
RE-Adapt: Reverse Engineered Adaptation of Large Language Models

William Fleshman
Johns Hopkins University
will.fleshman@jhu.edu

Benjamin Van Durme
Johns Hopkins University
vandurme@jhu.edu

Abstract

We introduce RE-Adapt, an approach to fine-tuning large language models on new domains without degrading any pre-existing instruction-tuning. We reverse engineer an adapter which isolates what an instruction-tuned model has learned beyond its corresponding pretrained base model. Importantly, this requires no additional data or training. We can then fine-tune the base model on a new domain and readapt it to instruction following with the reverse engineered adapter. RE-Adapt and our low-rank variant LoRE-Adapt both outperform other methods of fine-tuning, across multiple popular LLMs and datasets, even when the models are used in conjunction with retrieval-augmented generation.

1 Introduction

Large Language Models (LLMs) require a significant investment to develop and train, requiring resources available to only a limited number of organizations. For instance, Meta’s Llama-3 family of models was trained using two custom-built compute clusters, each containing 24,000 high-end GPUs (Meta, 2024). Parameter Efficient Fine Tuning (PEFT) enables resource efficient downstream customization of LLMs for new domains by adjusting a relatively small number of parameters while keeping the majority unchanged. However, an important distinction exists between the types of model used for further fine-tuning. It is common for LLM producers to release two versions of a model, one which is *pretrained* on a general task such as next-token prediction and an *instruct* version which is then continued trained on annotated data to learn how to follow instructions or respond to queries in a preferential manner (Touvron et al., 2023; Jiang et al., 2023; Almazrouei et al., 2023; Banks and Warkentin, 2024). The availability of both versions introduces a choice for organizations wanting to adapt a model to their custom task or domain. While an instruction-tuned model is generally more capable for popular tasks, the majority of data available for additional fine-tuning is unlabeled, lacking the annotations expected from instruct models. This poses a significant problem as annotation by the downstream organization can be too difficult, expensive, or error-prone (Fredriksson et al., 2020; Desmond et al., 2021). Additional fine-tuning can also degrade the performance of the instruction-tuned model outside of the new fine-tuning distribution (Kotha et al., 2024). On the other hand, pretrained models can be easily fine-tuned with unlabeled text but lack the additional capabilities of their instruct counterparts.

To address this dilemma, we seek the ability to fine-tune existing LLMs on unlabeled text while retaining the capabilities from pre-existing instruction-tuning. We draw inspiration from *adapters*, sets of learnable parameters added to an existing model for fine-tuning (Bapna and Firat, 2019; Houlsby et al., 2019). We make the key observation that **the difference in weights between an instruction-tuned and corresponding pretrained model is effectively an adapter**. Isolating the information learned from instruction-tuning into this *Reverse Engineered (RE)-Adapter* enables fine-tuning of the pretrained model, which can then be readapted with the instruction following capabilities (Figure 1). In this work we:

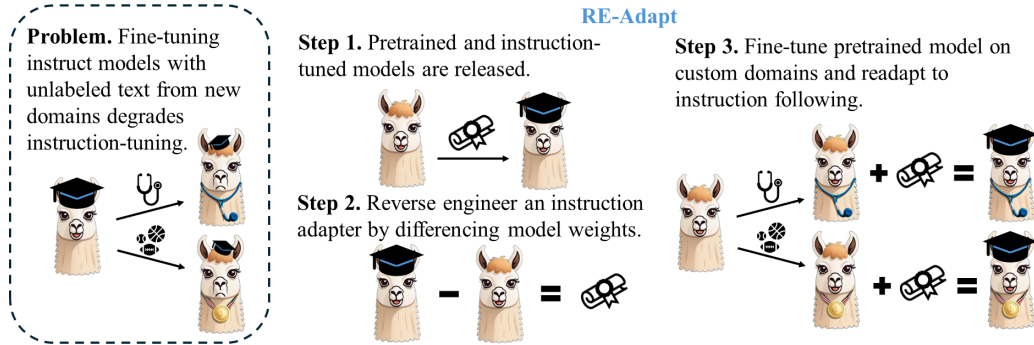


Figure 1: In RE-Adapt, an *instruction adapter* is isolated by differencing weights between instruct (🩺) and pretrained (🎓) versions of a model, which can be reapplied to the pretrained model after fine-tuning.

- Explore the differences in parameters between pretrained and instruct models and their use as instruction adapters;
- Quantify RE-Adapt’s effectiveness to leverage unstructured knowledge for question answering in new domains under both context-free and retrieval-augmented scenarios;
- Introduce *partial adaptation*, a technique for scaling the strength of adapters for fine-grain control of knowledge priorities; and
- Demonstrate that RE-Adapters are *effectively* low-rank, showing that low-rank RE-Adapters (LoRE-Adapters) are capable of similar performance using up to 5x fewer parameters.

2 Background

2.1 Adapters

Adapters (Bapna and Firat, 2019; Houlsby et al., 2019) have played an important role in the context of transfer learning for language models in recent years, particularly for fine-tuning pretrained models which are too large to fully train on commodity hardware. The concept introduced by Houlsby et al. (2019) provides a lightweight alternative to full fine-tuning through the augmentation of models with small modular sets of trainable parameters. Adapters have been useful for enabling the use of pretrained models on new tasks (Pfeiffer et al., 2021; Karimi Mahabadi et al., 2021; Rücklé et al., 2021), new domains (Malik et al., 2023; Schopf et al., 2023; Diao et al., 2023), and adapting to multiple languages (Chronopoulou et al., 2023b; Üstün et al., 2022; Parovic et al., 2023).

Low-Rank Adapters (LoRA) (Hu et al., 2022) are a particularly parameter efficient adaptation technique which adds a low-rank matrix to the weights of existing layers. Because the adapter is low-rank it can be represented as the product of two much smaller matrices, significantly lowering the number of trainable parameters. Weight-Decomposed Low-Rank Adaptation (DoRA) is an extension to LoRA with superior performance and similar efficiency (Liu et al., 2024). Liu et al. (2024) achieve this by decomposing the pretrained weights into both magnitude and direction components, applying LoRA for directional fine-tuning only. Important to this work, adapters learned with either LoRA or DoRA can be represented as a single matrix which captures the information learned during fine-tuning. The pretrained model is then adapted by simply adding the new matrix to the existing weights. We leverage DoRA to fine-tune our models on a new domain, and take inspiration from the additive nature of these techniques to derive our reverse engineered adapters.

Several approaches have been developed which utilize the mixing or combination of adapters to benefit from diverse tasks or domains Pfeiffer et al. (2021); Rücklé et al. (2021); Wang et al. (2022); Chronopoulou et al. (2023a); Fleshman et al. (2024); Zadouri et al. (2024) or for parameter efficient federated learning (Babakniya et al., 2023; Sun et al., 2024). One method to categorize these approaches is by the mechanism used for combining the adapters. Either a weighted combination of adapters is applied to the base model (Chronopoulou et al., 2023a; Fleshman et al., 2024; Babakniya et al., 2023; Sun et al., 2024) or another set of parameters are used to learn adapter interactions

(Pfeiffer et al., 2021; Rücklé et al., 2021; Wang et al., 2022; Zadouri et al., 2024). We focus on the former, as we reframe instruction-tuned models as the summation of a pretrained model with an instruction adapter. We add new knowledge by combining domain-specific and instruction adapters via linear combination. As highlighted by Sun et al. (2024), separate adapters can be incompatible when averaged. Chronopoulou et al. (2023a) and Fleshman et al. (2024) try to mitigate this by initializing adapters with the same random weights, and Sun et al. (2024) by doing the same through a data driven approach. Neither option is applicable here, as we have no control over the instruction adapter. This motivates our new approach for *partial adaptation* which we introduce in Section 3.

2.2 Instruct Models

Some of the most capable LLMs are *instruct* variants, pretrained on massive amounts of unannotated text and further trained on curated datasets with a combination of instruction-tuning (Mishra et al., 2022; Wei et al., 2022; Ouyang et al., 2022; Sanh et al., 2022) and Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017; Stiennon et al., 2020). For example, Llama-3 was pretrained on 15T tokens and the instruct version continued training with a combination of supervised fine tuning (SFT), rejection sampling, proximal policy optimization (PPO), and direct preference optimization (DPO) (Meta, 2024). Open-source LLM producers generally release both the instruct versions as well as the pretrained models from which they were derived (Jiang et al., 2023; Almazrouei et al., 2023; Banks and Warkentin, 2024; Meta, 2024). Access to the pretrained LLM allows users to customize the model to a new task or domain while taking advantage of the large investment required for pretraining. Fine-tuning the instruct model directly is generally avoided due to *catastrophic-forgetting*, a phenomenon where models lose previous abilities with subsequent rounds of continued training (McCloskey and Cohen, 1989; Kotha et al., 2024). This is unfortunate, as few organizations have the resources to conduct fine-tuning at the scale of the original instruction-tuned models. In this work, we explore methods of fine-tuning LLMs which take advantage of both the pretraining and instruction-tuning of existing LLMs. We specifically design our approach to minimize forgetting while fine-tuning instruction-capable models with unlabeled text.

2.3 Model Arithmetic

Previous works have looked at the ability to arithmetically manipulate models to isolate certain behaviors (Ilharco et al., 2023; Mitchell et al., 2024). Ilharco et al. (2023) constructed *task vectors* by differencing weights between a pretrained model and several corresponding models each fine-tuned for a particular task. They observed for their models that task vectors are almost orthogonal to each other, preventing interference and allowing combinations of the vectors for negating certain behaviors, improving multi-task performance, or performing well on new tasks via more complicated task analogies (Ilharco et al., 2023). We similarly solve for our reverse engineered adapter with a simple differencing, but using a single LLM fine-tuned for multi-task instruction-following. By effectively isolating instruction-tuning into an adapter, we allow further fine-tuning of pretrained models, maximizing knowledge acquisition before readapting their ability to follow instructions. We introduce an optional step for reducing the rank of our RE-Adapter, lowering memory requirements while maintaining or improving performance in some scenarios. Unlike individual task vectors, our multi-purpose RE-Adapters are not assumed to be orthogonal to new training domains. We introduce a technique for mitigating potential interference in Section 3 by controlling the adaptation strength.

Mitchell et al. (2024) developed an alternative approach for isolating pretraining knowledge from fine-tuned behaviors which they call *emulated fine-tuning*. Instead of differencing model weights, emulated fine-tuning considers the difference in outputs between pretrained and fine-tuned versions of a model. By combining this difference with the output of a larger pretrained model, Mitchell et al. (2024) found that they could benefit from the additional pretraining knowledge while still solving the task of the smaller model. Their technique could be extended to meet our goal but requires the storage and forward pass of multiple models for inference. Our approach isolates knowledge and instruction-following into adapters, merged into a single model at no extra cost.

3 Partial Adaptation

We detail our main methods in Section 4, but first we introduce a technique for controlling the strength of adaptation. Consider a model with weights \mathbf{W} and an adapter \mathbf{A} used to fine-tune the model on a

new domain. Using additive adapters such as LoRA or DoRA, the combined weights:

$$\hat{\mathbf{W}} = \mathbf{W} + \mathbf{A} \tag{1}$$

are then used for inference (Hu et al., 2022; Liu et al., 2024). We make the observation that the resulting model assigns equal weight to the original parameters and the new adapter, which is generally trained with significantly less data than the original weights. This potentially leads to overfitting in the new domain and degradation of performance in the general setting. These issues compound in situations where multiple adapters are combined. Both Chronopoulou et al. (2023a) and Fleshman et al. (2024) discuss complications arising from mixing adapters, especially if they were not initialized with the same values to encourage compatibility.

To mitigate these challenges we propose a technique for *partial adaptation* which introduces a post-hoc scaling factor for each fine-tuned adapter. Importantly, Equation 1 is still used during fine-tuning, but inference becomes:

$$\hat{\mathbf{W}} = \mathbf{W} + \lambda \mathbf{A} \tag{2}$$

where $0 \leq \lambda \leq 1$ is used to scale the strength of adaptation. In our experiments, we find that partial adaptation improves performance when using either single or multiple combined adapters.

4 Reverse Engineered Adaptation

Here we describe Reverse Engineered Adaptation (RE-Adapt), our approach to solve the challenge of updating an instruction-tuned model with unlabeled text without degrading the ability of the model to follow instructions. In Section 5, we demonstrate the effectiveness of this approach for closed-book and retrieval-augmented question answering.

4.1 RE-Adapters

First consider two language models: \mathbf{T}_Φ , which has been pretrained with parameters Φ ; and \mathbf{T}_Θ , having the same architecture as \mathbf{T}_Φ but with parameters Θ updated from the pretrained parameters Φ via instruction-tuning. Given these models, we can solve for the RE-Adapter parameters Δ using:

$$\Delta = \Theta - \Phi \tag{3}$$

to isolate the information learned during instruction-tuning. Next, we augment the pretrained model \mathbf{T}_Φ with a learnable adapter Ψ and fit $\mathbf{T}_{\Phi+\Psi}$ on a new domain by only updating the adapter weights Ψ . We refer to Ψ as the *knowledge adapter*. We utilize DoRA to fit Ψ in our experiments, but any fine-tuning approach is applicable. We construct our final model \mathbf{T}_Ω with parameters:

$$\Omega = \Phi + \alpha \Psi + \beta \Delta \tag{4}$$

where α and β are the scaling factors for the partial adaptation of Ψ and Δ respectively. We find that scaling down the strength of the knowledge adapter Ψ and RE-Adapter Δ with partial adaptation leads to better performance in instruction-based tasks related to the new domain while maintaining or slightly improving on the performance of the original instruction-tuned model out-of-domain.

4.2 LoRE-Adapters

Inspired by LoRA, we explore the intrinsic dimensionality of RE-Adapters and their ability to be represented by low-rank approximations. The Eckart-Young-Mirsky theorem establishes the truncated singular value decomposition (SVD) as the best low-rank approximation of matrices under the Frobenius norm (Eckart and Young, 1936). We compute the SVD of the RE-Adapter Δ from Equation 3 which yields $\Delta = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ with the diagonal of \mathbf{S} containing the singular values of Δ sorted by magnitude, with \mathbf{U} and \mathbf{V} the corresponding left and right singular vectors. We then compute the percentage of variance explained by each dimension by squaring the singular values and re-normalizing the results to sum to 1. The cumulative explained variance v at rank k is then:

$$v_k = \sum_{i=0}^k \frac{\sigma_i^2}{\sum \sigma_j^2} \tag{5}$$

where σ_i is the i th singular value. We replicate this analysis for multiple modern LLMs and find that the majority of total variation in parameters can be represented at low-rank. For example, Figure 2

displays the cumulative explained variance plots for three layers from the RE-Adapter derived from Llama-3: we see more than half of the variance in these layers can be captured by a rank 128 approximation. This suggests the potential for a low-rank RE-Adapter (LoRE-Adapter).

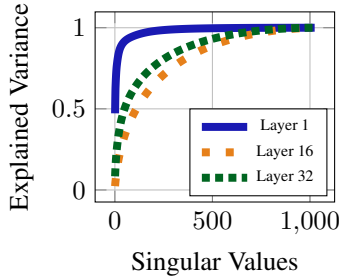


Figure 2: Cumulative explained variance for singular values from Llama-3 RE-Adapt k_{proj} layers.

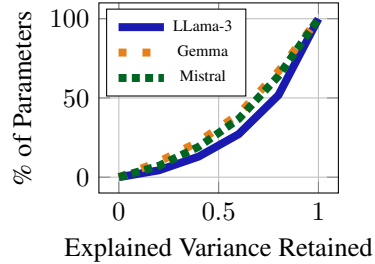


Figure 3: Percent of original model's parameter count used for LoRE-Adapt on Llama-3 with varying threshold of explained variance τ .

We can convert a RE-Adapter into a LoRE-Adapter using a similar approach as Sharma et al. (2024) by representing each layer with its truncated SVD. In our case, we truncate to the rank that captures a total explained variance above a user-defined threshold τ . Figure 3 shows the relationship between τ and the reduction in total parameters when using Llama-3 models to derive the adapter. As τ increases we maintain a higher percentage of the original parameters. We use LoRE-Adapters with $\tau = 0.5$ for the experiments in this work and see similar or better performance when compared to RE-Adapt while using up to 5x less parameters. Like LoRA, the savings in memory is beneficial in cases where several LoRE-Adapters are swapped in and out of the same model.

5 Experiments

We quantify the effectiveness of RE-Adapt using question answering (QA), a task for which instruction-tuned models should perform significantly better than their pretrained counterparts. Specifically, we want to see if RE-Adapt is better than alternatives for adding knowledge from data not annotated with question-answer pairs. We would like the resulting model to do well answering questions about the new domains, while maintaining the level of performance of the original instruction-tuned model when answering unrelated questions.

5.1 Models

We replicate all experiments using the pretrained and instruct versions from the Gemma-7B (Banks and Warkentin, 2024), Llama-3-8B (Meta, 2024), and Mistral-7B (Jiang et al., 2023) family of LLMs using the HuggingFace API (Wolf et al., 2020). We utilize the parameter efficient fine-tuning library (Mangrulkar et al., 2022) for adding DoRA (Liu et al., 2024) knowledge adapters to each of these models. We perform all fine-tuning and inference with a single 80GB A100 GPU. We include hyper-parameters and other details of our fine-tuning in Appendix A.

In Section 5 we compare RE-Adapt and LoRE-Adapt with the pretrained and instruct models of each family, as well as pretrained and instruct models fine-tuned with DoRA on the new domains. We perform experiments for closed-book QA as well as QA with retrieval-augmented generation (RAG).

5.2 Data

Kotha et al. (2024) showed that fine-tuning degrades performance outside of the fine-tuning distribution. We hypothesize that our approach mitigates this issue by isolating existing instruction-tuning from additional fine-tuning. We test this by measuring the changes in question-answering performance when various fine-tuning strategies are used to update models with unlabeled data. An optimal approach would benefit from the new knowledge when asked related questions, without losing the ability to answer unrelated questions.


			
	Instruct w/out News	Instruct w/ News added	News RE-Adapt
New knowledge: <i>Where was the Greg Mortimer Antarctic Cruise stranded on March 31, 2020?</i>	Antarctica ❌	Uruguay ✅	Uruguay ✅
Pretraining knowledge: <i>How many episodes are there in Dragon Ball Z?</i>	291 ✅	40 ❌	291 ✅

Figure 4: RE-Adapt enables the addition of new knowledge to an instruction-tuned model, without degrading capabilities on knowledge from pretraining.

We explore this hypothesis by fine-tuning models in two different settings. We use English WMT News Crawl (Kocmi et al., 2022) articles published in the year 2020 as our first fine-tuning distribution.¹ These articles provide non-annotated information which we capture through DoRA adapters trained for next-token-prediction. We evaluate how well this knowledge is acquired by using the resulting models to answer related questions from the StreamingQA dataset (Liška et al., 2022), which contains 21,681 QA pairs derived from our subset of articles.²

We use the evidence passages from RetrievalQA (Zhang et al., 2024) as our second fine-tuning distribution and measure performance on the corresponding questions from the same dataset.³ Zhang et al. (2024) curated the dataset by compiling the subset of questions from five other QA benchmarks for which GPT-4 (OpenAI et al., 2024) is unable to answer without access to external knowledge. The questions were selected with the goal of having the corresponding knowledge absent from current LLMs, making this dataset especially challenging in the closed-book setting.

To measure any performance degradation from fine-tuning, we also evaluate our models using a short-answer subset of the Natural Questions dataset (Kwiatkowski et al., 2019) which is unrelated to either fine-tuning distribution.⁴ We use these questions to measure performance before and after fine-tuning our models on the other domains. We would like our approach to result in improved performance when answering questions related to the fine-tuning data without a reduction in performance on the unrelated Natural Questions Figure 4.

5.3 Evaluation

We observe that instruction-tuned models will generally answer questions in long-form, often repeating the question and providing additional helpful context. An example of this behavior is shown in Table 1 where the model is asked for the number of episodes in a popular tv series. Here we see the reference answer is 291, which Llama-3 gets correct, but with a response containing full sentences and additional information to clarify its position.

Table 1: Example from Natural Questions with a truncated response. Llama-3’s full response includes more details per country.

Question	how many episodes are there in dragon ball z?
Answer	291
Llama-3	There are a total of 291 episodes in the original Japanese version of Dragon Ball Z. However, the episode count can vary depending on the version and the country.

Popular QA metrics such as Rouge-L (Lin, 2004) or exact match would penalize Llama-3 for not being precise. To alleviate this concern we evaluate using Rouge-L’s recall, which is the percentage

¹Available at <https://data.statmt.org/news-crawl/README> under CC0 license.

²Available at <https://github.com/google-deepmind/streamingqa> under CC-BY 4.0 license.

³Available at <https://huggingface.co/datasets/zihanz/RetrievalQA> under MIT license.

⁴Available at https://huggingface.co/datasets/natural_questions under CC-BY-SA 3.0 license.

of the longest common sub-sequence of the reference answer found in the model’s response. We additionally measure a version of exact match which looks for the exact reference answer anywhere in the response. In both cases, if the reference answer is in the response the score will be 1. If the answer is partially correct then exact match will be 0, but Rouge-L will provide partial credit.

5.4 Closed-Book QA

In our first experiment we conduct QA evaluation in a closed-book setting where the models must provide an answer given nothing but the question. We explore how RE-Adapt behaves in this setting with varying partial adaptation scaling factors. Figure 5 shows the QA performance of LLama-3 using a fixed-factor of 1.0 for the knowledge-adapter with varying scaling factors for the RE-Adapter. We find that partial adaptation with a factor of 0.5 for both the knowledge adapter and instruction adapter provides robust results across models and datasets when using both RE-Adapt and LoRE-Adapt.

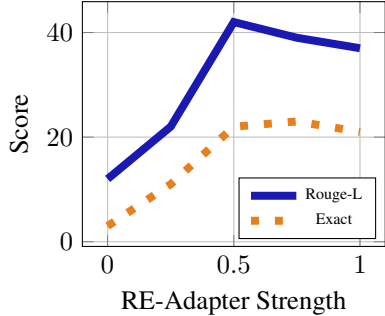


Figure 5: StreamingQA performance as RE-Adapter is added to fine-tuned Llama-3 model with varying strengths.

We use an explained variance threshold $\tau = 0.5$ for our LoRE-Adapters. The resulting percentage of original parameters for each model are: Llama-3 (19.2%), Gemma (30.2%), and Mistral (27.1%).

The closed-book performance of all models across datasets is shown in Table 2. Both RE-Adapt and LoRE-Adapt outperform the pretrained and instruction-tuned models on StreamingQA and RetrievalQA, even when those models are fine-tuned on the corresponding News Crawl or RetrievalQA passages. As expected, the pretrained models perform worse, although fine-tuning on the unlabeled data does improve the QA ability of both pretrained and instruct models in the domain used for adaptation. These in-domain results indicate that our approach is superior for knowledge acquisition. Next we will discuss the impact fine-tuning has on general QA performance by looking at results on the out of domain Natural Questions dataset.

Table 2: Closed-book QA performance. The QA dataset being evaluated is listed above the dataset used for fine-tuning DoRA adapters. R-L indicates Rouge-L and EM indicates exact match.

Model	StreamingQA News Crawl		RetrievalQA RQA Passages		Natural Questions				
	R-L	EM	R-L	EM	News Crawl R-L	News Crawl EM	RQA Passages R-L	RQA Passages EM	
Llama-3	Pretrained	9	0	1	0	10	3	10	3
	Pretrained + DoRA	12	3	3	2	10	4	14	7
	Instruct	33	19	5	3	46	34	46	34
	Instruct + DoRA	38	22	7	4	39	22	37	27
	LoRE-Adapt (Ours)	46	26	10	6	51	34	53	35
	RE-Adapt (Ours)	46	27	9	6	52	34	54	36
Gemma	Pretrained	11	2	1	0	10	3	10	3
	Pretrained + DoRA	19	4	1	0	7	1	10	2
	Instruct	20	9	2	1	26	12	26	12
	Instruct + DoRA	31	18	5	3	26	12	28	14
	LoRE-Adapt (Ours)	31	15	7	4	24	14	30	20
	RE-Adapt (Ours)	33	18	6	4	26	17	28	17
Mistral	Pretrained	17	5	2	0	14	5	14	5
	Pretrained + DoRA	22	8	2	1	14	5	15	6
	Instruct	29	16	4	2	33	22	33	22
	Instruct + DoRA	36	21	6	5	27	13	33	18
	LoRE-Adapt (Ours)	39	24	7	5	39	24	42	28
	RE-Adapt (Ours)	37	22	6	4	37	23	41	27

The closed-book results for the Natural Questions dataset on the right side of Table 2 demonstrate the issues with fine-tuning instruct models with non-annotated data, resulting in models that perform worse in their original setting. While fine-tuning on News Crawl or Retrieval QA passages improved the instruct models on the corresponding QA datasets, the majority of models saw a decrease in performance on Natural Questions. RE-Adapt alleviates this problem by using the data from the new domain to only fine-tune the pretrained model, keeping the instruction-tuning intact. Using our approach,

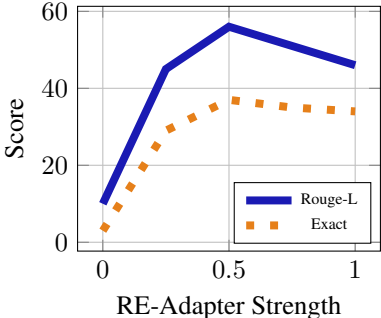


Figure 6: Natural Questions performance as the RE-Adapter is added to pretrained Llama-3 with varying strengths.

the resulting models performed significantly better on the fine-tuning distribution without a performance degradation out-of-domain. In fact, **RE-Adapt and LoRE-Adapt performed better than the original instruction-tuned models out-of-domain**. This improvement indicates that instruction-tuning likely degrades knowledge from pretraining; an issue our approach mitigates through partial adaptation. We confirm this suspicion by applying RE-Adapt to Llama-3 without any additional fine-tuning. This allows us to produce instruct models with instruction-tuning strengths ranging from 0 (the pretrained model) to 1 (the instruct model). We find that **we can improve existing instruct models with zero additional training by simply scaling down the strength of instruction-tuning** Figure 6. Combined, these results demonstrate the effectiveness of RE-Adapt for knowledge acquisition with minimal *forgetting*.

5.5 RE-Adapt with RAG

Retrieval-augmented generation (RAG) Lewis et al. (2020) is a popular alternative for utilizing new data with instruction-tuned models. Instead of altering the model directly, RAG maintains a database of all text and retrieves relevant documents to include in the prompt as context. This begs the question, is RE-Adapt still beneficial if the new data is already available via RAG?

Table 3: QA performance when using RAG with BM25 and (Oracle) retrievers.

	Model	StreamingQA		RetrievalQA	
		Rouge-L	Exact Match	Rouge-L	Exact Match
Llama-3	Pretrained	38 (59)	27 (48)	13 (16)	11 (14)
	Instruct	55 (57)	54 (58)	14 (30)	16 (32)
	LoRE-Adapt (Ours)	69 (74)	58 (64)	24 (37)	21 (31)
	RE-Adapt (Ours)	68 (71)	59 (64)	19 (36)	18 (30)
Gemma	Pretrained	39 (41)	28 (29)	4 (26)	3 (23)
	Instruct	52 (56)	48 (53)	17 (24)	16 (24)
	LoRE-Adapt (Ours)	46 (50)	49 (55)	12 (17)	18 (27)
	RE-Adapt (Ours)	50 (55)	50 (56)	21 (30)	18 (28)
Mistral	Pretrained	33 (38)	26 (30)	18 (12)	16 (10)
	Instruct	49 (52)	50 (56)	14 (23)	19 (28)
	LoRE-Adapt (Ours)	54 (58)	55 (61)	18 (23)	20 (28)
	RE-Adapt (Ours)	55 (58)	55 (60)	15 (24)	20 (29)

To answer this question, we replicate our experiments on StreamingQA and RetrievalQA, using a BM-25 index (Robertson and Zaragoza, 2009) to retrieve the most relevant passage to be used as context for the models. In practice, RAG setups can retrieve more than one document, but each question in our datasets can be answered from a single passage, and therefore we avoid known issues which RAG can face when too much context is provided to the models (Liu et al., 2023; Barnett et al., 2024; Gao et al., 2024). Because a poor retriever could bias results in our favor, we also repeat the experiment using an oracle retriever. Instead of performing a heuristic search, the oracle retriever directly selects the passages capable of answering the question as context. While this idealized

retriever is unrealistic in practice, it allows us to further isolate the benefit of combining RAG with fine-tuning by eliminating any impact from imperfect retrieval.

The RAG results are shown in Table 3. Again we see significant improvements when using RE-Adapt and LoRE-Adapt even in this RAG setting where the model should already have access to the relevant information needed to answer the questions. The BM-25 search retrieved the correct document with approximately 73% accuracy across models. Using RE-Adapt to incorporate the data outside of RAG alleviates the shortcomings of the retriever. However, RE-Adapt also improved results when using the oracle, suggesting that adding domain knowledge with an adapter also reduces incorrect interpretations of the context retrieved via RAG.

6 Discussion

Combined, our results demonstrate RE-Adapt’s effectiveness at incorporating new knowledge into existing LLMs without having to discard previous instruction-tuning. Our methods increase QA performance by a greater amount when compared to traditional fine-tuning strategies. We also find that our approach improves RAG based systems, even in the most optimistic case of perfect retrieval. Our improved results outside of the fine-tuning distribution suggest that we can recover additional pretraining knowledge by reducing the strength of instruction-tuning through partial adaptation. Importantly, an improvement is seen without any additional fine-tuning of the underlying models. These results encourage additional future research into controlling the competing priorities of knowledge acquisition and general problem solving capability.

Limitations. The limitations of our work are two-fold. First, instruction-tuned models perform better than pretrained models on a wide variety of tasks, but we limit our evaluations to the single task of question answering due to the large number of ablations required by our experiments and limited compute resources available. Second, we include the prompts used for instructing the models for QA in Appendix B but note that different prompting strategies could alter our results. We mitigate introducing bias in prompting by not optimizing the prompts for any particular method.

Societal Impact. We are unaware of any negative societal impacts likely to be caused by our contributions. We further amortize the costs of building open-source LLMs by enabling others to leverage existing instruction-tuning, hopefully decreasing the future energy consumption and environmental impacts caused by LLM customization.

7 Conclusion

In this work, we presented RE-Adapt, a new approach for adding knowledge to existing instruction-tuned models. RE-Adapt isolates the differences between an instruction-tuned model and its pretrained counterpart in order to preserve instruction-following capabilities during additional fine-tuning on unlabeled data. We demonstrated that our approach outperforms fine-tuning pretrained or instruction-tuned models directly, which otherwise causes performance to degrade outside of the new fine-tuning domain. Our findings are robust across three state of the art large language models.

We achieved our best performance using *partial adaptation*, a new method for controlling the strength of adaptation at inference time when using single or combined adapters. We found that partially adapting instruction-tuned models improved QA performance without any additional fine-tuning.

We also analyzed the spectrum of RE-Adapt’s weight matrices, constructing a low-rank variant of our approach, LoRE-Adapt, which captures the majority of variation in the instruction-tuning weights at a much lower rank. LoRE-Adapt performed similarly to RE-Adapt with occasional out-performance, while decreasing the number of parameters by as much as 5x in our experiments.

Finally, we demonstrated that RE-Adapt improves performance even when the information required to answer questions is available via retrieval augmented generation. Combined, our results suggest RE-Adapt is an effective approach for infusing new knowledge into already instruction-tuned LLMs.

References

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. The falcon series of open language models. *Preprint*, arXiv:2311.16867.
- Sara Babakniya, Ahmed Roushdy Elkordy, Yahya H. Ezzeldin, Qingfeng Liu, Kee-Bong Song, Mostafa El-Khamy, and Salman Avestimehr. 2023. Slora: Federated parameter efficient fine-tuning of language models. *Preprint*, arXiv:2308.06522.
- Jeanine Banks and Tris Warkentin. 2024. Gemma: Introducing new state-of-the-art open models. *Google*.
- Ankur Bapna and Orhan Firat. 2019. Simple, scalable adaptation for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1538–1548, Hong Kong, China. Association for Computational Linguistics.
- Scott Barnett, Stefanus Kurniawan, Srikanth Thudumu, Zach Brannelly, and Mohamed Abdelrazek. 2024. Seven failure points when engineering a retrieval augmented generation system. *Preprint*, arXiv:2401.05856.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Alexandra Chronopoulou, Matthew Peters, Alexander Fraser, and Jesse Dodge. 2023a. AdapterSoup: Weight averaging to improve generalization of pretrained language models. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2054–2063, Dubrovnik, Croatia. Association for Computational Linguistics.
- Alexandra Chronopoulou, Dario Stojanovski, and Alexander Fraser. 2023b. Language-family adapters for low-resource multilingual neural machine translation. In *Proceedings of the Sixth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2023)*, pages 59–72, Dubrovnik, Croatia. Association for Computational Linguistics.
- Michael Desmond, Evelyn Duesterwald, Kristina Brimijoin, Michelle Brachman, and Qian Pan. 2021. Semi-automated data labeling. In *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*, volume 133 of *Proceedings of Machine Learning Research*, pages 156–169. PMLR.
- Shizhe Diao, Tianyang Xu, Ruijia Xu, Jiawei Wang, and Tong Zhang. 2023. Mixture-of-domain-adapters: Decoupling and injecting domain knowledge to pre-trained language models’ memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5113–5129, Toronto, Canada. Association for Computational Linguistics.
- C. Eckart and G. Young. 1936. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218.
- William Fleshman, Aleem Khan, Marc Marone, and Benjamin Van Durme. 2024. Adapter-swap: Continuous training of llms with data removal and access-control guarantees. *Preprint*, arXiv:2404.08417.
- Teodor Fredriksson, David Issa Mattos, Jan Bosch, and Helena Holmström Olsson. 2020. Data labeling: An empirical investigation into industrial challenges and mitigation strategies. In *Product-Focused Software Process Improvement*, pages 202–216, Cham. Springer International Publishing.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey. *Preprint*, arXiv:2312.10997.

- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *Proceedings of the 11th International Conference on Learning Representations (ICLR 2023)*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. 2021. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 565–576, Online. Association for Computational Linguistics.
- Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Suhas Kotha, Jacob Springer, and Aditi Raghunathan. 2024. Understanding catastrophic forgetting in language models via implicit inference. In *NeurIPS 2023 Workshop on Distribution Shifts: New Frontiers with Foundation Models*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Adam Liška, Tomáš Kočiský, Elena Gribovskaya, Tayfun Terzi, Eren Sezener, Devang Agrawal, Cyprien de Masson d’Autume, Tim Scholtes, Manzil Zaheer, Susannah Young, Ellen Gilsonan-McMahon Sophia Austin, Phil Blunsom, and Angeliki Lazaridou. 2022. Streamingqa: A benchmark for adaptation to new knowledge over time in question answering models. *arXiv preprint arXiv:2205.11388*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: Weight-decomposed low-rank adaptation. *Preprint*, arXiv:2402.09353.

- Bhavivya Malik, Abhinav Ramesh Kashyap, Min-Yen Kan, and Soujanya Poria. 2023. UDAPTER - efficient domain adaptation using adapters. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2249–2263, Dubrovnik, Croatia. Association for Computational Linguistics.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>.
- Michael McCloskey and Neal J. Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation*, 24:109–165.
- Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.
- Eric Mitchell, Rafael Rafailov, Archit Sharma, Chelsea Finn, and Christopher D Manning. 2024. An emulator for fine-tuning large language models using small language models. In *The Twelfth International Conference on Learning Representations*.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,

- Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Marinela Parovic, Alan Ansell, Ivan Vulić, and Anna Korhonen. 2023. Cross-lingual transfer with target language-ready task adapters. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 176–193, Toronto, Canada. Association for Computational Linguistics.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. AdapterFusion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 487–503, Online. Association for Computational Linguistics.
- Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389.
- Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2021. AdapterDrop: On the efficiency of adapters in transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7930–7946, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Tali Bers, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. Multitask prompted training enables zero-shot task generalization. In *Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*.
- Tim Schopf, Dennis N. Schneider, and Florian Matthes. 2023. Efficient domain adaptation of sentence embeddings using adapters. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 1046–1053, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.
- Pratyusha Sharma, Jordan T Ash, and Dipendra Misra. 2024. The truth is in there: Improving reasoning in language models with layer-selective rank reduction. In *Proceedings of the 12th International Conference on Learning Representations (ICLR 2024)*.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. 2020. Learning to summarize from human feedback. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20*, Red Hook, NY, USA. Curran Associates Inc.
- Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. 2024. Improving lora in privacy-preserving federated learning. In *Proceedings of the 12th International Conference on Learning Representations (ICLR 2024)*.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288.

Ahmet Üstün, Arianna Bisazza, Gosse Bouma, and Gertjan van Noord. 2022. UDapter: Typology-based language adapters for multilingual dependency parsing and sequence labeling. *Computational Linguistics*, 48(3):555–592.

Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. 2022. AdaMix: Mixture-of-adaptations for parameter-efficient model tuning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5744–5760, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners. In *Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface’s transformers: State-of-the-art natural language processing. *Preprint*, arXiv:1910.03771.

Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Acyr Locatelli, and Sara Hooker. 2024. Extremely parameter efficient moe for instruction tuning. In *Proceedings of the 12th International Conference on Learning Representations (ICLR 2024)*.

Zihan Zhang, Meng Fang, and Ling Chen. 2024. Retrievalqa: Assessing adaptive retrieval-augmented generation for short-form open-domain question answering. *Preprint*, arXiv:2402.16457.

A Fine-Tuning Details

We include the settings for training our DoRA adapters in Table 4. All adapters were trained on a single NVIDIA A100 GPU with 80GB of memory.

B Prompts Used

Each LLM can use unique prompting roles and tokens when constructing prompts. We utilize the huggingface *tokenizers* library to ensure our prompts follow the correct template.

The Llama-3 instruct models use a combination of system, user, and assistant roles while Gemma and Mistral only use user and assistant. Our prompts were constructed using the following formats:

Llama-3 Closed-Book QA

system: *Answer the following question.*

user: *<question>?*

Llama-3 RAG

Table 4: Training details.

Setting	Value
LoRA Layers	all-linear
LoRA Rank	64
LoRA Alpha	128
LoRA Dropout	0.05
DoRA	True
Batch Size	20
Epochs News Crawl	10
Epochs RetrievalQA	3
Optimizer	AdamW
Learning Rate	0.0002
Schedule	Linear

system: *Answer the following question given this context: <context>.*
user: *<question>?*

Gemma and Mistral Closed-Book QA

user: *<question>?*

Gemma and Mistral RAG

user: *Answer the following question given this context: <context>\nQuestion: <question>?*