

SELF-TUNING: Instructing LLMs to Effectively Acquire New Knowledge through Self-Teaching

Anonymous ACL submission

Abstract

Large language models (LLMs) often struggle to provide up-to-date information due to their one-time training and the constantly evolving nature of the world. To keep LLMs current, existing approaches typically involve continued pre-training on new documents. However, they frequently face difficulties in extracting stored knowledge. Motivated by the remarkable success of the Feynman Technique in efficient human learning, we introduce SELF-TUNING, a learning framework aimed at improving an LLM’s ability to effectively acquire new knowledge from unseen raw documents through self-teaching. Specifically, we develop a SELF-TEACHING strategy that augments the documents with a set of knowledge-intensive tasks created in a self-supervised manner, focusing on three crucial aspects: *memorization*, *comprehension*, and *self-reflection*. Additionally, we introduce three Wiki-Newpages-2023-QA datasets to facilitate an in-depth analysis of an LLM’s knowledge acquisition ability concerning *memorization*, *extraction*, and *reasoning*. Extensive experimental results on various models, *e.g.*, LLAMA2-7B reveal that SELF-TUNING consistently exhibits superior performance across all knowledge acquisition tasks and excels in preserving previous knowledge.

1 Introduction

Armed with a wealth of factual knowledge acquired during the pre-training phase (Zhou et al., 2023a), LLMs (Touvron et al., 2023a; OpenAI, 2023) exhibit remarkable proficiency in numerous knowledge-intensive tasks (Cohen et al., 2023; Gekhman et al., 2024). Despite this, the knowledge stored in LLMs can quickly become outdated due to the one-time training of LLMs and the ever-changing nature of the world (Huang et al., 2023; Jiang et al., 2024c). These unavoidable knowledge limitations present notable obstacles to the trustworthiness of LLMs in real-world scenarios (Liu

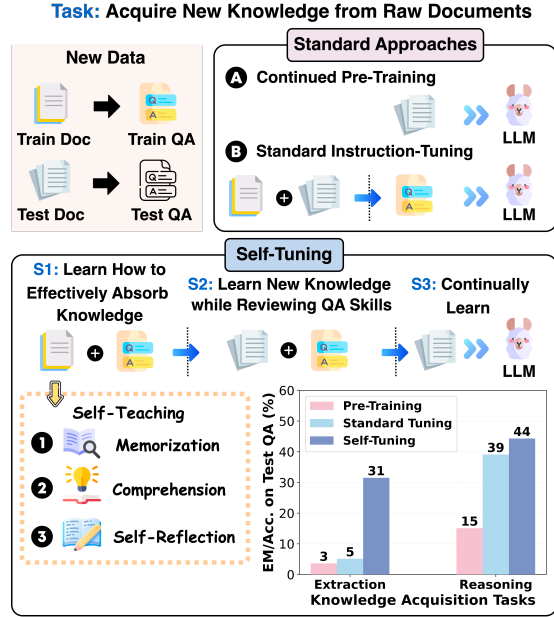


Figure 1: Illustration of the knowledge acquisition task with two standard knowledge injection approaches (in the upper part). Depiction of SELF-TUNING for effective knowledge acquisition from unseen raw documents, which significantly enhances factual accuracy compared to the standard approaches (in the lower part).

et al., 2023; Mecklenburg et al., 2024). Thus, it is essential to equip LLMs with new knowledge to keep them up-to-date.

In this paper, we focus on injecting new knowledge into the parameters of LLMs. As depicted in the upper part of Figure 1, a standard approach involves continued pre-training (A) on a raw corpus (here, test doc) containing new information (Jang et al., 2022). However, it struggles to extract the embedded knowledge, potentially due to the impaired question-answering (QA) capability (Allen-Zhu and Li, 2023; Cheng et al., 2024). Despite the assistance of subsequent instruction-tuning (B) (Wei et al., 2022; Ouyang et al., 2022a) on QA data, the knowledge retrieved from the LLMs remains notably constrained (Jiang et al., 2024c).

Recently, Jiang et al. (2024c) suggests fine-tuning on a mix of QA data and related documents before continuing pre-training, with the aim of teaching the model how to access knowledge from documents and answer questions. Although this method greatly outperforms standard approaches, our initial results suggest that its effectiveness in knowledge extraction remains limited.

Numerous studies (Ambion et al., 2020; Reyes et al., 2021) evidence the effectiveness of the Feynman Technique (Xiaofei et al., 2017) in promoting human learning and knowledge understanding. The remarkable success of this potent learning method is often attributed to its emphasis on “comprehension,” “self-reflection” (“identifying gaps and review”), rather than mere “memorization”. This encourages our exploration into its potential application in improving LLMs’ knowledge acquisition capabilities. As a result, we present SELF-TUNING, a framework that empowers an LLM to effectively internalize and recall new knowledge. As depicted in the lower part of Figure 1, SELF-TUNING consists of three stages: (i) Firstly, we train the model using a mix of training documents and associated QA data, equipping it with the ability to efficiently absorb knowledge from raw documents via self-teaching, as well as question-answering skills. Specifically, we design a SELF-TEACHING strategy to present the training documents as plain texts for *memorization* and a series of knowledge-intensive tasks derived from the documents in a self-supervised manner, without any mining patterns (van de Kar et al., 2022), for *comprehension* and *self-reflection*. (ii) Next, we deploy the model to apply the learning strategy for spontaneously acquiring knowledge from new documents while reviewing its QA skills. (iii) Finally, we continue training the model using only the new documents to ensure thorough acquisition of new knowledge.

In addition, we introduce three Wiki-Newpages-2023-QA datasets to conduct an in-depth study of how an LLM acquires new knowledge *w.r.t.*, *memorization*, *extraction*, and *comprehension* (in this study, *reasoning*) across single-domain, multi-domain, and cross-domain settings. These datasets are carefully curated to ensure minimal overlap with the LLM’s pre-training corpora, emphasizing two key knowledge-intensive tasks, *i.e.*, open-ended generation and natural language inference (NLI) tasks. Extensive experimental results on diverse models, *e.g.*, LLAMA2-7B (Touvron et al., 2023b), Qwen2-7B (Yang et al., 2024), and Mistral-

7B-v0.1 (Jiang et al., 2023) demonstrate that SELF-TUNING significantly outperforms all other compared methods on knowledge memorization and extraction tasks. In addition, SELF-TUNING consistently yields high accuracy on reasoning tasks, while the performance of the compared methods largely fluctuates in different scenarios. Inspiringly, SELF-TUNING exhibits exceptional performance in retaining previously acquired knowledge (*i.e.*, knowledge retention) concerning extraction and reasoning on two well-established benchmarks.

In summary, our contributions are three-fold:

- We present SELF-TUNING, a framework designed to improve an LLM’s knowledge acquisition capability via self-teaching.
- We introduce three Wiki-Newpages-2023-QA datasets to enable a comprehensive analysis of an LLM’s knowledge acquisition ability *w.r.t.*, memorization, extraction, and reasoning.
- We validate the efficacy of SELF-TUNING on three crucial knowledge acquisition tasks using the Wiki-Newpages-2023-QA datasets.

2 Related Work

Continual Knowledge Injection. The primary research approach for injecting new knowledge into LLMs (Xu et al., 2023; Ovadia et al., 2024; Mecklenburg et al., 2024) is through continued pre-training. This method entails the ongoing pre-training of LLMs on raw corpora containing new knowledge, carried out in a causal auto-regressive manner (Allen-Zhu and Li, 2023; Ibrahim et al., 2024; Ovadia et al., 2024). However, this straightforward approach often encounters hurdles in effectively enabling LLMs to extract the acquired knowledge during the inference phase (Allen-Zhu and Li, 2023; Jiang et al., 2024c; Cheng et al., 2024). To enhance knowledge extraction, instruction tuning on QA data after pre-training has been extensively employed (Wei et al., 2022; Ouyang et al., 2022b). Jiang et al. (2024c) suggests that the effectiveness of this method remains limited, and proposes fine-tuning the model on QA data before continued pre-training. This instructs the model on how to retrieve knowledge from raw corpora, thereby enhancing knowledge extraction. However, such an approach tends to underestimate the importance of comprehending the new knowledge.

Acknowledging the value of knowledge comprehension, Cheng et al. (2024) proposes converting raw corpora into reading comprehension texts.

Wiki-Newpages	Factual Knowledge	Open-Ended Generation (Train & Test Sets)		NLI (Test Set)	
		Statistics	Avg. # Tokens	Statistics	Answer Type
Wiki-Bio (Single-domain)	Birth Date, Profession, Education, <i>etc.</i>	Train: 6,136 (# QA); 1,136 (# Docs) Test: 663 (# QA); 127 (# Docs)	8.34 (Q) 4.24 (A) 59.64 (Doc)	729 (# QA) 127 (# Docs)	Yes (65.84%) No (33.47%) Impossible (0.69%)
Wiki-Multi (Multi-domain)	News, TV Series, Sports, <i>etc.</i>	Train: 10,004 (# QA); 1,823 (# Docs) Test: 1,502 (# QA); 281 (# Docs)	10.13 (Q) 5.70 (A) 69.25 (Doc)	1,627 (# QA) 281 (# Docs)	Yes (60.97%) No (36.63%) Impossible (2.40%)
Wiki-Film (Single-domain)	Genre, Language, Director, Released Time, <i>etc.</i>	Test: 955 (# QA); 169 (# Docs)	8.83 (Q) 4.61 (A) 58.10 (Doc)	1,387 (# QA) 169 (# Docs)	Yes (62.73%) No (26.53%) Impossible (10.74%)

Table 1: Statistical information of three Wiki-Newpages-2023-QA datasets, *i.e.*, Wiki-Bio, Wiki-Multi, and Wiki-Film. “Impossible”: “It’s impossible to say”. Details about token count distribution can be found in Appendix P.

This approach, however, focuses on domain adaptation and preserving general prompting abilities by mining a set of instruction-following tasks from the document content. In contrast, our work aims to equip the model with the ability to effectively absorb new knowledge from raw documents and employ the learned ability to unseen documents. Specifically, we develop a SELF-TEACHING strategy to present the raw document as plain texts for memorization, accompanied by a set of tasks for comprehension and self-reflection, which are created based on raw corpora in a self-supervised manner, without relying on any mining patterns.

Additionally, **knowledge editing** (Zhang et al., 2024a) and **retrieval-augmented generation** (Ovadia et al., 2024; Jeong et al., 2024) are recognized as two related research fields (Appendix A).

3 Wiki-Newpages-2023-QA: Datasets for Studying LLM Knowledge Acquisition

To explore the knowledge acquisition capabilities of LLMs from new documents, *w.r.t.*, memorization, extraction and reasoning, we introduce the Wiki-Newpages-2023-QA datasets (Table 1), which are carefully designed to minimize overlap with the initial pre-training corpus. These datasets comprise new document corpora for studying knowledge memorization and associated QA datasets for two vital knowledge-intensive tasks: open-ended generation and NLI for examining extraction and reasoning, respectively. Due to space constraints, we provide a brief overview of the dataset construction process here, with the complete version available in Appendix B.

3.1 Document Collection and QA Pair Generation

Document Collection. To construct the document corpus, we collect articles from September to October 2023 (4,257 articles in total) from

Wikipedia NewPages¹, which include new articles from various domains published after the pre-training cut-off time of the LLMs being evaluated.² Following Jiang et al. (2024c), we only use the first paragraph of each article, as it offers a comprehensive summary and contains a wealth of factual information.

QA Pair Generation. We gather QA pairs for generation and NLI tasks using our handcrafted prompts in Tables 23 and 24, aiming to cover all factual information within the given document.

3.2 Splitting

We construct three datasets for single-domain, multi-domain, and cross-domain analysis, splitting them into training and testing subsets while *ensuring zero factual knowledge overlap*.

Dataset Splitting. We generate three datasets: Wiki-Newpages-2023-10-Bio (Wiki-Bio), Wiki-Newpages-2023-10-Multi (Wiki-Multi), and Wiki-Newpages-2023-(9)10-Film (Wiki-Film) by randomly selecting 1,263 biographical documents, 2,104 multi-domain documents, and 955 film documents from the collected document corpus and their associated QA pairs.

Train-test Splitting. We divide Wiki-Bio and Wiki-Multi datasets into training and testing subsets for single-domain and multi-domain evaluations. We use Wiki-Film as the test set for cross-domain scenarios. Note that the training QA datasets only include open-ended generation task pairs, ensuring fair comparisons.

4 SELF-TUNING

In this section, we introduce the SELF-TUNING framework to improve the LLM’s capability to ac-

¹<https://en.wikipedia.org/wiki/Special:NewPages>

²The pre-training cut-off time for the LLAMA2 family models used in this study is 2022.

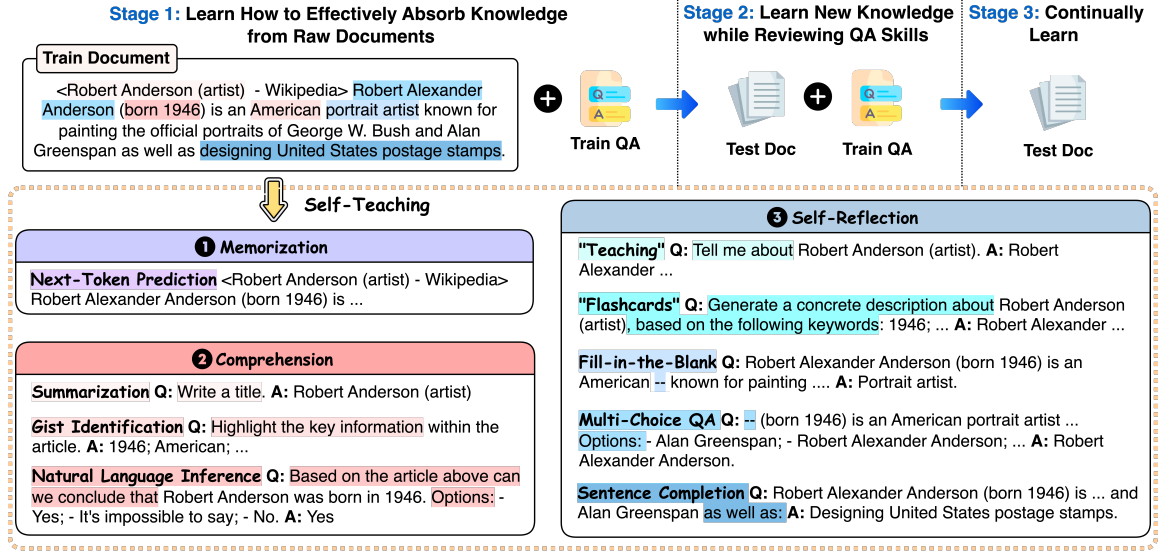


Figure 2: Illustration of the proposed SELF-TUNING. The framework consists of three stages (in the upper part): (i) Equipping the model with the ability to deeply absorb knowledge from raw documents using the proposed SELF-TEACHING strategy (in the lower part), along with question-answering capabilities; (ii) Applying the learning strategy acquired in Stage 1 to obtain new knowledge from unseen documents and refining QA skills; (ii) Continuously learning from unseen documents. See Appendix R for the full training document example in Stage 1.

quire knowledge from new documents, with the devised SELF-TEACHING strategy. We first give an overview of the training process for knowledge acquisition using the proposed SELF-TUNING in Section 4.1. Then, we delve into the SELF-TEACHING strategy in Section 4.2.

4.1 Overview

As depicted in Figure 2, the proposed SELF-TUNING comprises the following three stages.

Stage 1: Learn How to Effectively Absorb Knowledge from Raw Documents. Our objective is to equip an LLM M , parameterized by θ , with the ability to learn how to derive knowledge from raw documents. This is achieved by training the model using a combination of training document dataset D_{train}^{Doc} and associated training QA dataset D_{train}^{QA} , as depicted in the upper left part of Figure 2. To enhance effective knowledge absorption, we present D_{train}^{Doc} along with a series of knowledge-intensive tasks (*a.k.a.* self-teaching tasks) D_{train}^{Self} that are related to their content using the proposed SELF-TEACHING strategy (in the lower part of Figure 2). These tasks are generated in a self-supervised manner based on the contents of D_{train}^{Doc} , using the proposed SELF-TEACHING learning approach (Section 4.2). The multi-task training objective is:

$$L_{\theta}^{Stage1} = L_{\theta}(D_{train}^{Doc}) + L_{\theta}(D_{train}^{Self}) + L_{\theta}(D_{train}^{QA}) \quad (1)$$

Stage 2: Learn New Knowledge while Reviewing QA Skills. Our aim is to train the model M to apply the learned strategy for spontaneously extracting new knowledge from unseen documents (*i.e.*, raw test corpora D_{test}^{Doc}). In addition to training on D_{test}^{Doc} , we include D_{train}^{QA} , allowing the model M to review and refine its question-answering ability. The objective of this stage is:

$$L_{\theta}^{Stage2} = L_{\theta}(D_{test}^{Doc}) + L_{\theta}(D_{train}^{QA}) \quad (2)$$

Stage 3: Continually Learn. Our goal is to ensure that the model M thoroughly absorbs the new knowledge by conducting follow-up training on D_{test}^{Doc} (raw corpora). The objective is as follows:

$$L_{\theta}^{Stage3} = L_{\theta}(D_{test}^{Doc}) \quad (3)$$

4.2 SELF-TEACHING Learning Strategy

Motivated by the Feynman Technique, we aim to equip the model with systematic knowledge learning abilities from three perspectives: memorization, comprehension, and self-reflection, as shown in the lower part of Figure 2. Specifically, we devise a self-supervised SELF-TEACHING learning strategy that presents the raw documents D_{train}^{Doc} as plain texts for memorization and as a series of knowledge-intensive tasks in a question-answering format related to their content for comprehension and self-reflection (Table 21). This method *does not require any specific mining patterns, making it applicable to any raw texts.*

Memorization. To allow the model M to learn to memorize and capitalize on the factual information embedded in the raw texts, we execute the *next-token prediction* task on plain document texts.

Comprehension. Our goal is to facilitate the model’s ability to comprehend the factual knowledge within the document in a top-down manner. To achieve this, we conduct the following tasks:

(i) *Summarization* allows the model to learn to grasp the topic by using the prompt `Write a title:` to encourage the model to summarize the raw text, with the document title serving as the ground truth.

(ii) *Gist identification* improves the model’s ability to pinpoint the key elements (*i.e.*, entities) within the atomic facts. Specifically, we prompt the model with `Highlight the key information within the article:`, and use the entities within the document as gold answers, identified using Spacy³.

(iii) *Natural language inference* provides the model with the capability to determine whether a statement can be inferred from specific document contents (*i.e.*, “Yes,” “No,” or “It’s impossible to say”), thus avoiding misconceptions that may arise during knowledge acquisition. Specifically, we use a randomly sampled sentence (identified using NLTK⁴) within the document content as the true statement, and a corrupted version where one entity is replaced by an irrelevant entity from another sentence as the false statement. Then, we prompt the model with `Based on the article above can we conclude that` and the sampled sentence (either initial or corrupted), with the three relations as options and corresponding answers.

Self-Reflection. Our objective is to improve the model’s ability to memorize and recall acquired knowledge by “identifying and filling in the knowledge gaps.” To this end, we devise the following closed-book generation tasks:

(i) *“Teaching”* fosters the model’s ability to recall its acquired knowledge on a particular topic by “pretending to teach” others, using the prompt `Tell me about {topic}:` with the document content serving as the answer.

(ii) *“Flashcards”* imparts the model with the ability to recall its learned information based on the topic and associated keywords, using the prompt `Generate a concrete description`

about {topic} based on the following keywords:, with the document text as the answer.

(iii) *Fill-in-the-Blank* equips the model with the ability to conduct a detailed check on the acquired factual information. Specifically, we randomly replace one entity with a “-” symbol to form a cloze question, with the replaced entity serving as the corresponding answer.

(iv) *Multi-choice QA* helps the model learn to differentiate the correct answer from the available options and prevents confusion with irrelevant content. Specifically, we randomly replace one entity with a “-” symbol to form a cloze question, with the replaced entity and three other entities randomly sampled from the document forming the options, and the replaced entity serving as the correct choice.

(v) *Sentence completion* allows the model to develop its ability to focus on factual data found towards the end of a sentence. This is crucial since our initial observations indicate that the model frequently encounters difficulties when attempting to extract knowledge from later positions. Additionally, the model is anticipated to learn to emphasize not only entities but also phrase-level factual information. To achieve this, we first employ Spacy to pinpoint prepositions in a randomly chosen sentence from the document. Then, we store the phrase that follows the final preposition as the correct answer and the portion of the sentence preceding the phrase as the question. Comprehensive templates for each task can be found in Table 21.

5 Experiments

5.1 Setup

Datasets and Evaluation Metrics. We validate SELF-TUNING in both knowledge acquisition and retention for a well-rounded analysis.

We perform assessments on three **knowledge acquisition** tasks: (i) *Memorization*: We use test document datasets and report perplexity (PPL) (Jelinek et al., 1977). (ii) *Extraction*: We use test QA datasets for open-ended generation tasks and evaluate factual accuracy using exact match (EM), Recall, F1 (Kwiatkowski et al., 2019), Rouge-L (Lin, 2004), and accuracy. (iii) *Reasoning*: We use test QA datasets for NLI tasks and report accuracy.

We evaluate two aspects of **knowledge retention**: (i) *Extraction*: We assess the model’s performance in retaining general factual knowledge using Natural Questions (NQ) (Kwiatkowski et al., 2019), with EM and F1. (ii) *Reasoning*: We evaluate the

³<https://spacy.io/usage>

⁴A natural language toolkit. <https://www.nltk.org/>

Method	Training Data in Each Stage		
	Stage 1	Stage 2	Stage 3
Continued Pre-Training			① test doc
Standard Instruction-Tuning	① train doc & test doc		② train QA
PIT	① train QA train doc		② test doc
SELF-TUNING	① train QA & train doc w/ self-teaching tasks	② train QA & test doc	③ test doc
Variants of SELF-TUNING			
SELF-TUNING w/o Review	① train QA & train doc w/ self-teaching tasks		② test doc
SELF-TUNING via Read.	① train QA & train doc (reading-comprehension format (Cheng et al., 2024))		② test doc
SELF-TUNING w/ Pre-Review	① train QA & train doc w/ self-teaching tasks	② train QA & train doc	③ test doc

Table 2: Depiction of the training stages and datasets used in the compared methods. All approaches train on test documents for the same number of epochs. See Table 7 for the complete version.

Method	Wiki-Newpages-2023-QA (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Memorization	Extraction					Reason.	Extraction		Reasoning
	PPL (↓)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
Knowledge Acquisition on Wiki-Newpages-2023-10-Bio (Single-Domain Scenario)										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	8.41	55.20	31.83	64.48	75.55	62.10	7.96	-	-	-
Closed-book	8.41	4.68	2.87	14.63	16.98	15.07	7.96	16.05	24.67	53.40
<i>w/ Knowledge Injection</i>										
Cont. Pre-training	7.28	6.33	3.62	15.96	18.72	16.11	15.09	16.00	24.11	53.40
Standard Ins.-tuning	6.83	6.94	5.13	19.15	19.05	19.48	39.09	15.72	23.67	51.84
PIT	2.08	14.03	11.61	27.15	28.86	27.11	11.93	15.72	26.31	57.58
SELF-TUNING	1.11	37.25	31.52	50.83	52.62	50.61	44.31	16.45	25.67	66.01
Knowledge Acquisition on Wiki-Newpages-2023-10-Multi (Multi-Domain Scenario)										
Open-book w/ test doc	7.84	48.93	26.63	60.37	71.71	58.54	6.33	-	-	-
Closed-book	7.84	4.53	2.73	16.19	18.63	16.38	6.33	16.05	24.67	53.40
Cont. Pre-training	3.32	5.86	3.40	18.04	20.59	18.42	14.51	17.02	25.05	53.56
Standard Ins.-tuning	2.73	8.66	5.73	24.94	25.64	25.31	34.91	15.60	26.26	52.74
PIT	1.96	14.31	8.72	30.26	33.97	30.22	10.69	15.55	27.02	55.12
SELF-TUNING	1.13	22.30	16.51	39.94	41.02	39.89	50.65	16.34	25.85	69.29
Knowledge Acquisition on Wiki-Newpages-2023-(9)10-Film (Cross-Domain Scenario)										
Open-book w/ film doc	8.30	57.38	34.45	68.64	78.92	66.31	7.35	-	-	-
Closed-book	8.30	3.35	1.88	11.27	12.97	11.49	7.35	16.05	24.67	53.40
Cont. Pre-training	5.52	3.46	2.30	11.83	14.30	11.98	12.04	16.79	25.35	56.02
Standard Ins.-tuning	2.83	5.23	3.77	16.15	17.45	16.45	51.69	14.41	25.54	49.80
PIT	1.52	6.39	4.50	16.97	18.92	17.10	3.03	13.06	23.42	54.38
SELF-TUNING	1.10	22.51	16.44	35.58	36.60	35.43	44.92	16.77	26.44	66.34

Table 3: Five-shot evaluation results on LLAMA2-7B for knowledge acquisition and retention are presented across single-domain, multi-domain, and cross-domain scenarios. For the complete results, refer to Table 8 (Appendix C).

capability in retaining commonsense knowledge using CommonsenseQA (CSQA) (Talmor et al., 2019) and report accuracy. All evaluations are conducted in a closed-book setting (see Appendix S). **Compared Methods.** We compare SELF-TUNING with three representative approaches (Table 2): (1) Continued Pre-training (Ovadia et al., 2024), (2) Standard Instruction-tuning (Saito et al., 2024), and (3) PIT (Jiang et al., 2024c), which trains on D_{train}^{QA} and D_{train}^{Doc} with QA pairs positioned before their corresponding document texts. We also evaluate

their variants (Table 7). Results, averaged over three runs, show significant differences in means ($p < 0.001$), with details in Appendix T.

5.2 Main Results

Table 3 (top) presents the evaluation results on LLAMA2-7B in relation to knowledge acquisition and retention in the single-domain scenario using the Wiki-Bio dataset.

The curated dataset exhibits minimal overlap with the pre-training data of the LLMs. The extremely low performance in the closed-book setting

(e.g., with EM around 2% for knowledge extraction) indicates that the dataset has little in common with the pre-training data, thus ensuring the reliability of the evaluation results. The non-zero EM values might be due to a small number of collected Wikipedia articles that were initially published but underwent revisions after the cut-off time.

SELF-TUNING substantially improves the LLM’s knowledge acquisition ability. SELF-TUNING greatly enhances the performance of LLAMA2-7B across three dimensions: (i) reducing PPL to nearly 1, signifying effective memorization of the new documents; (ii) increasing EM by roughly 11.5% on the knowledge extraction task, attaining performance comparable to the open-book setting; (iii) achieving high accuracy among the compared methods for the reasoning task, demonstrating excellent understanding of the newly acquired knowledge. These results reinforce the importance of first training the model to acquire the ability to absorb new knowledge before training on test documents, aligning with findings in Jiang et al. (2024c). More importantly, SELF-TUNING significantly outperforms PIT, highlighting the role of comprehension and self-reflection beyond simple memorization, validating the effectiveness of SELF-TEACHING. Further analyses are shown in Appendices C, N, D, E, and F.

SELF-TUNING excels in knowledge retention. Unlike other methods that display fluctuating performance, SELF-TUNING shows a strong ability to maintain previously acquired knowledge in terms of both knowledge extraction and reasoning. The slight improvements in evaluation metrics, such as F1 (roughly 1% on extracting learned world knowledge) and accuracy (around 13% on commonsense reasoning), compared to the closed-book performance without knowledge injection, suggest that systematically learning new knowledge doesn’t necessarily lead to catastrophic forgetting.

We further validate SELF-TUNING’s efficacy by analyzing training efficiency (Appendix K).

5.3 Results in the Multi-Domain and Cross-Domain Scenarios

To explore the potential of SELF-TUNING for enhancing LLM’s knowledge acquisition and retention in real-world scenarios, we evaluate its performance in two challenging settings (Table 3): (i) the multi-domain scenario (in the middle part); (ii) the cross-domain scenario (in the bottom part), where the training data is from Wiki-Bio, while the test

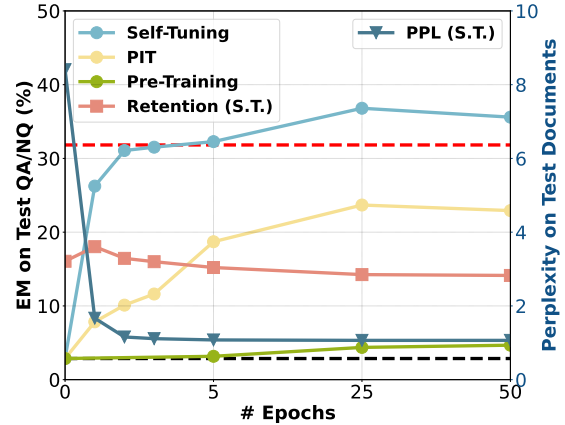


Figure 3: Training dynamics on LLAMA2-7B w.r.t., knowledge memorization, extraction, and retention across different numbers of training epochs. We present the EM scores on NQ datasets to evaluate knowledge retention. The black and red dashed lines represent the baseline closed-book and open-book performances for the knowledge extraction task, respectively.

data is from Wiki-Film.

SELF-TUNING shows strong potential in enhancing knowledge acquisition and retention across documents containing diverse new knowledge. In Table 3, SELF-TUNING consistently achieves the best performance in both settings.

The capacity to systematically absorb knowledge improves generalization ability. The substantial improvements over all compared methods in the cross-domain setting, highlight the value of equipping the model with the ability to effectively absorb knowledge from raw documents using the SELF-TEACHING strategy.

5.4 Training Dynamics

We analyze the training dynamics of SELF-TUNING during continued pre-training (beginning from Stage 2 in Figure 2) on the test documents by varying the number of training epochs for two main reasons: (i) to eliminate the possibility that the exceptional performance of SELF-TUNING in enhancing knowledge acquisition is merely a result of early fitting on the test documents, and (ii) to conduct an in-depth assessment of its long-term knowledge retention capability. Furthermore, we integrate the results of PIT and continued pre-training to offer a well-rounded evaluation.

The remarkable performance of SELF-TUNING in enhancing knowledge acquisition does not stem from early-fitting. In Figure 3, we observe that SELF-TUNING not only memorizes new knowledge more rapidly than the compared methods, lowering PPL to almost 1 within 3 epochs, but also consistently achieves the best performance

Method	Wiki-Bio (Acquisition)				
	Mem.	Extraction			Reason.
	PPL (\downarrow)	% Acc.	% EM	% Rouge	% Acc.
Cont. Pre-training	7.28	4.68	2.87	15.07	7.96
S.T. w/o Review	1.26	28.36	23.68	41.11	50.40
S.T. via Read.	1.46	20.97	17.65	34.55	39.37
S.T. w/ Pre-Review	1.28	29.86	25.94	43.31	46.91
SELF-TUNING	1.11	37.25	31.52	50.61	44.31

Table 4: Results of the SELF-TUNING variants on LLAMA2-7B on Wiki-Bio (Appendix H).

during long-term training. Remarkably, SELF-TUNING begins to outperform the open-book performance from the 5th epoch and reaches its peak at the 25th epoch with a 5% higher EM score on the knowledge extraction task.

SELF-TUNING performs well in preserving previously acquired knowledge, with only a small decline in EM of roughly 2-3% over the course of 50 training epochs. This suggests that SELF-TUNING has great potential for real-world applications.

5.5 Variants of SELF-TUNING

Setup. We investigate three variants of SELF-TUNING (in the lower part of Table 2): (i) **SELF-TUNING w/o Review**, where we continue training on test documents without the reviewing capability; (ii) **SELF-TUNING via Read.**, which displays the training documents in a reading-comprehension format (Cheng et al., 2024) (see Table 28); (iii) **SELF-TUNING w/ Pre-Review**, which trains on a mix of training documents and training QA in the second stage, before training on test documents.

Results. In Table 4, despite having lower performance than SELF-TUNING, all variants significantly enhance the model’s ability for knowledge acquisition compared to continued pre-training.

Reviewing QA ability aids in knowledge acquisition. Compared to SELF-TUNING, SELF-TUNING w/o Review exhibits inferior performance. Moreover, we suspect that the lower performance of SELF-TUNING w/ Pre-Review is because reviewing QA ability *during*, rather than before, the continuous learning of new knowledge is more effective in reducing distribution shift, thereby stabilizing the training process.

Decoupling the knowledge acquisition process into three perspectives is more effective than solely focusing on comprehension. The comparison between SELF-TUNING w/o Review and SELF-TUNING w/ Read. demonstrates that presenting the test document text from three distinct perspectives contributes more to knowledge memorization (1.26% vs. 1.46% on PPL), extraction (23.68%

Method	Acquisition				
	PPL (\downarrow)	% Acc.	% EM	% Recall	% Rouge
Varying Model (Qwen2-7B on WikiBio-2024)					
Closed-book	12.41	4.16	2.55	15.01	13.17
Stand. Ins.-tuning	2.77	11.29	9.36	25.45	24.83
PIT	1.97	11.41	9.53	25.98	25.64
SELF-TUNING	1.14	31.79	28.51	44.91	43.33
Varying Model (Mistral-7B-v0.1 on WikiBio-2023)					
Closed-book	8.45	6.64	4.37	19.51	17.25
Stand. Ins.-tuning	2.84	16.44	13.88	29.54	29.13
PIT	1.42	26.85	23.08	40.36	39.52
SELF-TUNING	1.08	41.63	36.50	55.32	52.87
Varying Corpora (LLAMA2-7B on WebNews-2023)					
Closed-book	11.20	9.04	6.30	24.22	17.99
Stand. Ins.-tuning	3.27	21.48	13.38	37.66	31.31
PIT	1.67	30.37	18.96	51.17	40.53
SELF-TUNING	1.10	37.48	28.74	56.26	48.21

Table 5: Results of Varying Models and Corpora.

vs. 17.65% on EM), and reasoning (50.40% vs. 39.37% on accuracy) than presenting the test document text with all constructed tasks as a whole.

5.6 Results of Varying Models and Corpora

Setup. We evaluate SELF-TUNING using diverse models, including Qwen2-7B (Yang et al., 2024), Mistral-7B-v0.1 (Jiang et al., 2023), and Gemma-7B (Team et al., 2024) (see Appendix L), as well as different corpora, such as WebNews-2023 (Tang and Yang, 2024), which consists of worldwide news articles from diverse websites (see Appendix U for further details).

Results. The results in Table 5 demonstrate that SELF-TUNING consistently achieves the best performance, highlighting its strong generalizability across both models and corpora. Further evaluation results for LLAMA2-13B and LLAMA2-7B-CHAT are available in Appendices I and J, respectively.

6 Conclusion

In this study, we introduce SELF-TUNING to enhance an LLM’s ability to effectively learn from raw documents through self-teaching. Specifically, we develop SELF-TEACHING, a self-supervised learning strategy that presents documents as plain texts along with various knowledge-intensive tasks derived directly from the documents. Additionally, we present three Wikipedia-Newpages-2023-QA datasets to enable a comprehensive evaluation of an LLM’s knowledge acquisition capabilities across three distinct scenarios. Our findings show that SELF-TUNING consistently yields superior performance on the knowledge acquisition tasks while showing impressive knowledge retention performance. These results suggest the potential for broader applications of SELF-TUNING.

Limitations

While our experimental results show promise, we consider these findings to be preliminary, as there are still many unexplored aspects in this field.

Combining with Continual Learning Approaches. Our study primarily focuses on enhancing a language model’s ability to effectively learn new knowledge from previously unseen raw corpora. Although experimental results on MCQA and NQ demonstrate that our SELF-TUNING method well preserves previously acquired knowledge, future research could explore integrating SELF-TUNING with continual learning approaches (Wang et al., 2024). For instance, regularization-based methods such as EWC (Kirkpatrick et al., 2017) and replay-based methods, like incorporating segments from general domain datasets (e.g., Wiki data (Zhang et al., 2024c)), could improve the model’s capacity to retain learned knowledge and skills while mitigating the risk of overfitting to new information.

In this study, we intentionally avoided using continual learning approaches to ensure a fair comparison of knowledge injection with previous methods. However, we present preliminary results of combining SELF-TUNING with a replay-based approach in Appendix M. These results confirm the strong potential of integrating SELF-TUNING with continual learning techniques to improve both knowledge acquisition and retention.

Performing More Comprehensive Evaluations of LLMs’ Knowledge Acquisition Capabilities.

In this study, we evaluate the knowledge acquisition capabilities of LLMs from three important perspectives: knowledge memorization, extraction, and reasoning. Future work could consider additional evaluation aspects, such as integrating factual knowledge with mathematical reasoning, to explore the model’s ability to utilize the learned factual knowledge in solving more complex real-world problems (Zheng et al., 2024).

Regarding Resource Demands. To verify the efficacy of SELF-TUNING, we provide a detailed analysis of training efficiency on the Wiki-News-2023-Bio dataset, conducted using 8 Tesla V100 GPUs (32GB) with LLAMA-7B, in Appendix K. This analysis demonstrates that our SELF-TUNING framework not only significantly outperforms the strongest baseline method but also

achieves this with reduced training time. Furthermore, the effectiveness of SELF-TUNING across three distinct scenarios highlights its ability to directly assimilate new knowledge from incoming test documents without requiring retraining on the original training corpus (*i.e.*, omitting the first stage). Notably, SELF-TUNING eliminates the need for any additional annotation costs. All experiments were conducted on 8 Tesla V100 GPUs (32GB), with training completing in just a few hours. Consequently, we anticipate minimal barriers to the adoption of SELF-TUNING, even for teams with limited computational resources.

Ethics Statement

Throughout our research, we have consistently adhered to ethical guidelines to uphold privacy, fairness, and the well-being of all individuals and groups involved. All benchmark datasets utilized in this study are used solely for research purposes and do not contain any personally identifiable information, thereby safeguarding privacy. During the QA data collection process for Wiki-News-2023-QA, WikiBio-2024, and WebNews-2023 using GPT-4, we meticulously crafted prompts to eliminate any language that might discriminate against specific individuals or groups. These measures were implemented to minimize potential negative effects on users’ well-being. Examples of these carefully designed prompts can be found in Table 23, Table 24, and Table 26. To further ensure the quality of the newly collected datasets, the authors manually reviewed them following the guidelines in Bai et al. (2022). These datasets were confirmed to be of high quality, free from offensive content, false information, and any personally identifiable information (Radharapu et al., 2023; Zhou et al., 2023b). Future research efforts could explore the OpenAI moderation API⁵ to systematically filter out inappropriate system responses. Additionally, GPT-4 is solely used to convert document knowledge into QA pairs—without introducing any additional information—and can be replaced by any other model.

⁵<https://platform.openai.com/docs/guides/moderation/overview>

References

- Zeyuan Allen-Zhu and Yuanzhi Li. 2023. [Physics of language models: Part 3.1, knowledge storage and extraction](#). *Preprint*, arXiv:2309.14316.
- Ronnel Ian A Ambion, Rainier Santi C De Leon, Alfonso Pio Angelo R Mendoza, and Reinier M Navarro. 2020. The utilization of the feynman technique in paired team teaching towards enhancing grade 10 anhs students’ academic achievement in science. In *2020 IEEE Integrated STEM Education Conference (ISEC)*, pages 1–3. IEEE.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.
- Daixuan Cheng, Shaohan Huang, and Furu Wei. 2024. [Adapting large language models via reading comprehension](#). In *The Twelfth International Conference on Learning Representations*.
- Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. 2023. [Crawling the internal knowledge-base of language models](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1856–1869, Dubrovnik, Croatia. Association for Computational Linguistics.
- ContextualAI. 2024. [Introducing rag 2.0](#).
- Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. [Does fine-tuning llms on new knowledge encourage hallucinations?](#) *Preprint*, arXiv:2405.05904.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *Preprint*, arXiv:2311.05232.
- Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Anthony, Timothée Lesort, Eugene Belilovsky, and Irina Rish. 2024. [Simple and scalable strategies to continually pre-train large language models](#). *Preprint*, arXiv:2403.08763.
- Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun KIM, Stanley Jungkyu Choi, and Minjoon Seo. 2022. [Towards continual knowledge learning of language models](#). In *International Conference on Learning Representations*.
- Frederick Jelinek, Robert L. Mercer, Lalit R. Bahl, and Janet M. Baker. 1977. [Perplexity—a measure of the difficulty of speech recognition tasks](#). *Journal of the Acoustical Society of America*, 62.
- Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C. Park. 2024. [Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity](#). *Preprint*, arXiv:2403.14403.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Jinhao Jiang, Junyi Li, Wayne Xin Zhao, Yang Song, Tao Zhang, and Ji-Rong Wen. 2024a. [Mix-cpt: A domain adaptation framework via decoupling knowledge learning and format alignment](#). *arXiv preprint arXiv:2407.10804*.
- Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjun Zhong, Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin Jiang, Lifeng Shang, Ruiming Tang, Qun Liu, and Wei Wang. 2024b. [Learning to edit: Aligning llms with knowledge editing](#). *Preprint*, arXiv:2402.11905.
- Zhengbao Jiang, Zhiqing Sun, Weijia Shi, Pedro Rodriguez, Chunting Zhou, Graham Neubig, Xi Victoria Lin, Wen tau Yih, and Srinivasan Iyer. 2024c. [Instruction-tuned language models are better knowledge learners](#). *Preprint*, arXiv:2402.12847.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2017. [Overcoming catastrophic forgetting in neural networks](#). *Proceedings of the National Academy of Sciences*, 114(13):3521–3526.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti,

888	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang,	949
889	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	and Yue Zhang. 2023. Survey on factuality in large	950
890	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	language models: Knowledge, retrieval and domain-	951
891	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	specificity . <i>Preprint</i> , arXiv:2310.07521.	952
892	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,		
893	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun	953
894	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	Zhu. 2024. A comprehensive survey of continual	954
895	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	learning: Theory, method and application . <i>Preprint</i> ,	955
896	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	arXiv:2302.00487.	956
897	Isabel Kloumann, Artem Korenev, Punit Singh Koura,		
898	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	957
899	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.	958
900	tiniet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Dai, and Quoc V Le. 2022. Finetuned language mod-	959
901	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	els are zero-shot learners . In <i>International Confer-</i>	960
902	stein, Rashmi Rungta, Kalyan Saladi, Alan Schelten,	<i>ence on Learning Representations</i> .	961
903	Ruan Silva, Eric Michael Smith, Ranjan Subrama-		
904	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	Kevin Wu, Eric Wu, and James Zou. 2024a. How	962
905	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	faithful are rag models? quantifying the tug-of-	963
906	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	war between rag and llms' internal prior . <i>Preprint</i> ,	964
907	Melanie Kambadur, Sharan Narang, Aurelien Ro-	arXiv:2404.10198.	965
908	driguez, Robert Stojnic, Sergey Edunov, and Thomas		
909	Scialom. 2023a. Llama 2: Open foundation and fine-	Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai	966
910	tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	Zhang, and Yanghua Xiao. 2024b. How easily do	967
911		irrelevant inputs skew the responses of large language	968
912	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	models? Preprint , arXiv:2404.03302.	969
913	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
914	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	Chong Xiang, Tong Wu, Zexuan Zhong, David Wag-	970
915	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	ner, Danqi Chen, and Prateek Mittal. 2024. Certifi-	971
916	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	ably robust rag against retrieval corruption . <i>Preprint</i> ,	972
917	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	arXiv:2405.15556.	973
918	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-		
919	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Wang Xiaofei, Chen Qing, Sun Yanyan, Tong Weifeng,	974
920	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	and Niu Wenzhi. 2017. The application of the feyn-	975
921	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	man technique for practical teaching of prosthodon-	976
922	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	tics. <i>Chinese Journal of Medical Education</i> ,	977
923	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	41(9):822.	978
924	tiniet, Todor Mihaylov, Pushkar Mishra, Igor Moly-		
925	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	Yan Xu, Mahdi Namazifar, Devamanyu Hazarika, Aish-	979
926	stein, Rashmi Rungta, Kalyan Saladi, Alan Schelten,	warya Padmakumar, Yang Liu, and Dilek Hakkani-	980
927	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	Tür. 2023. KILM: knowledge injection into encoder-	981
928	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	decoder language models . <i>CoRR</i> , abs/2302.09170.	982
929	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,		
930	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	983
931	Melanie Kambadur, Sharan Narang, Aurelien Ro-	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	984
932	driguez, Robert Stojnic, Sergey Edunov, and Thomas	Li, Dayiheng Liu, Fei Huang, Guanting Dong, Hao-	985
933	Scialom. 2023b. Llama 2: Open foundation and	ran Wei, Huan Lin, Jialong Tang, Jialin Wang,	986
934	fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	987
935		Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai,	988
936	Mozes van de Kar, Mengzhou Xia, Danqi Chen, and	Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Ke-	989
937	Mikel Artetxe. 2022. Don't prompt, search! mining-	qin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni,	990
938	based zero-shot learning with language models . In	Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize	991
939	<i>Proceedings of the 2022 Conference on Empirical</i>	Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan,	992
940	<i>Methods in Natural Language Processing</i> , pages	Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge,	993
941	7508–7520, Abu Dhabi, United Arab Emirates. As-	Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren,	994
942	sociation for Computational Linguistics.	Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing	995
943		Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan,	996
944	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry	Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang,	997
945	Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny	Zhifang Guo, and Zhihao Fan. 2024. Qwen2 techni-	998
946	Zhou, Quoc Le, and Thang Luong. 2023. Freshllms:	cal report . <i>Preprint</i> , arXiv:2407.10671.	999
947	Refreshing large language models with search engine		
948	augmentation . <i>Preprint</i> , arXiv:2310.03214.	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,	1000
		Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu	1001
	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru	Zhang. 2023. Editing large language models: Prob-	1002
	Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao,	lems, methods, and opportunities . In <i>Proceedings</i>	1003
	Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang,	<i>of the 2023 Conference on Empirical Methods in</i>	1004
		<i>Natural Language Processing</i> , pages 10222–10240,	1005

Singapore. Association for Computational Linguistics.

Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and HuaJun Chen. 2024a. *A comprehensive study of knowledge editing for large language models*. *Preprint*, arXiv:2401.01286.

Tianjun Zhang, Shishir G. Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E. Gonzalez. 2024b. *Raft: Adapting language model to domain specific rag*. *Preprint*, arXiv:2403.10131.

Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024c. *Self-alignment for factuality: Mitigating hallucinations in LLMs via self-evaluation*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1946–1965, Bangkok, Thailand. Association for Computational Linguistics.

Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. *Can we edit factual knowledge by in-context learning?* In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 4862–4876. Association for Computational Linguistics.

Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H. Chi, Quoc V Le, and Denny Zhou. 2024. *Take a step back: Evoking reasoning via abstraction in large language models*. In *The Twelfth International Conference on Learning Representations*.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. *Lima: Less is more for alignment*. *Preprint*, arXiv:2305.11206.

Jingyan Zhou, Minda Hu, Junan Li, Xiaoying Zhang, Xixin Wu, Irwin King, and Helen Meng. 2023b. *Re-thinking machine ethics – can llms perform moral reasoning through the lens of moral theories?* *Preprint*, arXiv:2308.15399.

A Additional Efforts for Knowledge Injection

Knowledge editing (Zheng et al., 2023; Yao et al., 2023; Jiang et al., 2024b; Zhang et al., 2024a) and retrieval-augmented generation (RAG) (Lewis et al., 2021; Ovadia et al., 2024; Jeong et al., 2024) are recognized as two related research initiatives in the field of knowledge injection.

(i) *Knowledge editing* (Mitchell et al., 2022; Zheng et al., 2023; Yao et al., 2023; Jiang et al., 2024b; Zhang et al., 2024a) concentrates on rectifying outdated or inaccurate factual knowledge stored in the model, without affecting other facts. In contrast, our focus lies in enabling LLMs to efficiently acquire knowledge from raw corpora.

(ii) *Retrieval-augmented generation (RAG)* (Lewis et al., 2021; Vu et al., 2023; Ovadia et al., 2024; Jeong et al., 2024) equips LLMs with new knowledge by augmenting off-the-shelf LLMs with retrieved knowledge from external sources. However, its performance is vulnerable to irrelevant or malicious information in the retrieval results (ContextualAI, 2024), potentially leading to inaccurate responses (Zhang et al., 2024b; Wu et al., 2024b; Xiang et al., 2024). Moreover, recent findings (Wu et al., 2024a) emphasize an underlying tension between a model’s prior knowledge and the information presented in retrieved documents. Consequently, this paper primarily focuses on exploring the injection of knowledge into the parameters of LLMs.

B Wiki-Newpages-2023-QA: Datasets for Studying LLM Knowledge Acquisition

To explore the knowledge acquisition capabilities of LLMs from new documents, *w.r.t.*, memorization, extraction and reasoning, we introduce the Wiki-Newpages-2023-QA datasets, which are carefully designed to minimize overlap with the initial pre-training corpus. These datasets comprise new document corpora for studying knowledge memorization and associated QA datasets for two vital knowledge-intensive tasks: open-ended generation and NLI for examining extraction and reasoning, respectively. We provide the details on dataset construction in the following subsections.

B.1 Document Collection

Given the well-structured nature of Wikipedia articles, which encompass extensive factual information and cover a wide range of topics across various domains, we gather documents from Wikipedia NewPages⁶. This collection includes new articles from diverse domains published after the pre-training cut-off time of the LLMs being evaluated, allowing us to largely ensure that the models have not been exposed to these facts. To construct the

⁶<https://en.wikipedia.org/wiki/Special:NewPages>

Document: <Sawyer Gipson-Long - Wikipedia> Alec Sawyer Gipson-Long (born December 12, 1997) is an American professional baseball pitcher for ...
QA Pair Example for Generation Task
Question: When was Sawyer Gipson-Long born? Answer: December 12, 1997.
QA Pair Example for NLI Task
Question: Based on the paragraph above can we conclude that <Alec Sawyer Gipson-Long> Sawyer Gipson-Long was born in December 1997. Options: -Yes; -It's impossible to say; -No Answer: Yes

Table 6: A simplified example of a document and its associated QA pair for the open-ended generation task. Factual information related to the QA pairs is denoted in blue.

document corpus, we specifically gather articles from September to October 2023, resulting in a total of 4,257 articles.⁷ Following Jiang et al. (2024c), we only utilize the first paragraph of each article, which provides a comprehensive summary and sufficient factual information.

B.2 QA Pair Generation

To gather QA pairs, we utilize GPT-4 (OpenAI, 2023) along with our manually curated prompts to generate a variety of questions and their corresponding answers, aiming to cover all factual information within the given document. Note that GPT-4 is solely used to convert document knowledge into QA pairs *without introducing any additional information*. It can be replaced by any other model.

Specifically, we construct QA datasets for the open-ended generation and NLI tasks by employing the prompts shown in Table 23 and Table 24, respectively. A simplified example document with associated QA pairs is provided in Table 6. More detailed examples can be found in Appendix O.

B.3 Splitting

To enable a comprehensive analysis in single-domain, multi-domain, and cross-domain situations, we develop three datasets and divide them into training and testing subsets, *ensuring zero factual knowledge overlap*.

Dataset Splitting. We create three datasets: Wiki-Newpages-2023-10-Bio (Wiki-Bio), Wiki-

Newpages-2023-10-Multi (Wiki-Multi), and Wiki-Newpages-2023-(9)10-Film (Wiki-Film). Specifically, we randomly select 1,263 biographical documents to curate Wiki-Bio, choose 2,104 documents covering various topics for constructing Wiki-Multi, and compile 955 film documents for producing Wiki-Film, using the assembled document corpus along with their associated QA pairs.

Train-test Splitting. We partition the Wiki-Bio and Wiki-Multi datasets, comprising the document corpus and the derived QA datasets, into training and testing subsets for conducting evaluations in single-domain and multi-domain contexts. We directly utilize the Wiki-Film dataset as the test set for the cross-domain scenario. It is crucial to note that the training QA datasets only contain the QA pairs from open-ended generation tasks, ensuring a fair comparison with existing knowledge injection approaches. We provide extensive statistical information for the three datasets in Table 1 and a thorough analysis of the QA types in Appendix Q.

C Evaluation Results on LLAMA2-7B

To thoroughly assess our proposed SELF-TUNING method, we compare its efficiency against three other notable approaches: standard instruction-tuning without forgetting, PIT⁺⁺, and mixed training, as shown in Table 7. The evaluation results, presented in Table 8, demonstrate that SELF-TUNING consistently outperforms the alternatives. For instance, it improves EM by 11% in the knowledge extraction task.

Combining the results from Table 8 with the training strategies in Table 7, we emphasize the importance of first training the model to develop the ability to absorb new knowledge before training on test documents. This finding aligns with the conclusions of Jiang et al. (2024c).

Notably, SELF-TUNING enables the model to absorb new knowledge from incoming test documents more efficiently. Unlike mixed training—which requires retraining on both training documents, training QA, and test documents—SELF-TUNING leverages the capabilities acquired in the initial training stage to directly learn from test documents, needing only a review of QA ability. This makes it significantly more training-efficient over time.

⁸To ensure a fair comparison, all compared approaches train on the test documents for 3 epochs in total, regardless of the number of training stages. For continued pre-training, which is observed to struggle in grasping new knowledge, we train the models for 5 epochs.

⁷The pre-training cut-off for the LLAMA2 family models used in this study is 2022.

Method	Training Data in Each Stage		
	Stage 1	Stage 2	Stage 3
Continued Pre-training		① test doc	
Standard Ins.-tuning	① train doc & test doc	② train QA	
PIT	① train QA train doc	② test doc	
SELF-TUNING	① train QA & train doc w/ self-teaching tasks	② train QA & test doc	③ test doc
Variants of SELF-TUNING			
SELF-TUNING w/o Review	① train QA & train doc w/ self-teaching tasks	② test doc	
SELF-TUNING via Read.	① train QA & train doc (reading-comp. format)	② test doc	
SELF-TUNING w/ Pre-Review	① train QA & train doc w/ self-teaching tasks	② train QA & train doc	③ test doc
Additional Compared Methods			
Standard Ins.-Tuning w/o Forget.	① train doc & test doc	② train QA & test doc	
PIT ⁺⁺	① train QA	② train QA train doc	③ test doc
Mixed Training	① train doc & train QA & test doc		

Table 7: Depiction of the training stages and associated datasets employed in the compared methods. “Train doc w/ self-teaching tasks”: the training documents presented together with the self-teaching tasks. “Reading-comp. format”: reading-comprehension format. “Forget.”: “Forgetting”.⁸

D Fine-grained Comparison

Setup. To fully understand how the ability to systematically acquire knowledge aids in the knowledge extraction task, we conduct fine-grained comparisons of PIT and SELF-TUNING on generated answers for 100 randomly sampled questions from the Wiki-bio dataset. This subset includes 56 QA types in total. Furthermore, we categorize the questions based on the fact types they contain: (i) the top-5 most common (accounting for 37%), which includes birthdate, affiliation, nationality, profession, and position/sport; (ii) time-related (accounting for 27%), such as birthdate, event date, and time period; (iii) multiple-facts (accounting for 10%), which ask about more than one fact, for example, inquiring both birth date and place; and we report the evaluation results separately. We assess the factual accuracy using exact match.

Results. As shown in Table 9, we observe that SELF-TUNING consistently outperforms PIT in the overall evaluation and the fine-grained evaluations related to different QA types. These findings underscore the importance of equipping the model with the ability to systematically acquire new knowledge. Furthermore, we present a qualitative comparison between the answers generated by PIT and SELF-TUNING in Appendix E. To gain insights into potential enhancements for SELF-TUNING, we also conduct a detailed error analysis on the types of factual errors that remain challenging after implementing SELF-TUNING in Appendix G.

E Qualitative Analysis

In Table 10, we provide a qualitative comparison between the answers generated by PIT and SELF-TUNING on the Wiki-Bio test set. We observe that SELF-TUNING performs better in answering questions that inquire about multiple facts and time-related facts, as indicated in the top part of Table 10. Furthermore, as shown in the lower part, PIT tend to fail to recall and extract facts at the end of the documents, *i.e.*, suffering from “positional bias”. This observation is consistent with the findings in Saito et al. (2024). Encouragingly, our proposed SELF-TUNING aids in recalling and extracting factual knowledge embedded at the end of the documents. These findings align with the automatic evaluation results, underscoring the effectiveness of SELF-TUNING in enhancing the LLM’s knowledge acquisition capability, particularly in knowledge extraction.

F Ablation Study

Setup. We conduct a comprehensive analysis of how comprehension and self-reflection tasks within the self-teaching tasks contribute to enhancing the LLM’s knowledge acquisition ability. We focus on two vital aspects: knowledge memorization and extraction. Specifically, we calculate the percentage of the constructed examples for each task type and systematically remove certain tasks to study their impacts.

Results. In Figure 4, we observe the following: (i)

Method	Wiki-Newpages-2023-QA (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Mem.	Extraction				Reason.	Extraction		Reasoning	
	PPL (↓)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
Knowledge Acquisition on Wiki-Newpages-2023-10-Bio (Single-Domain Scenario)										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	8.41	55.20	31.83	64.48	75.55	62.10	7.96	-	-	-
Closed-book	8.41	4.68	2.87	14.63	16.98	15.07	7.96	16.05	24.67	53.40
<i>w/ Knowledge Injection</i>										
Cont. Pre-training	7.28	6.33	3.62	15.96	18.72	16.11	15.09	16.00	24.11	53.40
Standard Ins.-tuning	6.83	6.94	5.13	19.15	19.05	19.48	39.09	15.72	23.67	51.84
Standard Ins.-Tuning w/o Forget.	2.82	9.35	7.09	21.25	21.72	21.51	36.08	16.05	24.88	54.30
PIT	2.08	14.03	11.61	27.15	28.86	27.11	11.93	15.72	26.31	57.58
PIT++	1.78	22.78	20.06	37.11	37.62	37.06	42.25	16.39	25.67	57.00
Mixed Training	1.42	24.13	20.67	38.82	39.95	38.66	55.69	19.33	28.40	58.97
SELF-TUNING	1.11	37.25	31.52	50.83	52.62	50.61	44.31	16.45	25.67	66.01
Knowledge Acquisition on Wiki-Newpages-2023-10-Multi (Multi-Domain Scenario)										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	7.84	48.93	26.63	60.37	71.71	58.54	6.33	-	-	-
Closed-book	7.84	4.53	2.73	16.19	18.63	16.38	6.33	16.05	24.67	53.40
<i>w/ Knowledge Injection</i>										
Cont. Pre-training	3.32	5.86	3.40	18.04	20.59	18.42	14.51	17.02	25.05	53.56
Standard Ins.-tuning	2.73	8.66	5.73	24.94	25.64	25.31	34.91	15.60	26.26	52.74
PIT	1.96	14.31	8.72	30.26	33.97	30.22	10.69	15.55	27.02	55.12
SELF-TUNING	1.13	22.30	16.51	39.94	41.02	39.89	50.65	16.34	25.85	69.29
Knowledge Acquisition on Wiki-Newpages-2023-(9)10-Film (Cross-Domain Scenario)										
<i>w/o Knowledge Injection</i>										
Open-book w/ film doc	8.30	57.38	34.45	68.64	78.92	66.31	7.35	-	-	-
Closed-book	8.30	3.35	1.88	11.27	12.97	11.49	7.35	16.05	24.67	53.40
<i>w/ Knowledge Injection</i>										
Cont. Pre-training	5.52	3.46	2.30	11.83	14.30	11.98	12.04	16.79	25.35	56.02
Standard Ins.-tuning	2.83	5.23	3.77	16.15	17.45	16.45	51.69	14.41	25.54	49.80
PIT	1.52	6.39	4.50	16.97	18.92	17.10	3.03	13.06	23.42	54.38
SELF-TUNING	1.10	22.51	16.44	35.58	36.60	35.43	44.92	16.77	26.44	66.34

Table 8: Five-shot evaluation results on LLAMA2-7B for knowledge acquisition and retention in three scenarios: single-domain (top), multi-domain (middle), and cross-domain (bottom). Following (Jiang et al., 2024c), we also report results for: (i) closed-book, where base LLMs are prompted with open-ended questions related to new knowledge in the test documents, and (ii) open-book w/ test doc, where base LLMs are prompted with questions along with relevant gold knowledge snippets from the test documents. Results that fall below the baseline performance are highlighted in red.

Method	Q&A Types (% EM)			
	Total	Top-5 (37%)	Time-Related (27%)	Multiple (10%)
PIT	7.00	10.81	3.70	0
SELF-TUNING	32.00	37.84	40.74	20.00

Table 9: Fine-grained evaluation results on the open-ended generation task, using the Wiki-Bio test dataset concerning the fact types of QA pairs.

comprehension tasks. These findings confirm the efficacy of the developed SELF-TEACHING strategy, underscoring the crucial role of comprehension and self-reflection in learning new knowledge for LLMs.

G Error Analysis

The examples of self-reflection tasks account for a slightly higher ratio than comprehension tasks among the self-teaching tasks. (ii) Both comprehension and self-reflection tasks benefit overall performance on the knowledge acquisition tasks. Notably, removing the examples of self-reflection tasks results in a more significant drop in performance, aligning with its higher percentage over

In order to gain insights into potential enhancements for SELF-TUNING, we outline four common errors that persist as challenges after implementing SELF-TUNING. We offer an in-depth analysis of these errors in Table 11, using EM as the evaluation metric.

Case study 1: Questions requesting information on multiple facts.

Document: <Helmut Moritz - Wikipedia> Helmut Moritz (1 November 1933 - 21 October 2022) was an Austrian physical geodesist. He was a member of the Austrian Academy of Sciences and of many other international academies and societies. He became internationally known with a fundamental work on Error propagation in Geodesy. From 1991 to 1995, he was president of the International Union of Geodesy and Geophysics (IUGG).

Question: When was Helmut Moritz born and when did he pass away?

Gold Answer: Born on November 1, 1933, passed away on October 21, 2022.

Model Answers

PIT's Answer: Information not provided.

SELF-TUNING's Answer: Born on november 1, 1933, passed away on october 21, 2022.

Case study 2: Questions inquiring about time-related details.

Document: <Brad Smiley - Wikipedia> Brad Smiley (born June 19, 1973) is an American college football coach. He is the head football coach for Southern Arkansas University; a position he has held since 2022. He also was the head coach for Trinity Valley Community College from 2007 to 2017. He also coached for Baylor, Northwestern State, and Tulane.

Question: Since when has Brad Smiley been the head football coach for Southern Arkansas University?

Gold Answer: Since 2022.

Model Answers

PIT's Answer: Since 2016.

SELF-TUNING's Answer: Since 2022.

Case study 3: Questions inquiring about facts encoded in the end of the document, *i.e.*, “positional bias”.

Document: <Nathan Saliba - Wikipedia> Nathan-Dylan Saliba (born February 7, 2004) is a Canadian professional soccer player who plays for Major League Soccer club CF Montréal.

Question: Which Major League Soccer club does Nathan Saliba play for?

Gold Answer: CF Montréal.

Model Answers

PIT's Answer: San jose earthquakes.

SELF-TUNING's Answer: CF Montréal.

Table 10: Qualitative analyses comparing the answers produced by PIT and SELF-TUNING on the open-ended generation task using the Wiki-News-2023-10-Bio test dataset. The false answers and correct answers are highlighted in red and blue, respectively.

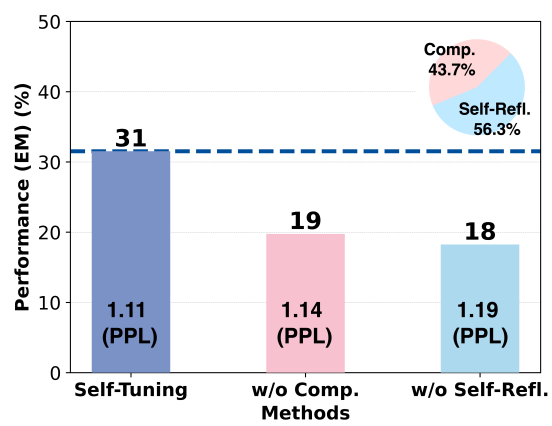


Figure 4: Ablation analysis exploring the impact of removing comprehension and self-reflection tasks from the self-teaching tasks for knowledge memorization and acquisition. The proportion of each task type among the self-teaching tasks in the training documents is shown in the upper right corner.

H Evaluation Results on SELF-TUNING Variants

Setup. To further investigate the effectiveness of SELF-TUNING, we present three variants, as depicted in Table 7: (1) **SELF-TUNING w/o Review**, where we continue training on test documents without the reviewing capability; (2) **SELF-TUNING via Read.**, which displays the training documents in a reading-comprehension format (Cheng et al., 2024) (an example is shown in Table 28); (3) **SELF-TUNING w/ Pre-Review**, which trains on a combination of training documents and training QA in the second stage, before training on test documents.

Results. In Table 12, despite having lower performance than SELF-TUNING, all variations significantly enhance the model’s ability for knowledge acquisition compared to continued pre-training, which further validates the effectiveness of SELF-TUNING in improving knowledge acquisition.

Reviewing QA ability aids in both knowledge acquisition and retention. Compared to SELF-

TUNING, SELF-TUNING w/o Review also displays inferior performance on the knowledge retention task.

I Evaluation Results on LLAMA2-13B in the Single-domain Scenario

Table 13 presents the evaluation results on LLAMA2-13B concerning knowledge acquisition and retention in the single-domain scenario using the Wiki-Bio dataset. We make the following observations:

SELF-TUNING consistently demonstrates superior performance in enhancing the model’s knowledge acquisition and retention abilities as the model size scales. As the model size scales, SELF-TUNING continues to achieve the highest performance across all evaluation metrics on memorization and acquisition tasks, consistently outperforming the compared methods by a significant margin (*e.g.*, improving EM score by 20% on the extraction task). On the reasoning task, SELF-TUNING consistently attains high accuracy. Additionally, SELF-TUNING consistently exhibits strong performance on knowledge retention tasks. These findings confirm the effectiveness of SELF-TUNING, suggesting the potential and robustness of SELF-TUNING for applications on larger-scale models.

Continued pre-training for knowledge acquisition proves challenging across all three dimensions. We find that continuing pre-training on new documents may result in a decline in knowledge extraction performance on LLAMA2-13B, compared to the baseline performance. This could be due to the fact that merely continuing pre-training might adversely affect its question-answering capability, even when equipped with new knowledge, as demonstrated by the lowered PPL. This observation is consistent with the findings in [Cheng et al. \(2024\)](#). Moreover, the marginal improvements in memorization (reducing PPL by 2%) and reasoning (increasing accuracy by 2%) suggest that such a naive approach fails to help the model memorize and capitalize on new knowledge. This highlights the importance of evaluating the model’s knowledge acquisition ability comprehensively across multiple dimensions.

J Evaluation Results on LLAMA2-7B-CHAT in the Single-domain Scenario

In this section, we showcase the evaluation outcomes for LLAMA2-7B-CHAT in Table 14. We find that even after extensive instruction-following training ([Ouyang et al., 2022a](#)), LLAMA2-7B-CHAT faces difficulty in extracting newly acquired knowledge after simply continuing pre-training on test documents. Almost all high-performing approaches struggle with knowledge retention, indicating that to incorporate new knowledge, it is preferable to train a base model rather than the version fine-tuned via RLHF (reinforcement learning from human feedback) ([Ouyang et al., 2022a](#)), despite its remarkable instruction-following capability. More significantly, SELF-TUNING consistently surpasses all other compared methods by a considerable margin on knowledge acquisition tasks. These promising outcomes further validate the effectiveness of SELF-TUNING. The results imply a potential foundation for exploring the domain of enhancing knowledge acquisition for various models.

K Training Efficiency Analysis

To ensure a fair comparison, all methods for knowledge injection presented in Table 3 were trained on raw test documents for 3 epochs, as detailed in Appendix S. Additionally, we conducted a detailed analysis of training efficiency on the Wiki-Newpages-2023-Bio dataset using 8 Tesla V100 GPUs (32G) with LLAMA2-7B:

- Continued pre-training: 112.91 seconds
- Standard instruction-tuning: 1661.06 seconds
- PIT: 6205.52 seconds
- SELF-TUNING: 5220.50 seconds

Our SELF-TUNING significantly outperforms the most competitive baseline, PIT, on the knowledge acquisition task and is more time-efficient.

L Evaluation Results on Gemma-7B in the Single-Domain Scenario

We present the evaluation results for Gemma-7B in Table 15. Our SELF-TUNING method consistently achieves the best performance, significantly outperforming the baseline methods by a substantial

margin. This observation aligns with the results reported across all other evaluation scenarios in Section 5.6.

M Evaluation Results with Continual Learning Techniques

This section explores the potential of integrating SELF-TUNING with continual learning techniques. Table 16 presents the evaluation results of combining SELF-TUNING with a representative continual learning strategy: the replay-based method. In this approach, 500 training QA pairs were randomly sampled from the Wiki QA datasets, curated prior to the pre-training cutoff date of LLAMA2-7B (Zhang et al., 2024c). These QA pairs were included throughout the training process to reinforce general domain knowledge.

The evaluation results on MCQA and NQ confirm that SELF-TUNING effectively preserves previously acquired knowledge. Moreover, integrating SELF-TUNING with replay-based continual learning (SELF-TUNING+Replay) further enhances model performance, demonstrating the following benefits:

1. **Effective knowledge acquisition:** The model successfully learns new knowledge from previously unseen raw documents. This aligns with findings in Appendix H, which highlight that the reviewing QA ability facilitates knowledge acquisition.
2. **Efficient knowledge retention:** The replay-based approach ensures that previously learned knowledge is preserved, mitigating catastrophic forgetting during the learning of new tasks.

These findings underscore the significant potential of integrating SELF-TUNING with continual learning techniques. By combining SELF-TUNING’s strengths with replay-based strategies, the model not only excels in acquiring new knowledge but also maintains strong retention of existing information, making this approach an effective solution for long-term knowledge management.

N Evaluation Results Comparing with a Baseline Utilizing Document-Based QA Generation on Test Corpora

Constructing QA data for every newly introduced raw corpus is infeasible in real-world scenarios

due to cost and scalability constraints. Our goal is to enable the model to autonomously acquire new knowledge from previously unseen raw test documents. Therefore, approaches that rely on constructing and training on QA pairs specifically tailored to test documents, as explored in prior studies (Mecklenburg et al., 2024; Jiang et al., 2024a), are not aligned with the objectives of this work.

Nevertheless, for completeness and to provide additional context, we include the results of a baseline that simultaneously trains on test documents and QA pairs generated using our proposed SELF-TEACHING strategy.

As presented in Table 17, SELF-TUNING consistently demonstrates superior performance. This comparison highlights the significant advantages of our approach in fostering autonomous knowledge acquisition without depending on pre-constructed QA pairs tailored to the test data.

O In-depth Sample Documents and Corresponding QA Pairs for Open-Ended Generation and Natural Language Inference Tasks

We present detailed sample documents along with their corresponding QA pairs for open-ended generation and natural language inference tasks in Table 18 and Table 19, respectively.

P Token Count Distribution for the Open-ended Generation Task Across the Three Datasets

The distribution of token counts for the open-ended generation task across the three datasets is depicted in Figure 5, Figure 6, and Figure 7, respectively.

Q Examination of QA Types in Open-ended Generation QA Datasets

We perform a detailed analysis of the QA types associated with the factual information in the open-ended generation QA datasets, as displayed in Table 20, by using the prompt in Table 26 with GPT-4.

R Detailed Templates used in the SELF-TEACHING Strategy

We provide the detailed templates employed in the SELF-TEACHING strategy in Table 21 and a complete example of a training document accompanied by its associated SELF-TEACHING tasks in Table 22.

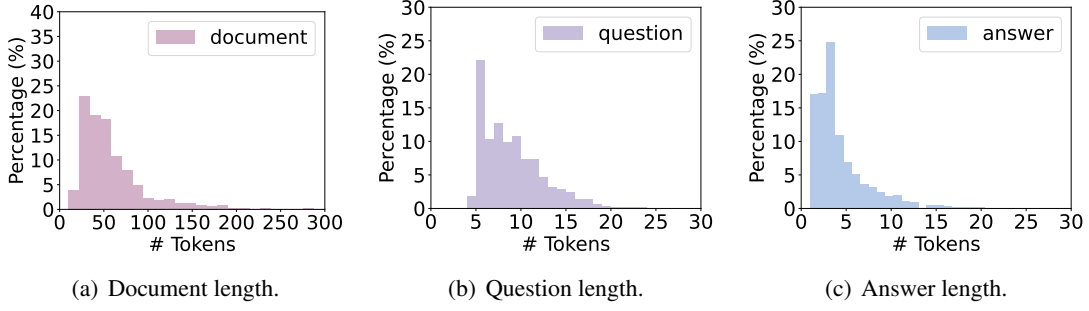


Figure 5: Distribution histogram of the token count in a document, a question, and an answer for the open-ended generation task from the Wiki-Newpages-2023-10-Bio dataset, respectively.

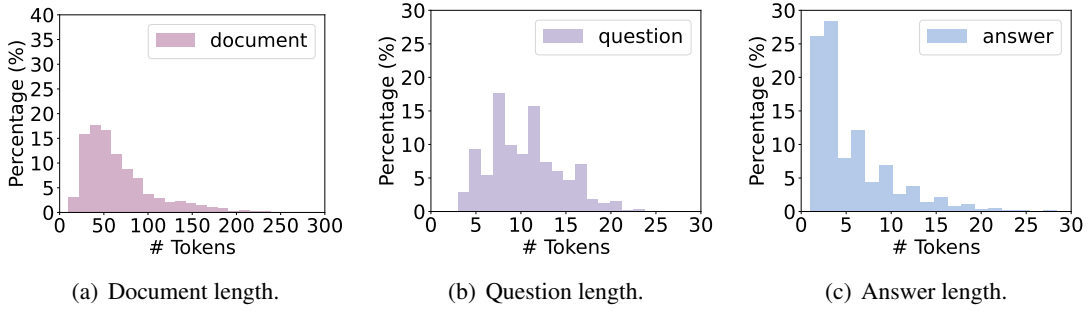


Figure 6: Distribution histogram of the token count in a document, a question, and an answer for the open-ended generation task from the Wiki-Newpages-2023-10-Multi dataset, respectively.

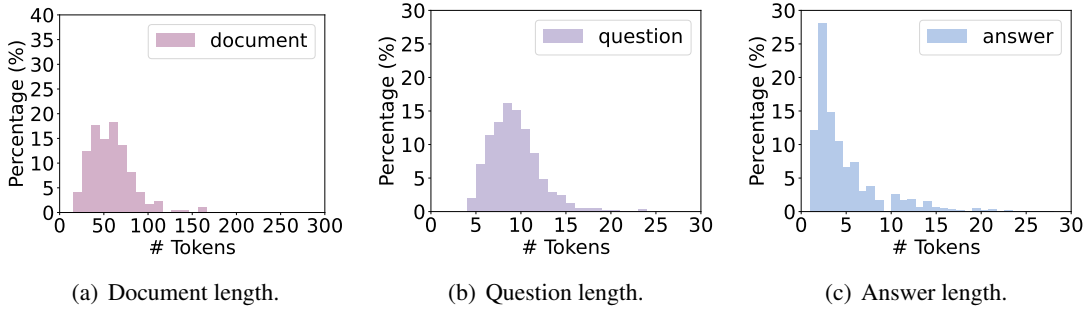


Figure 7: Distribution histogram of the token count in a document, a question, and an answer for the open-ended generation task from the Wiki-Newpages-2023-(9)10-Film dataset, respectively.

S Datasets and Evaluation Metrics

Evaluation on Knowledge Acquisition. We assess the effectiveness of SELF-TUNING in enhancing the model’s knowledge acquisition capabilities on the curated Wiki-Newpages-QA datasets, concentrating on memorization, extraction, and reasoning. (i) For memorization, we utilize test document datasets and report perplexity (Jelinek et al., 1977), which measures how well a language model predicts a text sample. (ii) For extraction, we employ test QA datasets for open-ended generation tasks. To evaluate the factual accuracy of the generated responses, we use exact match (EM), Re-

call, and F1 over words in the answer(s), following Kwiatkowski et al. (2019). Additionally, we report Rouge-L (Lin, 2004) to measure the overlap of n-grams between the generated and gold answers, accounting for minor lexical variations, following Jiang et al. (2024c). We also assess accuracy by comparing each response’s factual correctness to the gold answer, using the bidirectional entailment approach with the Deberta-Large-MNLI model (He et al., 2021). We report the five-shot evaluation results on the open-ended generation tasks using the prompt in Table 25. (iii) Concerning reasoning, we utilize the test QA datasets for NLI tasks and

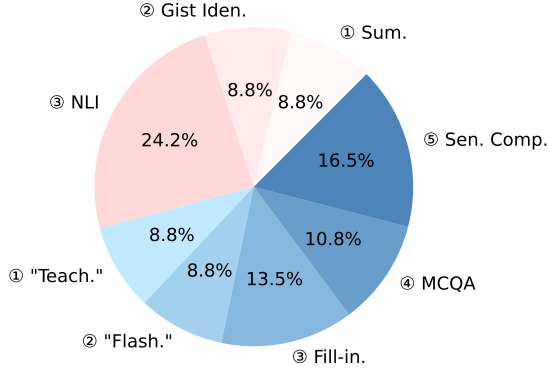


Figure 8: The percentage of constructed examples of each task type in the self-teaching tasks on training documents in Wiki-News-2023-10-Bio dataset.

report the accuracy by comparing the generated option with the gold option using EM. We present the zero-shot evaluation results on NLI tasks.

Evaluation on Knowledge Retention. It is well-known that knowledge acquisition is often accompanied by catastrophic forgetting (Allen-Zhu and Li, 2023; Wang et al., 2023). Therefore, we also provide the knowledge retention performance for a comprehensive investigation. Specifically, (i) we verify the knowledge extraction performance on world knowledge using natural questions (NQ) (Kwiatkowski et al., 2019) (*i.e.*, NQ-open (Min et al., 2021) in the closed-book setting) and report EM and F1 scores. We report the five-shot evaluation results using the first five QA pairs in the dev sets as prompts. (ii) we assess the reasoning capability on Commonsense knowledge using CommonsenseQA (CSQA) (Talmor et al., 2019), employing accuracy to assess the correctness of the selected option, calculated by comparing the generated option against the gold option using EM. We present the five-shot performance on the dev sets, as the test set does not contain golden annotations, and use the first five multi-choice QA pairs in the training set as prompts. We use these two datasets because they were curated before the cut-off time of LLAMA2 family models (*i.e.*, year 2022), making it likely that the models have obtained relevant knowledge in these datasets during the pre-training stage, as evidenced by Touvron et al. (2023a).

T Implementation Details

Training Details. We utilize LLAMA2-7B for our investigation and provide analyses on Qwen2-7B, Mistral-7B-v0.1, Gemma-7B, LLAMA2-13B,

and LLAMA2-7B-CHAT for a comprehensive understanding. We use the following training objectives: (i) for training on document data D^{Doc} , we compute the standard next-token prediction loss by averaging over all tokens in the document d (Equation 4); (ii) for training on QA data D^{QA} , we compute the average negative log-likelihood loss only on tokens in the answer a given the question q (Equation 5), where $|d|$ and $|a|$ refer to the length of the tokenized document sequence and answer sequence, respectively.

$$L_{\theta}(D^{Doc}) = -\frac{1}{|d|} \sum_t \log p_{\theta}(d_t | d_{<t}) \quad (4)$$

$$L_{\theta}(D^{QA}) = -\frac{1}{|a|} \sum_t \log p_{\theta}(a_t | q, a_{<t}) \quad (5)$$

We train LLAMA2-7B, Qwen2-7B, Mistral-7B-v0.1, Gemma-7B, and LLAMA2-7B-CHAT on 8 32GB Tesla V100 GPUs using a batch size of 8 and a learning rate of 5e-6. Additionally, we train LLAMA2-13B on 8 A100-SXM4-40GB GPUs with a batch size of 8 and a learning rate of 5e-6. To ensure a fair comparison, all compared approaches train on the test documents for 3 epochs in total, regardless of the number of training stages. For continued pre-training, which is observed to struggle in grasping new knowledge, we train the models for 5 epochs. The specific number of training epochs used for each approach in Table 7 are as follows:

- **Continued Pre-training** trains the model on the D_{test}^{Doc} dataset for 5 epochs.
- **Standard Instruction-tuning** first trains on both D_{train}^{Doc} and D_{test}^{Doc} datasets, then fine-tunes on D_{train}^{QA} dataset for 3 epochs.
- **PIT** (Jiang et al., 2024c) first trains on D_{train}^{QA} and D_{train}^{Doc} datasets for 3 epochs, positioning the QA pairs right before the corresponding document texts, then trains on the D_{test}^{Doc} data for 3 epochs.
- **SELF-TUNING** (ours) first trains on D_{train}^{QA} and D_{train}^{Doc} with the created instruction-following dataset D_{train}^{Self} (in the QA format) using the SELF-TEACHING strategy for 2 epochs, then continues training on D_{test}^{Doc} data while reviewing the D_{train}^{QA} data for 1 epoch, and finally continues training on D_{test}^{Doc} data for 2 epochs. In addition, we provide the percentage of SELF-TEACHING task examples on training

documents in Wiki-Newpages-2023-10-Bio dataset in Figure 8.

Specifically, in the cross-domain setting, where there is a substantial difference between the domains of the training data and test documents, we continue training on D_{test}^{Doc} data while reviewing the D_{train}^{QA} data for 2 epochs after the initial training stage, followed by further training on D_{test}^{Doc} data for 1 epoch. Furthermore, we adopt the same training strategy when dealing with LLAMA2-7B-CHAT, where the process of knowledge injection poses a significant challenge, as demonstrated by our experimental results. In accordance with Jiang et al. (2024c), for PIT and SELF-TUNING, we include 64 examples and 128 examples randomly sampled from D_{train}^{QA} datasets, respectively, during the final training stages when solely training on the D_{test}^{Doc} data, to prevent the model from losing its question-answering capabilities. It is important to note that all evaluation results are reported at the temperature $T = 1$.

Training Details for SELF-TUNING Variants.

- **SELF-TUNING w/o Review** first trains on D_{train}^{QA} and D_{train}^{Doc} with the created instruction-following dataset D_{train}^{Self} (in the QA format) using the SELF-TEACHING strategy for 2 epochs, then continues training on D_{test}^{Doc} data for 3 epochs.
- **SELF-TUNING via Read.** initially trains on D_{train}^{QA} and D_{train}^{Doc} (in the read-comprehension format, as shown in Table 28 for 3 epochs, then trains on the D_{test}^{Doc} data for 3 epochs.
- **SELF-TUNING w/ Pre-Review** first trains on D_{train}^{QA} and D_{train}^{Doc} with the created instruction-following dataset D_{train}^{Self} (in the QA format) using the SELF-TEACHING strategy for 2 epochs, then continues training on D_{train}^{Doc} and D_{train}^{QA} data for 1 epoch, and finally continues training on D_{test}^{Doc} data for 3 epochs.

Training Details for Additional Compared Methods.

- **Standard Instruction-Tuning w/o Forgetting** initially trains on the mixture of D_{train}^{Doc} and D_{test}^{Doc} for 3 epochs, then on D_{train}^{QA} and D_{test}^{QA} datasets for 1 epoch.
- **PIT⁺⁺** (Jiang et al., 2024c) initially trains on D_{train}^{QA} for 1 epoch, then on D_{train}^{QA} and

D_{train}^{Doc} datasets for 3 epochs, with the QA pairs placed right before the corresponding document texts, and finally, it trains on the D_{test}^{Doc} data for 3 epochs.

- **Mixed Training** trains on mixture of the D_{train}^{Doc} , D_{test}^{Doc} and D_{train}^{QA} datasets simultaneously for 3 epochs.

Future research could explore the inclusion of segments from general domain datasets, such as Wiki data (Zhang et al., 2024c) and the Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2021), which were compiled prior to the pre-training cut-off date. Adopting this strategy may improve the model’s capacity to retain learned knowledge and skills while reducing the risk of overfitting to novel information. In our current study, we deliberately avoid integrating extra data to ensure a precise assessment of knowledge injection, thereby preventing any biases that might arise from the inclusion of additional sources.

Prompts Employed in this Study. The prompts used for constructing the QA datasets for open-ended generation and NLI tasks are presented in Table 23 and Table 24, respectively. The prompt used during the evaluation process is displayed in Table 25. The prompt used by GPT-4 for annotating QA types in the open-ended generation tasks of the Wiki-Newpages-2023-QA datasets is presented in Table 26.

U Implementation Details of Evaluation on Varying Models and Corpora

U.1 Evaluation on Different Models

For the evaluation using different models, specifically Qwen2-7B (Yang et al., 2024), we collected articles published on Wikipedia NewPages from June 2024 to September 2024 to minimize overlap with the pre-training corpus. We randomly selected 146 biographies from the collected articles, following the data construction pipeline described in Appendix B, to create a new question-answering dataset for an open-ended generation task. This resulted in a test set, named WikiBio-2024, comprising 146 documents and a total of 827 QA pairs.

U.2 Evaluation on Varied Corpora

For the evaluation using varied corpora, we utilized the news data collected by Tang and Yang

(2024) using the mediastack API⁹. Specifically, this dataset includes articles published from September 26, 2023, to December 26, 2023, which is beyond the pre-training cutoff time of LLAMA2-7B. The dataset covers a range of news categories, such as entertainment, business, technology, and science.

For each factual sentence extracted from the original articles by Tang and Yang (2024), we concatenated the article title with the fact to create a knowledge snippet. Following Tang and Yang (2024), we used GPT-4o (OpenAI, 2024) (version dated 2024-02-01) to first paraphrase these snippets to make them clearer and more concise, and then generate relevant QA pairs. The prompt utilized can be found in Table 27. Using the data construction pipeline described in Appendix B, we generated a new training set, *i.e.*, WebNews-2023, consisting of 1,800 training documents and 6,038 QA pairs, as well as a testing set with 400 testing documents and 1,350 QA pairs.

V A Sample Training Document in the Reading-Comprehension Format

Drawing inspiration from Cheng et al. (2024), we restructure the training document in the reading-comprehension text format. Each raw text is enriched with a series of tasks related to its content, constructed using our proposed SELF-TEACHING strategy. An example of a training document is provided in Table 28.

⁹<https://mediastack.com/>

Type	Fraction	Example			
		Document	Question	Gold Answer	Model Answer
Wrong answer	76.47%	<Jalen Mack - Wikipedia> Jalen Mack (born August 5, 2005) is an American professional stock car racing driver who competes part-time in the ARCA Menards Series and ARCA Menards Series East, driving the No. 43 Chevrolet for Tamayo Cosentino Racing . He also competes part time in the ARCA Menards Series West, driving the No. 83 Chevrolet for Mack Motorsports in conjunction with Bill McAnally Racing.	Which team does Jalen Mack drive for in the ARCA Menards Series and ARCA Menards Series East?	Tamayo Cosentino Racing.	Venturini motorsports.
Higher granularity	7.35%	<Andriyko Olha Fedorivna - Wikipedia> Andriyko Olha Fedorivna (born January 28, 1945, Voronkiv, Kyiv region) is a Doctor of Law, Professor , Head of the Department of Constitutional, Administrative and Financial Law of the Kyiv University of Law of the National Academy of Sciences of Ukraine, and Deputy Head of the Department of State and Legal Problems of Management of the V. M. Koretsky Institute of State and Law of the National Academy of Sciences of Ukraine.	What are Andriyko Olha Fedorivna’s academic and professional titles?	Doctor of Law, Professor.	Doctor of law, professor, head of the department of constitutional, administrative, and financial law of the kyiv university of law of the national academy of sciences of ukraine.
Lower granularity	5.88%	<Mike Babcock (American football) - Wikipedia> Michael Babcock (born February 13, 1979) is an American college football coach. He is the head football coach for McKendree University; a position he has held since 2013 . He also coached for UCLA, Colorado, San Diego, and CSU Pueblo. He played college football for UCLA as a linebacker.	Since when has Mike Babcock (American football) held the head coach position at McKendree University?	Since 2013.	2013.
Paraphrase	10.29%	<Lil Tay - Wikipedia> Tay Tian (born July 29, 2009), known professionally as Lil Tay, is an American-born Canadian internet personality and singer. In 2018, she gained prominence online for a period of three months, proclaiming herself to be the “youngest flexer of the century”. During her brief career, she posted rap videos on YouTube and Instagram which garnered tens of millions of views. Her career ended in mid-2018, after her father applied to the superior court of Canada for full custody and control of her career. According to court documents, he was abusive and largely an absentee.	What is Lil Tay’s nationality?	American-born Canadian.	Canadian-American.

Table 11: Analysis on the types of factual errors that remain challenging after applying SELF-TUNING.

Method	Wiki-Newpages-2023-10-Bio (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Mem.	Extraction					Reason.	Extraction		Reasoning
	PPL (\downarrow)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
Continued Pre-training	7.28	4.68	2.87	14.63	16.98	15.07	7.96	16.05	24.67	53.40
SELF-TUNING w/o Review	1.26	28.36	23.68	41.29	41.93	41.11	50.40	15.55	24.20	65.11
SELF-TUNING via Read.	1.46	20.97	17.65	34.54	39.19	34.55	39.37	18.43	27.99	62.74
SELF-TUNING w/ Pre-Review	1.28	29.86	25.94	43.46	44.96	43.31	46.91	16.28	24.80	65.11
SELF-TUNING	1.11	37.25	31.52	50.83	52.62	50.61	44.31	16.45	25.67	66.01

Table 12: Five-shot evaluation results of the SELF-TUNING variants on LLAMA2-7B in the single-domain scenario. Results that fall below the baseline closed-book performance (previously shown in Table 3) are highlighted in red.

Method	Wiki-Newpages-2023-10-Bio (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Memorization	Extraction					Reason.	Extraction		Reasoning
	PPL (\downarrow)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
LLAMA2-13B										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	8.27	58.97	37.41	70.38	78.64	68.09	3.57	-	-	-
Closed-book	8.27	6.33	4.68	17.45	19.37	17.58	3.57	19.84	28.71	66.34
<i>w/ Knowledge Injection</i>										
Con. Pre-training	6.35	4.98	3.77	17.12	18.95	17.04	5.49	21.25	30.35	66.34
Standard Ins.-tuning	3.00	12.67	10.11	26.79	27.42	27.00	52.43	19.95	30.95	65.77
PIT	1.70	22.93	19.61	36.50	36.99	36.25	59.40	19.05	31.02	70.93
SELF-TUNING	1.09	44.19	39.37	58.31	60.47	57.90	54.18	20.69	31.62	71.50

Table 13: Five-shot evaluation results on LLAMA2-13B for knowledge acquisition and retention in the single-domain scenario. Results that are inferior to closed-book performance without knowledge injection are indicated in red.

Method	Wiki-Newpages-2023-10-Bio (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Memorization	Extraction					Reason.	Extraction		Reasoning
	PPL (\downarrow)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
LLAMA2-7B-CHAT										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	12.36	71.34	43.74	75.11	88.38	73.74	31.14	-	-	-
Closed-book	12.36	5.58	4.07	16.05	17.63	16.19	31.14	18.20	26.84	67.16
<i>w/ Knowledge Injection</i>										
Con. Pre-training	8.12	5.73	3.32	15.89	18.60	15.81	24.83	18.32	27.01	65.19
Standard Ins.-tuning	2.99	12.67	10.56	25.13	25.41	25.38	67.76	14.81	23.72	58.07
PIT	1.85	15.54	13.12	29.03	29.47	29.45	39.51	14.92	23.38	62.33
SELF-TUNING	1.10	33.03	29.41	46.94	47.90	47.00	72.29	13.57	22.28	64.21

Table 14: Five-shot evaluation results on LLAMA2-7B-CHAT for knowledge acquisition and retention in the single-domain scenario. Results that are inferior to closed-book performance without knowledge injection are indicated in red.

Method	Knowledge Acquisition					
	PPL (\downarrow)	% Acc.	% EM	% F1	% Recall	% Rouge
<i>w/o Knowledge Injection</i>						
Closed-book	12.41	7.09	4.68	17.60	18.10	17.65
<i>w/ Knowledge Injection</i>						
Continued Pre-training	3.99	8.14	6.33	19.91	20.97	19.82
Standard Instruction-tuning	10.13	10.41	8.60	24.06	23.86	24.16
PIT	4.19	8.14	5.88	20.87	20.78	20.68
Self-Tuning	1.09	41.93	36.80	56.95	57.41	56.30

Table 15: Evaluation results of different methods applied to Gemma-7B on the Wiki-Newpages-2023-Bio dataset.

Method	Wiki-Newpages-2023-QA (Acquisition)						NQ (Reten.)		CSQA (Reten.)	
	Memorization	Extraction				Reason.	Extraction		Reasoning	
	PPL (\downarrow)	% Acc.	% EM	% F1	% Rec.	% Rouge	% Acc.	% EM	% F1	% Acc.
Knowledge Acquisition on Wiki-Newpages-2023-10-Bio (Single-Domain Scenario)										
<i>w/o Knowledge Injection</i>										
Open-book w/ test doc	8.41	55.20	31.83	64.48	75.55	62.10	7.96	-	-	-
Closed-book	8.41	4.68	2.87	14.63	16.98	15.07	7.96	16.05	24.67	53.40
<i>w/ Knowledge Injection</i>										
PIT	2.08	14.03	11.61	27.15	28.86	27.11	11.93	15.72	26.31	57.58
SELF-TUNING	1.11	37.25	31.52	50.83	52.62	50.61	44.31	16.45	25.67	66.01
SELF-TUNING+Replay	1.03	44.49	39.82	58.44	60.58	58.00	56.24	22.67	33.86	73.55

Table 16: Five-shot evaluation results of LLAMA2-7B combined with continual learning techniques for knowledge acquisition and retention in the single-domain scenario.

Method	Knowledge Acquisition (Wiki-Newpages-2023-Bio)					
	PPL (\downarrow)	% Acc.	% EM	% F1	% Recall	% Rouge
<i>w/o Knowledge Injection</i>						
Open-book w/ test doc	8.41	55.20	31.83	64.48	75.55	62.10
Closed-book	8.41	4.68	2.87	14.63	16.98	15.07
<i>w/ Knowledge Injection</i>						
Continued Pre-training	7.28	6.33	3.62	15.96	18.72	16.11
Training on test doc w/ QA pairs	1.08	15.84	12.07	28.58	31.06	28.07
SELF-TUNING	1.11	37.25	31.52	50.83	52.62	50.61

Table 17: Evaluation results comparing SELF-TUNING to training on test documents with constructed QA pairs using LLAMA2-7B for knowledge acquisition on the Wiki-Newpages-2023-Bio dataset.

Sample document and associated QA pairs for open-ended generation tasks

Dataset: Wiki-Newpages-2023-10-Bio

Document: <Helmut Moritz - Wikipedia> Helmut Moritz (1 November 1933 - 21 October 2022) was an Austrian physical geodesist. He was a member of the Austrian Academy of Sciences and of many other international academies and societies. He became internationally known with a fundamental work on Error propagation in Geodesy. From 1991 to 1995, he was president of the International Union of Geodesy and Geophysics (IUGG).

Question: When was Helmut Moritz born and when did he pass away?

Answer: Born on November 1, 1933, passed away on October 21, 2022.

Question: What was Helmut Moritz's profession?

Answer: Austrian physical geodesist.

Question: Which academies and societies was Helmut Moritz a member of?

Answer: Austrian Academy of Sciences, many other international academies, and societies.

Question: What work made Helmut Moritz internationally known?

Answer: A fundamental work on Error propagation in Geodesy.

Question: What position did Helmut Moritz hold from 1991 to 1995?

Answer: President of the International Union of Geodesy and Geophysics (IUGG).

Dataset: Wiki-Newpages-2023-10-Multi

Document: <2018 California Proposition 71 - Wikipedia> Proposition 71, also known as Prop 71, was a California ballot proposition and proposed state constitution amendment to change the effective date of passed ballot measures from the day after the election to the fifth day after the Secretary of State certified the results.\n\n Stated goals of the measure was to ensure results were official before new measures were implemented. Opposers fearing a delay in urgent measures. Kevin Mullin supported the amendment. The California Democratic Party endorsed the amendment. Rural County Representatives of California also endorsed the amendment.

Question: What was the 2018 California Proposition 71, also known as Prop 71?

Answer: A California ballot proposition, proposed state constitution amendment, change effective date of passed ballot measures.

Question: What was the proposed change in the effective date of passed ballot measures in the 2018 California Proposition 71?

Answer: From the day after the election, to the fifth day after the Secretary of State certified the results.

Question: What were the stated goals of the 2018 California Proposition 71?

Answer: To ensure results were official before new measures were implemented.

Question: What concern did opposers of the 2018 California Proposition 71 have?

Answer: A delay in urgent measures.

Question: Who supported the 2018 California Proposition 71 amendment?

Answer: Kevin Mullin.

Question: Which organizations endorsed the 2018 California Proposition 71 amendment?

Answer: The California Democratic Party, Rural County Representatives of California.

Dataset: Wiki-Newpages-2023-(9)10-Film

Document: <Krazy House (film) - Wikipedia> Krazy House is an upcoming Dutch comedy film. It is written, directed, and co-produced by Steffen Haars and Flip van der Kuil in their English-language feature debut. Shot on location in Amsterdam, the film stars Nick Frost, Kevin Connolly and Alicia Silverstone. Maarten Swart is producer for Kaap Holland Films.

Question: What is Krazy House (film)?

Answer: An upcoming Dutch comedy film.

Question: Who are the writers, directors, and co-producers of Krazy House (film)?

Answer: Steffen Haars, Flip van der Kuil.

Question: What is significant about Steffen Haars and Flip van der Kuil's involvement in Krazy House (film)?

Answer: It is their English-language feature debut.

Question: Where was Krazy House (film) shot?

Answer: On location in Amsterdam.

Question: Who is the producer of Krazy House (film) and which production company is involved?

Answer: Maarten Swart, Kaap Holland Films.

Table 18: Sample document and associated QA pairs for open-ended generation tasks in Wiki-Newpages-2023-10-Bio, Wiki-Newpages-2023-10-Multi, and Wiki-Newpages-2023-(9)10-Film datasets.

Sample document and associated QA pairs for natural language inference tasks

Dataset: Wiki-Newpages-2023-10-Bio

Document: <Sawyer Gipson-Long - Wikipedia> Alec Sawyer Gipson-Long (born December 12, 1997) is an American professional baseball pitcher for the Detroit Tigers of Major League Baseball (MLB). He made his MLB debut in 2023.

Question: Based on the paragraph above can we conclude that <Alec Sawyer Gipson-Long> Sawyer Gipson-Long was born in December 1997. Options: -Yes; -It's impossible to say; -No

Answer: Yes

Question: Based on the paragraph above can we conclude that <Alec Sawyer Gipson-Long> Sawyer Gipson-Long is a professional football player. Options: -Yes; -It's impossible to say; -No

Answer: No

Question: Based on the paragraph above can we conclude that <Alec Sawyer Gipson-Long> Sawyer Gipson-Long plays for the Detroit Tigers in Major League Baseball. Options: -Yes; -It's impossible to say; -No

Answer: Yes

Question: Based on the paragraph above can we conclude that <Alec Sawyer Gipson-Long> Sawyer Gipson-Long made his MLB debut in 2020. Options: -Yes; -It's impossible to say; -No

Answer: No

Dataset: Wiki-Newpages-2023-10-Multi

Document: <2023 Astana Open 2013 Singles - Wikipedia> Novak Djokovic was the reigning champion, but chose not to compete this year.Seeds.

Question: Based on the paragraph above can we conclude that <2023 Astana Open 2013 Singles> Novak Djokovic won the previous Astana Open singles tournament.Options: -Yes; -It's impossible to say; -No

Answer: Yes

Question: Based on the paragraph above can we conclude that <2023 Astana Open 2013 Singles> Novak Djokovic is participating in the 2023 Astana Open singles tournament.Options: -Yes; -It's impossible to say; -No

Answer: No

Question: Based on the paragraph above can we conclude that <2023 Astana Open 2013 Singles> The 2023 Astana Open is a tennis tournament.Options: -Yes; -It's impossible to say; -No

Answer: It's impossible to say

Question: Based on the paragraph above can we conclude that <2023 Astana Open 2013 Singles> Novak Djokovic was injured and could not compete in the 2023 Astana Open singles tournament.Options: -Yes; -It's impossible to say; -No

Answer: It's impossible to say

Dataset: Wiki-Newpages-2023-(9)10-Film

Document: <Unstoppable (2023 film) - Wikipedia> Unstoppable is a 2023 comedy-drama film directed by Diamond Ratnababu and produced by Rajith Rao under AB2 Productions. The film was released theatrically worldwide on 9 June 2023.

Question: Based on the paragraph above can we conclude that<Unstoppable (2023 film)> Unstoppable is a film that combines elements of comedy and drama.Options: -Yes; -It's impossible to say; -No

Answer: Yes

Question: Based on the paragraph above can we conclude that<Unstoppable (2023 film)> Diamond Ratnababu is the producer of the film Unstoppable.Options: -Yes; -It's impossible to say; -No

Answer: No

Question: Based on the paragraph above can we conclude that<Unstoppable (2023 film)> Unstoppable was released in theaters worldwide.Options: -Yes; -It's impossible to say; -No

Answer: Yes

Question: Based on the paragraph above can we conclude that<Unstoppable (2023 film)> The film Unstoppable was released before June 2023.Options: -Yes; -It's impossible to say; -No

Answer: No

Question: Based on the paragraph above can we conclude that<Unstoppable (2023 film)> The film Unstoppable was distributed by Diamond Ratnababu.Options: -Yes; -It's impossible to say; -No

Answer: It's impossible to say

Table 19: Sample document and associated QA pairs for natural language inference tasks in Wiki-Newpages-2023-10-Bio, Wiki-Newpages-2023-10-Multi, and Wiki-Newpages-2023-(9)10-Film test datasets.

Dataset	QA Type Instances	QA Types		QA Types w/ Multiple Facts	
		Statistics	Top-5 Types	Statistics	Top-5 Types
Wiki-Newpages-2023-10-Bio (Single-domain)					
Train	Birth Date, Achievements, Position, <i>etc.</i>	2014 (# Types); 6073 (# Counts)	Birth Date (11.24%) Nationality (5.37%) Profession (5.15%) Team/Affiliation (3.05%) Role/Position (2.56%)	158 (# Types); 265 (# Counts)	Birth & Death Dates (0.93%) Birth Date & Place (0.44%) Death Date & Place (0.12%) Nationality & Profession (0.10%) Current Position & Tenure (0.08%)
Test	Full Name, Affiliation, Residence, <i>etc.</i>	281 (# Types); 655 (# Counts)	Birth Date (13.11%) Profession (6.18%) Nationality (5.62%) Team/Affiliation (4.49%) Role/Position (3.00%)	16 (# Types); 30 (# Counts)	Birth Date & Place (1.31%) Birth & Death Dates (1.12%) Death Date & Place (0.56%) Car Number & Manufacturer (0.37%) Current Club & League (0.19%)
Within the train and test sets, there are 63 and 8 answers labeled as “Information not provided/missing,” respectively.					
Wiki-Newpages-2023-10-Multi (Multi-domain)					
Train	Album Source, Location, Season Number, <i>etc.</i>	4813 (# Types); 9973 (# Counts)	Birth Date (3.37%) Profession (1.76%) Nationality (1.47%) Location (1.39%) Release Date (1.27%)	303 (# Types); 371 (# Counts)	Birth & Death Dates (0.32%) Birth Date & Place (0.14%) Event Date & Location (0.06%) Death Date & Place (0.06%) Nationality & Profession (0.05%)
Test	Legacy/Impact, Purpose, Leadership, <i>etc.</i>	924 (# Types); 1498 (# Counts)	Birth Date (3.06%) Release Date (1.80%) Profession (1.57%) Nationality (1.25%) Team/Affiliation (1.02%)	57 (# Types); 66 (# Counts)	Birth & Death Dates (0.31%) Birth Date & Place (0.31%) Death Date & Place (0.16%) Job Titles & Affiliations (0.16%) Language & Genre (0.16%)
Within the train and test sets, there are 31 and 4 answers labeled as “Information not provided/missing,” respectively.					
Wiki-Newpages-2023-(9)10-Film (Single-domain)					
Test	Director, Actor, Music Composer, <i>etc.</i>	339 (# Types); 955 (# Counts)	Director (9.07%) Release Date (7.23%) Genre (6.96%) Cast (3.55%) Language (2.76%)	13 (# Types); 15 (# Counts)	Title & Release Year (0.39%) Milestone & Historical Comparison (0.13%) Profession & Industry (0.13%) Cast & Roles (0.13%) Producer & Production Banner (0.13%)

Table 20: A comprehensive analysis of QA types related to factual information in open-ended generation QA datasets from Wiki-News-2023-10-Bio (Wiki-Bio), Wiki-News-2023-10-Multi (Wiki-Multi), and Wiki-News-2023-(9)10-Film (Wiki-Film).

Type	Task	Template
Memorization		
Next-Token Prediction	Text-to-Text	<Document>
Comprehension		
① Summarization	Text-to-Topic	Question: Write a title: <Document>. Answer: <Title>.
② Gist Identification	Text-to-Word	Question: Highlight the key information within the article: <Document>. Answer: <Entity1>, <Entity2>, etc.
③ Natural Language Inference	Text-to-Option	Question: <Document> Based on the article above can we conclude that <Sentence>. Options: -Yes; -It's impossible to say; -No. Answer: Yes/It's impossible to say/No.
Self-Reflection		
① "Teaching"	Topic-to-Text	Question: Tell me about <Title>. Answer: <Document>.
② "Flashcards"	Word-to-Text	Question: Generate a concrete description about <Title> based on the following keywords: <Entity>, etc. Answer: <Document>.
③ Fill-in-the-Blank	Cloze Entity Sentence-to-Entity	Question: <Title> <Sentence_Part1> – <Sentence_Part2> (w/o <Entity>). Answer: <Entity>.
④ Multi-Choice QA	Cloze Sentence (w/ options)-to-Entity	Question: <Title> <Sentence_Part1> – <Sentence_Part2> (w/o <Entity>) Options: - <Entity1>; - <Entity2>, etc. Answer: <Entity>.
⑤ Sentence Completion	Text-to-Text	Question: <Title> <Sentence_Part1>: Answer: <Sentence_Part2>.

Table 21: The detailed templates for each task used in the SELF-TEACHING learning strategy.

Type	Example
Memorization	
Next-Token Prediction	<Robert Anderson (artist) - Wikipedia> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.
Comprehension	
① Summarization	Question: Write a title: <Robert Anderson (artist) ... stamps. Answer: Robert Anderson (artist).
② Gist Identification	Question: Highlight the key information within the article: <Robert Anderson (artist) ... stamps. Answer: United States; American; Alan Greenspan; George W. Bush; Robert Alexander Anderson; 1946
③ Natural Language Inference	Question: <Robert Anderson (artist) ... stamps. Based on the article above can we conclude that <Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps. Options: - Yes - It's impossible to say - No Answer: Yes
Self-Reflection	
① "Teaching"	Question: Tell me about Robert Anderson (artist). Answer: Robert Alexander Anderson (born 1946) is ... stamps.
② "Flashcards"	Question: Generate a concrete description about Robert Anderson (artist), based on the following keywords: United States; American; Alan Greenspan; George W. Bush; Robert Alexander Anderson; 1946 Answer: Robert Alexander Anderson (born 1946) is ... stamps.
③ Fill-in-the-Blank	Question: <Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American – known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps. Answer: Portrait artist.
④ Multi-Choice QA	Question: <Robert Anderson (artist)> - (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps. Options: - Alan Greenspan - 1946 - Robert Alexander Anderson - George W. Bush Answer: Robert Alexander Anderson.
⑤ Sentence Completion	Question: <Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as: Answer: Designing United States postage stamps.

Table 22: An example of a training document from the Wiki-Newpages-2023-10-Bio train set, accompanied by related self-teaching tasks.

The prompt utilized by GPT-4 for building QA datasets for open-ended generation tasks

Below is a paragraph about the 51st International Emmy Awards ceremony. Your task is to formulate a detailed list of questions and corresponding answers that encompass all the information within the paragraph. To ensure clarity, each question should explicitly mention the 51st International Emmy Awards ceremony. Answers should be concise, consisting of a few short phrases separated by commas. For instance:

Paragraph: The 51st International Emmy Awards ceremony, presented by the International Academy of Television Arts and Sciences (IATAS), occurred on November 20, 2023, at the New York Hilton Midtown in New York City. It was held to acknowledge the best television programs initially produced and aired outside the United States in 2022. Nominations were announced on September 26, 2023.

Question: When was the 51st International Emmy Awards ceremony held?

Answer: November 20, 2023.

Question: Who was responsible for presenting the 51st International Emmy Awards ceremony?

Answer: The International Academy of Television Arts and Sciences (IATAS).

Question: Where was the 51st International Emmy Awards ceremony held?

Answer: The New York Hilton Midtown in New York City.

Question: What was the purpose of the 51st International Emmy Awards ceremony?

Answer: To recognize the best television programs initially produced and aired outside the United States in 2022.

Question: When were the nominations for the 51st International Emmy Awards announced?

Answer: September 26, 2023.

Below is a paragraph about {topic}. Your task is to formulate a detailed list of questions and corresponding answers that encompass all the information within the paragraph. To ensure clarity, each question should explicitly mention {topic}. Answers should be concise, consisting of a few short phrases separated by commas. For instance:

Paragraph: {paragraph}

Question:

Table 23: The prompt utilized by GPT-4 for building QA datasets for open-ended generation tasks based on the gathered Wiki-Newpages documents.

The prompt utilized by GPT-4 for building QA datasets for natural language inference tasks

Below is a paragraph about Luis Hugo Hernán Palma Pérez. Your task is to formulate a detailed list of natural language inference tasks with questions and corresponding answers based on the paragraph. For instance:

Paragraph: Luis Hugo Hernán Palma Pérez (born November 3, 1958) is a Chilean surgeon and politician, founding member of the Humanist Party of Chile. He is a deputy for the period 2022-2026, after being elected in the 2021 Chilean parliamentary elections.

Question: Based on the paragraph above can we conclude that Luis Hugo Hernán Palma Pérez was born in November.

Options:

- Yes
- It's impossible to say
- No

Answer: Yes

Question: Based on the paragraph above can we conclude that Luis Hugo Hernán Palma Pérez is a deputy for the period 2020-2024.

Options:

- Yes
- It's impossible to say
- No

Answer: No

Question: Based on the paragraph above can we conclude that The Humanist Party of Chile is a political party in Chile.

Options:

- Yes
- It's impossible to say
- No

Answer: Yes

Question: Based on the paragraph above can we conclude that Luis Hugo Hernán Palma Pérez is a dentist.

Options:

- Yes
- It's impossible to say
- No

Answer: No

Question: Based on the paragraph above can we conclude that Luis Hugo Hernán Palma Pérez was elected in the 2021 Chilean parliamentary elections.

Options:

- Yes
- It's impossible to say
- No

Answer: Yes

Below is a paragraph about {topic}. Your task is to formulate a detailed list of natural language inference tasks with questions and corresponding answers based on the paragraph. For instance:

Paragraph: {paragraph}

Question:

Table 24: The prompt utilized by GPT-4 for building QA datasets for natural language inference tasks based on the gathered Wiki-Newpages documents.

The five-shot prompt used for assessing open-ended generation tasks

Question: Which animated film is included in the list of characters in the Zootopia franchise?

Answer: The animated film "Zootopia" (2016).

Question: Who were the coaches in The Voice Generations (Philippine TV series)?

Answer: Billy Crawford, Chito Miranda, Julie Anne San Jose, and Stell of SB19.

Question: Who is Cyrelle Saut?

Answer: A futsal and football player who has been associated with Tuloy Foundation and the Azkals Development team.

Question: What team does the 2023 Southern Miss Golden Eagles football team represent?

Answer: The University of Southern Mississippi.

Question: When was Kenneth Mitchell (basketball) born?

Answer: October 1, 1975.

Table 25: The five-shot prompt used for assessing open-ended generation tasks, which is derived from the gathered Wiki-Newpages-2024-03 documents.

The prompt used by GPT-4 for annotating QA types in the open-ended generation tasks of the Wiki-News-2023-QA datasets

Below is a paragraph along with corresponding question and answer pairs. Your task is to analyze the paragraph and the question-answer pairs by categorizing the type of information they inquire about or provide. Use concise phrases to describe each category. For example:

Paragraph: <Andrew Turner (rugby union, born 2002) - Wikipedia> Andrew Turner (born 16 February 2002) is an English rugby union player, currently playing for the and . His preferred position is prop.

Question: When was Andrew Turner (rugby union, born 2002) born?

Answer: February 16, 2002.

Question: What nationality is Andrew Turner (rugby union, born 2002)?

Answer: English.

Question: What sport does Andrew Turner (rugby union, born 2002) play?

Answer: Rugby union.

Analysis: Types of question-answer pairs: (1) Birth date, (2) Nationality, (3) Sport/Profession.

Types of the paragraph: Biography - Biographical information about Andrew Turner, a rugby union player born in 2002, including his birth date, nationality, sport, and preferred position.

Below is a paragraph along with corresponding question and answer pairs. Your task is to analyze the paragraph and the question-answer pairs by categorizing the type of information they inquire about or provide. Use concise phrases to describe each category. For example:

Paragraph: {paragraph}

{QA}

Analysis:

Table 26: The prompt used by GPT-4 for annotating QA types in the open-ended generation tasks of the Wiki-News-2023-QA datasets.

The prompt utilized by GPT-4o for building QA datasets for open-ended generation tasks

Your task is to rephrase the paragraph below to make it clearer and more concise. Then, create a detailed list of questions and corresponding answers that cover the factual information in the revised content. Answers should be concise, consisting of a few short phrases separated by commas. For example:

Paragraph:

6 VCs explain how startups can capture and defend marketshare in the AI era. Ninety-four percent of business leaders agree AI will be critical to all businesses' success over the next five years, and total global spending on AI is expected to reach \$154 billion by the end of this year, a 27% increase from 2022.

Revised Content:

Six venture capitalists (VCs) explain how startups can capture and defend market share in the AI era. Ninety-four percent of business leaders agree that AI will be critical to the success of all businesses over the next five years. Additionally, total global spending on AI is expected to reach \$154 billion by the end of this year, representing a 27% increase from 2022.

Simple Question-Answering Pairs:

Question: How many VCs explain how startups can capture and defend market share in the AI era?

Answer: Six venture capitalists (VCs).

Question: What percentage of business leaders agree that AI will be critical to the success of all businesses over the next five years?

Answer: Ninety-four percent.

Question: Over what period do business leaders believe AI will be critical to all businesses' success?

Answer: Over the next five years.

Question: How much is the total global spending on AI expected to reach by the end of this year?

Answer: \$154 billion.

Question: By what percentage is the global spending on AI expected to increase from 2022?

Answer: Twenty-seven percent.

Your task is to rephrase the paragraph below to make it clearer and more concise. Then, create a detailed list of simple questions and corresponding answers that cover the information in the revised content. Answers should be concise, consisting of a few short phrases separated by commas. For example:

Paragraph:

{paragraph}

Revised Content:

Table 27: The prompt utilized by GPT-4o for building QA datasets for open-ended generation tasks based on the gathered WebNews documents.

A training document example in the reading-comprehension format

<Robert Anderson (artist) - Wikipedia> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Answer the questions based on the article:

Question: Write a title:

Answer: Robert Anderson (artist)

Question: Highlight the key information within the article:

Answer: United States; American; Alan Greenspan; George W. Bush; Robert Alexander Anderson; 1946

Question: Based on the article above can we conclude that

<Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Options:

- Yes
- It's impossible to say
- No

Answer: Yes

Question: Tell me about Robert Anderson (artist).

Answer: Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Question: Generate a concrete description about Robert Anderson (artist) based on the following keywords:

United States; American; Alan Greenspan; George W. Bush; Robert Alexander Anderson; 1946

Answer: Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Question: <Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American – known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Answer: Portrait artist.

Question: <Robert Anderson (artist)> - (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as designing United States postage stamps.

Options:

- Alan Greenspan
- 1946
- Robert Alexander Anderson
- George W. Bush

Answer: Robert Alexander Anderson

Question: <Robert Anderson (artist)> Robert Alexander Anderson (born 1946) is an American portrait artist known for painting the official portraits of George W. Bush and Alan Greenspan as well as:

Answer: designing United States postage stamps

Table 28: An example of a training document from the Wiki-Newpages-2023-10-Bio train set, presented in a reading-comprehension format.