the exception of IPT. This performance degradation highlights the caveats of using RAG-based approaches, where the retrieved context may *confuse* the model on out-of-domain tasks.

Surprisingly, all continued pre-training methods evaluated result in an increase in performance when training with a small number of tokens. This overall performance improvement was observed across all datasets except BoolQ and WinoGrande. One possible explanation for this result is that we start from an instruction-tuned model and evaluate it by computing log probabilities on multiple-choice tasks using the Language Model Evaluation Harness framework [60]. This evaluation approach may favor base models over instruction-tuned models due to probability calibration. Thus, this short continued pre-training on the next-token prediction task may resemble the original pre-training objective, leading to performance gains on some control datasets. However, we leave a deeper investigation of this phenomenon to future work.

The results indicate that, as the amount of training tokens increase, the average performance tends to drop. This result is consistent with the catastrophic forgetting phenomenon. This is not entirely true for the Para method, which oscillated on similar performances, but a more informed conclusion would require more tokens to check if the trend holds, since it is the method that resulted on the lowest amount of tokens.

Despite the difference in the in-domain accuracy, i.e., the accuracy on the recent news dataset, both RTW methods achieved similar performance on the control dataset, even though the dataset with all four types of prompt resulted on a slightly higher amount of generated tokens.

The IPT method exhibited the largest drop in performance, suggesting that incorporating synthetic instructions generated

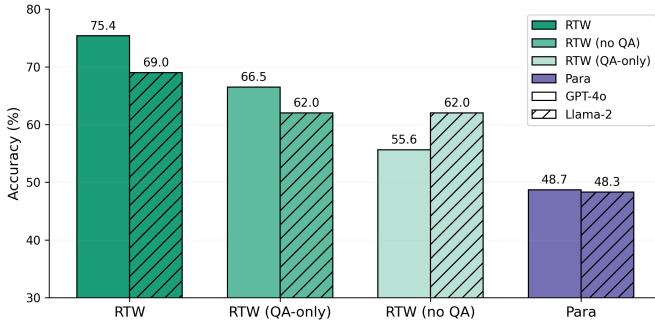


Fig. 3. Performance on the TiEBe test set of recent events using different augmentation methods that leverage external models to generate synthetic examples. We present variants of the RTW algorithm and Para, utilizing GPT-4o and Llama-2-7B. Specifically, we highlight that the model used for training can also generate synthetic data to augment its own training.

by its synthesizer accelerated the model’s forgetting of previously learned capabilities. Moreover, IPT showed low performance both in-domain and out-of-domain, indicating that the one-shot generation approach used in our experiments may be more detrimental than beneficial in our evaluated scenario.

### C. Can models augment themselves?

Previous results show that models can effectively learn new information by continuous pre-training on new data and benefits from synthetically augmented data. In this section, we investigate the role of the generator model used to augment the dataset. We used the RTW and Para techniques, since the IPT uses their specific synthesizer model, instead of a generic LLM.

Fig. 3 shows the comparison of the different generator models: GPT-4o and Llama-2-7B. For simply paraphrasing the content (Para), there was no significant difference of using a frontier model or the model itself to generate the synthetic training data. For the RTW, which uses varied prompts to augment the data, it is inconclusive whether one model performs better than the other because using all four prompts the GPT-4o lead to better results, but without the QA-style prompt, the model itself leads to better results.

One potential caveat is that GPT-4o was the same model used to generate the questions for the original TiEBe dataset [36]. Thus, an unwanted indirect contamination may explain it leading to better performance, since the model might have generated similar questions for the knowledge injection methods and for our test set. This hypothesis also explains why removing the QA prompt from RTW resulted in a large drop in performance when using GPT-4o as a generator. However, this drop is smaller when using LLaMA-2-7B as the generator.

Even though the results of the QA-only variant are higher with the GPT-4o generator, there is no difference when using the model itself as the generator. This highlights the performance achieved by the no-QA variant, which avoids the risk of exposing the model to similar QA pairs in the test, while showing the generalization of the learned information due to

solely training the model on rephrased versions of the original text.

Achieving comparable results with the smaller LLaMA-2-7B model and a state-of-the-art model is particularly noteworthy. It shows that the model can generate synthetic data to enhance its own capabilities, suggesting the possibility of continuous or iterative self-improvement through the ingestion of newly generated data.

## V. CONCLUSION

In this work, we investigated methods for injecting small-scale, unstructured knowledge into LLMs and examined the tradeoff between learning new facts and retaining prior knowledge. We found that simple continued pre-training yields modest improvements, while RAG can be effective but often degrades performance on unrelated tasks.

Our results highlight the importance of synthetic data augmentation: models trained on diverse rephrasings (e.g., RTW) learn new information more effectively while avoiding catastrophic forgetting. Notably, models can generate their own augmentation data, opening avenues for self-improving updates without external supervision.

Our findings emphasize the need for diverse training inputs to enhance knowledge acquisition while minimizing the degradation of previously learned information. We hope our released datasets and code will support future research on efficient knowledge injection.

## REFERENCES

- [1] F. Petroni, T. Rocktäschel, S. Riedel, P. Lewis, A. Bakhtin, Y. Wu, and A. Miller, “Language models as knowledge bases?” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 2463–2473.
- [2] B. Heinzlerling and K. Inui, “Language models as knowledge bases: On entity representations, storage capacity, and paraphrased queries,” in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, P. Merlo, J. Tiedemann, and R. Tsarfaty, Eds. Online: Association for Computational Linguistics, Apr. 2021, pp. 1772–1791.
- [3] C. Wang, P. Liu, and Y. Zhang, “Can generative pre-trained language models serve as knowledge bases for closed-book QA?” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, C. Zong, F. Xia, W. Li, and R. Navigli, Eds. Online: Association for Computational Linguistics, Aug. 2021, pp. 3241–3251.
- [4] Z. Zhong, D. Friedman, and D. Chen, “Factual probing is [MASK]: Learning vs. learning to recall,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tur, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, Eds. Online: Association for Computational Linguistics, Jun. 2021, pp. 5017–5033.
- [5] P. Youssef, O. Koraş, M. Li, J. Schlöterer, and C. Seifert, “Give me the facts! a survey on factual knowledge probing in pre-trained language models,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 15 588–15 605.

- [6] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith, "Don't stop pretraining: Adapt language models to domains and tasks," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds. Online: Association for Computational Linguistics, Jul. 2020, pp. 8342–8360. [Online]. Available: <https://aclanthology.org/2020.acl-main.740/>
- [7] T. Wu, L. Luo, Y.-F. Li, S. Pan, T.-T. Vu, and G. Haffari, "Continual learning for large language models: A survey," *arXiv preprint arXiv:2402.01364*, 2024.
- [8] A. Lewkowycz, A. J. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, Y. Wu, B. Neyshabur, G. Gur-Ari, and V. Misra, "Solving quantitative reasoning problems with language models," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022.
- [9] Z. Azerbayev, H. Schoelkopf, K. Paster, M. D. Santos, S. M. McAleer, A. Q. Jiang, J. Deng, S. Biderman, and S. Welleck, "Llemma: An open language model for mathematics," in *The Twelfth International Conference on Learning Representations*, 2024.
- [10] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl *et al.*, "Large language models encode clinical knowledge," *Nature*, vol. 620, no. 7972, pp. 172–180, 2023.
- [11] Q. Xie, Q. Chen, A. Chen, C. Peng, Y. Hu, F. Lin, X. Peng, J. Huang, J. Zhang, V. Kelothe, H. He, L. Ohno-Machido, Y. Wu, H. Xu, and J. Bian, "Me llama: Foundation large language models for medical applications," 2024.
- [12] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. C. Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve, "Code llama: Open foundation models for code," 2024. [Online]. Available: <https://arxiv.org/abs/2308.12950>
- [13] M. Ciosici, J. Cecil, D.-H. Lee, A. Hedges, M. Freedman, and R. Weischedel, "Perhaps PTLMs should go to school – a task to assess open book and closed book QA," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 6104–6111.
- [14] Z. Ke, Y. Shao, H. Lin, T. Konishi, G. Kim, and B. Liu, "Continual pre-training of language models," in *The Eleventh International Conference on Learning Representations*, 2023.
- [15] O. Ovadia, M. Brief, M. Mishaali, and O. Elisha, "Fine-tuning or retrieval? comparing knowledge injection in llms," *arXiv preprint arXiv:2312.05934*, 2023.
- [16] D. Cheng, S. Huang, and F. Wei, "Adapting large language models via reading comprehension," in *The Twelfth International Conference on Learning Representations*, 2024.
- [17] D. Biderman, J. Portes, J. J. G. Ortiz, M. Paul, P. Greengard, C. Jennings, D. King, S. Havens, V. Chiley, J. Frankle, C. Blakeney, and J. P. Cunningham, "LoRA learns less and forgets less," *Transactions on Machine Learning Research*, 2024, featured Certification.
- [18] Z. Gekhman, G. Yona, R. Aharoni, M. Eyal, A. Feder, R. Reichart, and J. Hertz, "Does fine-tuning LLMs on new knowledge encourage hallucinations?" in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds. Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 7765–7784. [Online]. Available: <https://aclanthology.org/2024.emnlp-main.444/>
- [19] L. Berglund, M. Tong, M. Kaufmann, M. Balesni, A. C. Stickland, T. Korbak, and O. Evans, "The reversal curse: LLMs trained on "a is b" fail to learn "b is a"," in *The Twelfth International Conference on Learning Representations*, 2024.
- [20] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [21] M. Riemer, I. Cases, R. Ajemian, M. Liu, I. Rish, Y. Tu, , and G. Tesaro, "Learning to learn without forgetting by maximizing transfer and minimizing interference," in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=B1gTShAc7>
- [22] A. Kleiman, J. Frankle, S. M. Kakade, and M. Paul, "Predicting task forgetting in large language models," in *ICML 2023 Workshop on Deployable Generative AI*, 2023. [Online]. Available: <https://openreview.net/pdf?id=0BMg0OgNTP>
- [23] Y. Luo, Z. Yang, F. Meng, Y. Li, J. Zhou, and Y. Zhang, "An empirical study of catastrophic forgetting in large language models during continual fine-tuning," *ArXiv*, vol. abs/2308.08747, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:261031244>
- [24] S. Kotha, J. M. Springer, and A. Raghunathan, "Understanding catastrophic forgetting in language models via implicit inference," in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=VrHiF2hsrm>
- [25] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," ser. *Psychology of Learning and Motivation*, G. H. Bower, Ed. Academic Press, 1989, vol. 24, pp. 109–165. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0079742108605368>
- [26] R. M. French, "Catastrophic forgetting in connectionist networks," *Trends in Cognitive Sciences*, vol. 3, no. 4, pp. 128–135, 1999. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1364661399012942>
- [27] I. J. Goodfellow, M. Mirza, X. Da, A. C. Courville, and Y. Bengio, "An empirical investigation of catastrophic forgetting in gradient-based neural networks," in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2014. [Online]. Available: <http://arxiv.org/abs/1312.6211>
- [28] C. Zhu, A. S. Rawat, M. Zaheer, S. Bhojanapalli, D. Li, F. Yu, and S. Kumar, "Modifying memories in transformer models," 2020.
- [29] K. Meng, D. Bau, A. J. Andonian, and Y. Belinkov, "Locating and editing factual associations in gpt," in *Advances in Neural Information Processing Systems*, 2022.
- [30] K. Meng, A. S. Sharma, A. J. Andonian, Y. Belinkov, and D. Bau, "Mass-editing memory in a transformer," in *The Eleventh International Conference on Learning Representations*, 2023.
- [31] D. Cheng, Y. Gu, S. Huang, J. Bi, M. Huang, and F. Wei, "Instruction pre-training: Language models are supervised multitask learners," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Y. Al-Onaizan, M. Bansal, and Y.-N. Chen, Eds. Miami, Florida, USA: Association for Computational Linguistics, Nov. 2024, pp. 2529–2550. [Online]. Available: <https://aclanthology.org/2024.emnlp-main.148/>
- [32] A. Ibrahim, B. Thérien, K. Gupta, M. L. Richter, Q. G. Anthony, E. Belilovsky, T. Lesort, and I. Rish, "Simple and scalable strategies to continually pre-train large language models," *Transactions on Machine Learning Research*, 2024. [Online]. Available: <https://openreview.net/forum?id=DimPeeCxKO>
- [33] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive nlp tasks," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 9459–9474.
- [34] A. Balaguer, V. Benara, R. L. de Freitas Cunha, R. de M. Estevão Filho, T. Hendry, D. Holstein, J. Marsman, N. Mecklenburg, S. Malvar, L. O. Nunes, R. Padilha, M. Sharp, B. Silva, S. Sharma, V. Aski, and R. Chandra, "Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture," 2024.
- [35] N. Mecklenburg, Y. Lin, X. Li, D. Holstein, L. Nunes, S. Malvar, B. Silva, R. Chandra, V. Aski, P. K. R. Yannam, T. Aktas, and T. Hendry, "Injecting new knowledge into large language models via supervised fine-tuning," 2024.
- [36] T. S. Almeida, G. K. Bonás, J. G. A. Santos, H. Abonizio, and R. Nogueira, "Tiebe: A benchmark for assessing the current knowledge of large language models," 2025.
- [37] R. Eldan and M. Russinovich, "Who's harry potter? approximate unlearning in llms," 2023.
- [38] N. Zhang, Y. Yao, B. Tian, P. Wang, S. Deng, M. Wang, Z. Xi, S. Mao, J. Zhang, Y. Ni, S. Cheng, Z. Xu, X. Xu, J.-C. Gu, Y. Jiang, P. Xie, F. Huang, L. Liang, Z. Zhang, X. Zhu, J. Zhou, and H. Chen, "A comprehensive study of knowledge editing for large language models," 2024.

- [39] T. Hartvigsen, S. Sankaranarayanan, H. Palangi, Y. Kim, and M. Ghassemi, "Aging with GRACE: Lifelong model editing with discrete key-value adaptors," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [40] O. Levy, M. Seo, E. Choi, and L. Zettlemoyer, "Zero-shot relation extraction via reading comprehension," in *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, 2017, pp. 333–342.
- [41] R. Cohen, E. Biran, O. Yoran, A. Globerson, and M. Geva, "Evaluating the ripple effects of knowledge editing in language models," *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 283–298, 2024.
- [42] S. Wu, M. Peng, Y. Chen, J. Su, and M. Sun, "Eva-kellm: A new benchmark for evaluating knowledge editing of llms," 2023.
- [43] J. Zhang, W. Cui, Y. Huang, K. Das, and S. Kumar, "Synthetic knowledge ingestion: Towards knowledge refinement and injection for enhancing large language models," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 21456–21473.
- [44] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, and H. Wang, "Retrieval-augmented generation for large language models: A survey," 2024.
- [45] P. Finardi, L. Avila, R. Castaldoni, P. Gengo, C. Larcher, M. Piau, P. Costa, and V. Caridà, "The chronicles of rag: The retriever, the chunk and the generator," 2024.
- [46] Z. Zhong, H. Liu, X. Cui, X. Zhang, and Z. Qin, "Mix-of-granularity: Optimize the chunking granularity for retrieval-augmented generation," 2024.
- [47] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang, "Lost in the middle: How language models use long contexts," *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 157–173, 2024.
- [48] J. Wei, M. Bosma, V. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," in *International Conference on Learning Representations*, 2022.
- [49] E. Wu, K. Wu, and J. Zou, "Finetunebench: How well do commercial fine-tuning apis infuse knowledge into llms?" 2024.
- [50] Z. Yang, N. Band, S. Li, E. Candès, and T. Hashimoto, "Synthetic continued pretraining," 2024.
- [51] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom, "Llama 2: Open foundation and fine-tuned chat models," 2023.
- [52] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica, "Judging LLM-as-a-judge with MT-bench and chatbot arena," in *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [53] J. Gu, X. Jiang, Z. Shi, H. Tan, X. Zhai, C. Xu, W. Li, Y. Shen, S. Ma, H. Liu, Y. Wang, and J. Guo, "A survey on llm-as-a-judge," 2025.
- [54] T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal, "Can a suit of armor conduct electricity? a new dataset for open book question answering," in *EMNLP*, 2018.
- [55] P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord, "Think you have solved question answering? try ARC, the AI2 reasoning challenge," *arXiv:1803.05457v1*, 2018.
- [56] K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi, "WinoGrande: An adversarial winograd schema challenge at scale," *arXiv preprint arXiv:1907.10641*, 2019.
- [57] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi, "HellaSwag: Can a machine really finish your sentence?" in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, A. Korhonen, D. Traum, and L. Marquez, Eds. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 4791–4800.
- [58] Y. Bisk, R. Zellers, R. L. Bras, J. Gao, and Y. Choi, "PIQA: Reasoning about physical commonsense in natural language," in *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [59] C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, and K. Toutanova, "BoolQ: Exploring the surprising difficulty of natural yes/no questions," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 2924–2936.
- [60] L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu, A. Le Noac'h, H. Li, K. McDonell, N. Muenighoff, C. Ociepa, J. Phang, L. Reynolds, H. Schoelkopf, A. Skowron, L. Sutawika, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou, "A framework for few-shot language model evaluation," 07 2024.
- [61] X. Ma, G. Fang, and X. Wang, "LLM-pruner: On the structural pruning of large language models," in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [62] P. Zhang, G. Zeng, T. Wang, and W. Lu, "Tinyllama: An open-source small language model," 2024.
- [63] R. Eldan and M. Russinovich, "Who's harry potter? approximate unlearning for LLMs," 2024.
- [64] S. E. Robertson, S. Walker, S. Jones, M. Hancock-Beaulieu, and M. Gatford, "Okapi at trec-3," in *Text Retrieval Conference*, 1994.
- [65] G. M. Rosa, R. C. Rodrigues, R. Lotufo, and R. Nogueira, "Yes, bm25 is a strong baseline for legal case retrieval," 2021.
- [66] P. Maini, S. Seto, R. Bai, D. Grangier, Y. Zhang, and N. Jaitly, "Rephrasing the web: A recipe for compute and data-efficient language modeling," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 14044–14072.
- [67] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *International Conference on Learning Representations*, 2019.