

LoRA-Augmented Generation (LAG) for Knowledge-Intensive Language Tasks

William Fleshman, Benjamin Van Durme

Johns Hopkins University
3400 N. Charles Street
Baltimore, MD 21218 USA
will.fleshman@jhu.edu, vandurme@jhu.edu

Abstract

The proliferation of fine-tuned language model experts for specific tasks and domains signals the need for efficient selection and combination methods. We propose LoRA-Augmented Generation (LAG) for leveraging large libraries of knowledge and task-specific LoRA adapters. LAG requires no additional training or access to data, and efficiently filters, retrieves, and applies experts on a per-token and layer basis. We evaluate LAG on various knowledge-intensive tasks, achieving superior performance over existing data-free methods. We explore scenarios where additional data is available, demonstrating LAG’s compatibility with alternative solutions such as retrieval-augmented generation (RAG).

Introduction

Modern language models (LMs) pretrained on large general-purpose text collections have led to rapid gains in task performance. Significant research is now focused on methods for effectively deploying them in novel data domains and for specialized applications (Lewis et al. 2020; Hu et al. 2022; Ilharco et al. 2023; Fleshman and Van Durme 2024; Yang et al. 2025). One prominent approach involves leveraging techniques such as Retrieval-Augmented Generation (RAG), which dynamically injects relevant information from a collection of documents at inference time to guide the model output (Lewis et al. 2020). While highly effective, RAG necessitates the availability and efficient retrieval of these documents during the inference phase.

Alternatively, LMs can be adapted to specific tasks or custom data domains through fine-tuning, a process that adjusts the model parameters using task-specific or domain-specific datasets (Bapna and Firat 2019; Houlsby et al. 2019; Wei et al. 2022; Mangrulkar et al. 2022). A particularly efficient and popular fine-tuning method is Low-Rank Adaptation (LoRA), which introduces a small number of trainable parameters, known as LoRA adapters, alongside the frozen pretrained model (Hu et al. 2022). This approach significantly reduces computational costs and storage requirements compared to full fine-tuning. The success and efficiency of LoRA has led to a rapid proliferation of these adapters, with numerous variations and specialized versions openly shared

and readily accessible in public repositories such as Hugging Face (Wolf et al. 2020). The abundance of these readily available LoRA adapters, each potentially specialized for different tasks, domains, or styles, presents a unique opportunity. However, it also introduces a critical challenge: how to effectively select or combine these diverse adapters at inference time to achieve optimal performance.

Many existing methods for leveraging LoRA adapters rely on access to the original training data corresponding to each adapter or require additional training to merge or select them (Pfeiffer et al. 2021; Wang et al. 2022; Caccia et al. 2023; Ponti et al. 2023; Fleshman et al. 2024; Huang et al. 2024; Zadouri et al. 2024). For example, *Parametric-RAG (PRAG)* uses a RAG-like retrieval mechanism over training documents, but loads an adapter corresponding to the selected document instead of including the content of the document in the prompt (Su et al. 2025). The dependency on data or additional training can be a significant bottleneck, especially in scenarios where the data is proprietary, unavailable, or too large to manage. These issues have led to the recent development of unsupervised routing methods. Arrow routing constructs rank-1 prototypes directly from the LoRA weights and uses them to efficiently select adapters on-the-fly (Ostapenko et al. 2024). Spectral routing (SpectR) improves the accuracy of Arrow at the cost of higher computational complexity (Fleshman and Van Durme 2025).

We introduce LoRA-Augmented Generation (LAG), an approach developed to address the challenge of effectively utilizing existing LoRA adapters without requiring their corresponding data or any additional training. LAG delivers a flexible and efficient mechanism for leveraging the collective knowledge and capabilities embedded within a large set of LoRA adapters, enabling dynamic LoRA selection on a per-token and per-layer basis (Figure 1). Specifically we:

- Develop an efficient LoRA selection and routing procedure to outperform other training-free approaches.
- Leverage a library of over 1000 LoRAs to demonstrate our improved results on knowledge-intensive tasks; and
- Compare and combine LAG with alternative methods for augmenting LMs with additional knowledge.

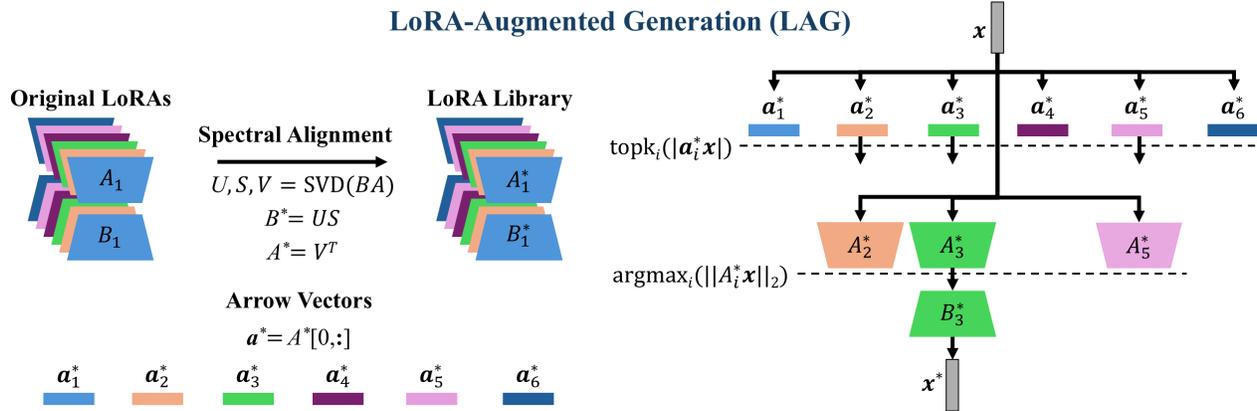


Figure 1: Overview of LAG. LoRA adapters are converted offline via SVD to align representations and extract *arrows*. The token vector x is processed in two stages: (1) *Arrow routing* is used to efficiently filter the large library of adapters to a smaller set of k potential LoRAs, and (2) *Spectral routing* is used to rank the filtered selection by measuring the length of the token representation in the basis of each adapter. The best adapter completes the new token representation x^* .

Background

This section provides an overview of the key concepts and prior works motivating our problem setting and proposed LAG framework. We discuss existing paradigms for incorporating new knowledge into language models, focusing on retrieval-augmented generation and parameter-efficient fine-tuning methods, particularly LoRA. We then delve into recent advancements in knowledge acquisition, adapter retrieval, and unsupervised routing techniques, highlighting the strengths and limitations of each approach.

Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm for grounding LMs with external knowledge, thereby mitigating issues like hallucination and enabling access to up-to-date information (Lewis et al. 2020). The core idea behind RAG is to augment the LM’s input with relevant passages retrieved from a knowledge base. A typical RAG pipeline involves several stages: a retriever component fetches relevant documents or passages based on a query, an optional reranker then reorders these retrieved passages to select the most pertinent ones, and finally, these selected passages are concatenated with the user query and fed into the generative model (Glass et al. 2022). The retriever can be based on sparse methods like BM25 or dense methods using embedding models (Robertson and Zaragoza 2009; Karpukhin et al. 2020; Gao et al. 2024). While highly effective, RAG systems face several challenges. Managing large and dynamic knowledge bases, ensuring retrieval relevance for diverse queries, and fitting retrieved context within the LM’s finite context window are all significant hurdles (Liu et al. 2023; Gao et al. 2024; Barnett et al. 2024). Moreover, the performance of a RAG solution heavily depends on the quality and freshness of the underlying knowledge base, making robust RAG difficult to implement, especially in rapidly evolving domains where the dataset itself is dynamic and constantly updated. Critically, RAG solutions re-

quire the external data to be available and retrieved at inference time, which might not always be feasible. In this work, we explore the case where knowledge is only available parametrically via adapters, but also compare to scenarios where relevant documents can be retrieved via RAG.

Parameter-Efficient Fine-Tuning (PEFT)

Beyond RAG, fine-tuning is another prominent approach to adapt LMs to new tasks or domains. However, full fine-tuning of large models is computationally expensive and memory-intensive, leading to the development of Parameter-Efficient Fine-Tuning (PEFT) methods (Mangrulkar et al. 2022). PEFT techniques aim to update only a small subset of model parameters while keeping the majority of the pre-trained weights frozen, significantly reducing computational cost and storage. Popular PEFT methods include prompt tuning, prefix tuning, and adapter-based approaches (Bapna and Firat 2019; Houlsby et al. 2019; Lester, Al-Rfou, and Constant 2021; Li and Liang 2021).

Among these, Hu et al. (2022)’s Low-Rank Adaptation (LoRA) has gained significant traction due to its effectiveness and simplicity. LoRA fine-tunes LMs by injecting trainable low-rank decomposition matrices into the transformer layers. Specifically, for a weight matrix $W \in \mathbb{R}^{m \times n}$, LoRA introduces two smaller matrices $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$ where $r \ll \min(m, n)$ is the LoRA rank such that the update to the original weight matrix is represented as $W + BA$. During training, W is frozen, and only A and B are optimized. At inference, the full matrix $W + BA$ can be computed and used, or A and B can be dynamically loaded and applied (Mangrulkar et al. 2022). This approach dramatically reduces the number of trainable parameters, enables faster training, and allows for the easy storage and swapping of multiple adapters for different tasks or knowledge domains without requiring modifications to the base model (Hu et al. 2022; Chronopoulou et al. 2023; Fleshman et al. 2024). *Task-vectors* isolate the difference in weights regardless of fine-tuning strategy (Ilharco et al. 2023), and Fleshman and

Van Durme (2024) demonstrates that the SVD can be used to convert these differences into LoRA-like adapters, enabling broader applicability for our LoRA-specific approach.

Knowledge Acquisition beyond RAG

While RAG augments models with external information at inference, several methods have been developed to embed new knowledge into existing language models (Lu, Dou, and Nguyen 2021; Zhang et al. 2023; Wang et al. 2024; Yang et al. 2025; Caccia et al. 2025). These approaches often involve training the model to internalize specific facts or domains. For instance, *Synthetic continued pretraining* constructs a knowledge graph from target data and uses LMs to generate additional training data by sampling relationships from the graph (Yang et al. 2025). Allen-Zhu and Li (2024) find that data augmentation is essential for LMs to extract knowledge, a finding which adapter-based methods have also leveraged (Caccia et al. 2025; Su et al. 2025).

Adapter Retrieval

The proliferation of LoRA adapters necessitates effective methods for selecting or combining them at inference time, especially when dealing with a vast library of adapters each specialized for a particular task or knowledge domain. Recent works leverage data associated with adapters to aid in retrieval (Zhao et al. 2024; Su et al. 2025). *Parametric-RAG (PRAG)* proposes selecting knowledge adapters based on the similarity between a query and the training data used for each adapter (Su et al. 2025). They employ BM25 over training documents and match the selected document to the associated adapter. Like traditional RAG, this approach can be effective for conditioning models on external knowledge if the corresponding data is available. Similarly, *LoRARetriever* constructs task representations using examples from the training data of each adapter and then trains a retriever to select the most appropriate task adapter for a given query (Zhao et al. 2024). While promising, both PRAG and LoRARetriever leverage the associated training data. In this work, we focus on data- and training-free adapter retrieval.

Unsupervised Adapter Routing

To overcome the data and training dependencies of adapter retrieval methods, recent research has explored unsupervised routing techniques that leverage the inherent properties of the adapter weights themselves (Ostapenko et al. 2024; Fleshman and Van Durme 2025). *Arrow routing* is an efficient method that projects adapter weights into a single principal component. The routing decision is then made based on the similarity of the query embedding to these one-dimensional projections. While computationally efficient, relying on a single dimension can lead to inaccuracies due to the inherent loss of information from higher dimensions. *Spectral routing (SpectR)* addresses this issue by utilizing the full spectral properties of the adapter weights, providing a more accurate representation for routing decisions (Fleshman and Van Durme 2025). However, this increased accuracy comes at the cost of higher computational complexity, as it involves working with the full rank of

the adapters. This extra cost is compounded by the number of available adapters, making SpectR untenable with large LoRA libraries. LAG leverages the complementary strengths of these two methods, yielding efficient performance with LoRA libraries too large for SpectR processing.

Problem Setting

The number of pretrained LMs is increasing, and so are repositories of specialized LoRA adapters. These are often fine-tuned for specific tasks or imbued with domain-specific knowledge. This motivates a need for methods to select and apply these adapters at inference time. Specifically, we explore the setting where a large library of *knowledge adapters* is available, each adapter trained on a specific body of knowledge. Likewise, we assume a library of *task adapters*, where each adapter is trained for a knowledge-intensive task. Our core objective is to develop an approach for selecting and applying the knowledge and task adapters most relevant to a specific query without access to the corresponding data or additional training for learning to route.

LoRA-Augmented Generation

Our method introduces an efficient framework for leveraging large libraries of existing LoRA adapters during LM inference, which we term LoRA-Augmented Generation (LAG). We assume access to two distinct sets of LoRA adapters trained for different purposes: (1) a task library \mathcal{T} comprising adapters specialized for specific tasks (fact checking, question answering, etc.), and (2) a knowledge adapter library \mathcal{K} containing LoRAs that encode domain or document-specific information (e.g., Wikipedia articles).

To perform inference with a base LM augmented by these adapter libraries, we propose a two-stage routing strategy for dynamically selecting adapters on a per-token and per-layer basis. Our approach addresses the challenge of scaling inference to thousands of available adapters while maintaining computational efficiency and task-specific performance.

Spectral Alignment

First, we perform offline processing of our adapter libraries following Fleshman and Van Durme (2025)’s *spectral alignment* procedure. For LoRA weight matrices A and B with rank r , we calculate the rank- r singular value decomposition (SVD) of the matrix product BA :

$$U, S, V = \text{SVD}_r(BA), \quad (1)$$

with $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$ the left and right singular vectors and $S \in \mathbb{R}^{r \times r}$ the diagonal matrix of singular values such that:

$$USV^T = BA. \quad (2)$$

The adapter is stored as new matrices B^* and A^* where:

$$B^* = US, \text{ and} \quad (3)$$

$$A^* = V^T. \quad (4)$$

Importantly, the SVD is performed once offline and results in an adapter of equivalent size and function. Ostapenko et al. (2024) store adapters in their original form but use a^* ,

the row vector from A^* associated with the largest singular value, as the adapter prototype for their Arrow routing algorithm. The matrix A^* contains the eigenvectors of the covariance matrix of the LoRA parameters, which represent orthogonal directions of maximum variation induced by the adapter in the space of input vectors $\mathbf{x} \in \mathbb{R}^n$ (Ostapenko et al. 2024; Fleshman and Van Durme 2025). The vector \mathbf{a}^* is the direction capturing the most variation and produces the largest adapter activations of all unit-length input vectors:

$$\mathbf{a}^* = \operatorname{argmax}_{\mathbf{x}, \|\mathbf{x}\|_2=1} \|B A \mathbf{x}\|_2. \quad (5)$$

SpectR improves routing accuracy at the cost of additional computation by using the entire A^* matrix for routing decisions (Fleshman and Van Durme 2025). We leverage both representations in LAG and perform adapter routing in two stages: Arrow-based retrieval and SpectR reranking.

Arrow Retrieval

Once adapters are aligned, the new representation can be leveraged to perform routing without additional training or data access. For a specific LM layer l , there may be associated adapters from one or both of the adapter libraries \mathcal{K} and \mathcal{T} . If both libraries apply, routing is performed separately for task adapters and knowledge adapters, and the selected adapter from each set is used. SpectR becomes infeasible with a large adapter library due to the memory requirements of storing and computing with the full matrices (Fleshman and Van Durme 2025). We therefore perform an efficient first pass retrieval to select the k adapters \mathcal{A} producing the largest magnitude product between their associated arrow vector and the input vector \mathbf{x} :

$$\mathcal{A} = \operatorname{topk}_i |\mathbf{a}_i^* \mathbf{x}|. \quad (6)$$

The parameter k can be chosen based on memory or time budgets, with a lower k being more efficient and a higher k reducing the impact of inaccurate arrow rankings.

SpectR Reranking

We use SpectR scores on the subset of retrieved adapters \mathcal{A} , to more accurately compare and select the adapter most aligned with the input vector \mathbf{x} :

$$\hat{A} = \operatorname{argmax}_{A^* \in \mathcal{A}} \|A^* \mathbf{x}\|_2. \quad (7)$$

The selected adapter (\hat{B} , \hat{A}) is then used with the target layer weights W_l to produce the output representation:

$$\mathbf{h} = W_l \mathbf{x} + \hat{B}(\hat{A} \mathbf{x}), \quad (8)$$

with the low-rank vector $\hat{A} \mathbf{x}$ precomputed and reused from Equation 7. This two-stage process is used to select a single adapter for each position in the sequence. Figure 2 broadly compares LAG with the mechanisms of RAG and PRAG.

Experiments

In this section, we experiment in practical scenarios where a large number of adapters are available for a diverse body of knowledge and tasks. We compare the theoretical requirements for each of the training and data-free approaches and

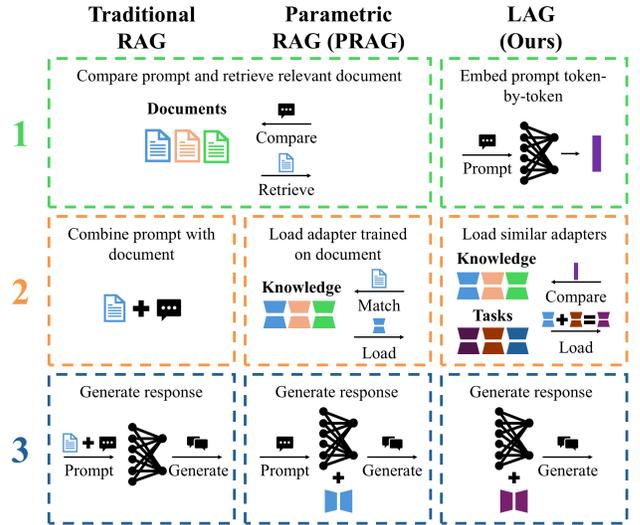


Figure 2: RAG and PRAG both identify the document most similar to the prompt. RAG includes the document in the prompt, while PRAG loads the corresponding adapter. LAG uses internal representations of the prompt to select and load knowledge and task adapters during generation.

confirm that SpectR is computationally infeasible for the large adapter library in our experiments. We explore our core setting with zero access to additional data or training, empirically demonstrating the benefits of our approach. We measure the impact of more aggressive filtering on LAG’s performance, and compare LAG with alternative techniques for leveraging the data associated with our knowledge library.

Data and Metrics

We leverage the KILT benchmark, a set of *knowledge-intensive language task* datasets (Petroni et al. 2021). KILT contains five tasks: fact checking, entity linking, slot filling, question answering (QA), and dialog generation (chat); all of which are grounded in knowledge from a shared collection of Wikipedia articles (Petroni et al. 2021). No existing LoRAs correspond to these articles, so we must train a large adapter library ourselves. We choose a library of size $n = 1000$, which is 2 orders of magnitude larger than previous SpectR experiments, yet still manageable to train on a single GPU. We filter the Wikipedia pages to the top 1000 referenced as provenance across the benchmark tasks. We use this set as our knowledge base and filter the task-specific datasets to only those examples grounded in these articles. Two datasets contain too few samples from the resulting selection and are dropped. Our final evaluation dataset is comprised of **fact checking**: FEVER (Thorne et al. 2011); **entity linking**: AIDA CoNLL-YAGO (Hof-fart et al. 2011), WNED-WIKI, and WNED-CWEB (Alani, Guo, and Barbosa 2018); **slot filling**: Zero Shot RE (Levy et al. 2017) and T-REx (Elsahar et al. 2018); **QA**: Natural Questions (Kwiatkowski et al. 2019) and TriviaQA (Joshi et al. 2017); and **chat**: Wizard of Wikipedia (Dinan et al. 2019). KILT prescribes different evaluation metrics depend-

Dataset	Task	Metric	Size
WoW	Chat	F1	28546
FEV	Fact	Acc	24886
AY2	Link	Acc	9056
WnCW	Link	Acc	995
WnWi	Link	Acc	237
zsRE	Slot	Acc	116
T-REx	Slot	Acc	280
NQ	QA	EM	3855
TQA	QA	EM	4583

Table 1: Filtered dataset information organized by task.

ing on the dataset and task being evaluated (Petroni et al. 2021). These include *Accuracy* for tasks with a discrete output (fact checking, entity linking, and slot filling), *Exact Match* for extractive QA, and *F1 Score* for chat. We also introduce a normalized performance metric to more easily compare across tasks and to control for the differences in difficulty between datasets. Let $f_D(M)$ represent the score for the dataset-specific metric of model M evaluated on a dataset D . We produce a normalized task score S_T using the performance of M compared to a strong reference model R across all datasets in the same task T :

$$S_T = \sum_{D \in T} \frac{|D| f_D(M)}{|T| f_D(R)}, \quad (9)$$

where $|\cdot|$ is the number of samples in D or T . S_T represents the average percentage of the reference model performance on task T achieved by the model under evaluation. We refer to our reference model as the *Oracle* model, which is the LM augmented with the ground-truth knowledge and task adapters for each query. For example, when answering a question about the first U.S. President, the Oracle model would use the *QA* and *George Washington* LoRA adapters. Table 1 includes a summary of the datasets used with their associated task, metric, and number of samples.

Model and Adapter Libraries

We use the Llama-3.2-3B-Instruct model as the LM in our experiments due to its generally good performance and efficient size amenable to running experiments using 1000s of LoRAs on a single GPU (Grattafiori et al. 2024).

Evaluating LAG requires both task and knowledge LoRA libraries corresponding to our evaluation data. We therefore construct libraries with a LoRA adapter for each Wikipedia article and task. No hyperparameter tuning is performed since LAG is designed to work with externally sourced adapters. See Appendix for training and adapter details.

Task Library For each of the original five tasks, we fit a task-specific adapter using samples with provenance outside of our selected Wikipedia articles. The training data contains only the prompt and answer, not the corresponding articles.

Knowledge Library We use synthetic continued pretraining to fit our knowledge adapters in a task-agnostic manner (Yang et al. 2025). We train one adapter per Wikipedia

	Arrow	SpectR	LAG
Disk	$2nhr + nh$	$2nhr$	$2nhr$
GPU	$2khr + nh$	$2nhr$	$2khr + nh$
FLOPs	$2nh$	$2nhr$	$2h(n + rk)$

Table 2: Overview of approximate FLOPs and best-case storage requirements in terms of parameters needed on disk and GPU assuming a library of n rank- r adapters with a hidden dimension of h and top- k filtering. LAG inherits the best between Arrow and SpectR for storage and is on par with Arrow in terms of computational efficiency.

article. The KILT representation of each article includes a set of *anchor* elements corresponding to hyperlinked entities contained in the text (Petroni et al. 2021). The base LM is prompted to rewrite each article once per entity, emphasizing the focused entity in relation to the article contents. The original article is combined with these synthetically generated documents to increase the amount of training data used to fit each LoRA adapter. We train our knowledge adapters using the pretrained version of the LM to mitigate negatively impacting the existing instruction-tuning (Fleshman and Van Durme 2024). The pretrained model is only used during LoRA training, and the adapters are applied to the instruction-tuned model during evaluation.

Theoretical Efficiency

Before empirically demonstrating the benefits of our approach, we discuss the expected efficiency gains of LAG in terms of required disk space, GPU memory, and computation. We summarize these comparisons in Table 2.

Disk Space Given a linear layer with input and output dimension h and a library of n adapters with rank r : Arrow, SpectR, and LAG all require the storage of $2nhr$ LoRA parameters. Each of the n adapters is composed of two $h \times r$ matrices. Arrow routing uses adapters in their original form and requires an additional nh parameters to store the arrow vector from each adapter in the library (Ostapenko et al. 2024). LAG also uses arrow vectors, but because the adapters are stored in their aligned representation, the arrow vectors are captured by the first row of the A^* matrix from Equation 4, preventing the need for extra storage. Like LAG, SpectR only requires the aligned adapter library as it performs routing using the A^* matrix instead of using arrow vectors (Fleshman and Van Durme 2025).

GPU Memory Depending on available memory, all the parameters can be loaded onto the GPU to eliminate data transfer. If memory is limited, Arrow and LAG can take advantage of their ability to filter down to only $k \ll n$ adapters per token using arrow vectors and load only the selected LoRAs into memory. For a token sequence of length s , the nh arrow parameters are used to choose the sk adapters for the sequence (Ostapenko et al. 2024). In the worst case, all n adapters are required in memory if the sequence is longer than n/k tokens and positions in the sequence share little overlap in the relevant knowledge or tasks needed. However, Fleshman and Van Durme (2025) demon-

	Instr	Arrow	LAG	Oracle
WoW	12.5	17.9	17.7	18.1
FEV	66.0	89.2	89.6	90.2
AY2	21.0	45.2	62.4	68.2
WnCw	15.8	23.1	47.2	67.6
WnWi	13.5	32.5	61.2	73.0
zsRE	24.1	42.2	44.8	49.1
T-REx	18.9	45.4	46.1	53.6
NQ	27.0	26.1	30.9	38.8
TQA	50.8	44.2	46.5	50.8

Table 3: LAG outperforms other data and training-free approaches for the majority of the evaluated datasets.

	Chat	Fact	Link	Slot	QA	AVG
Instr	69.3	73.2	29.8	39.4	86.0	59.5
Arrow	99.2	99.0	62.6	85.0	78.0	84.8
LAG	98.2	99.4	89.2	87.5	86.0	92.1

Table 4: LAG produces an average gain of 7 points in normalized performance over Arrow across all tasks.

strated the intuition that sequences do result in high overlap, which reduces the number of adapters needed in memory. In the best case, the same k adapters would be selected for all tokens, resulting in a GPU storage requirement of $nh + 2khr$ parameters for both Arrow and LAG¹. In the unlikely worst case, the GPU requirements are the same as the disk space.

Computation We compare the additional floating-point operations (FLOPs) required, assuming each method selects a single adapter per-token. SpectR requires n matrix-vector products using the A^* matrices of each adapter and an additional product using the B^* matrix from the chosen adapter (Fleshman and Van Durme 2025). This results in $2nh + 2hr = 2hr(n + 1) \approx 2nh$ FLOPs. Arrow is the most efficient method, requiring only $2nh$ FLOPs to choose an adapter using the n arrow vector dot products (Ostapenko et al. 2024). Another $2hr$ FLOPs are used to multiply by each of the LoRA matrices of the chosen adapter, incurring a total of $2nh + 2hr + 2hr = 2h(n + 2r) \approx 2nh$ FLOPs. LAG inherits the same $2nh$ FLOPs from Arrow to choose the top- k adapters and then requires $2khr$ more to process those k LoRAs using SpectR, a total of $2nh + 2khr = 2h(n + rk)$ FLOPs. In our experiments, $rk \ll n$, yielding computation requirements for LAG similar to Arrow. Next, we use these efficiency gains for our experiments in scenarios where the LoRA library is too large for SpectR.

Data and Training Free

We explore our core objective of leveraging large LoRA libraries in a scenario where access to the corresponding adapter data or additional training for learning to route is unavailable. Following Table 2, SpectR is roughly $8\times$ more expensive per-token than Arrow and LAG using the

¹E.g. 1 million rank-6 adapters with $k = 20$ and $h = 4096$ would require 49B extra parameters w/ SpectR versus 4B w/ LAG.

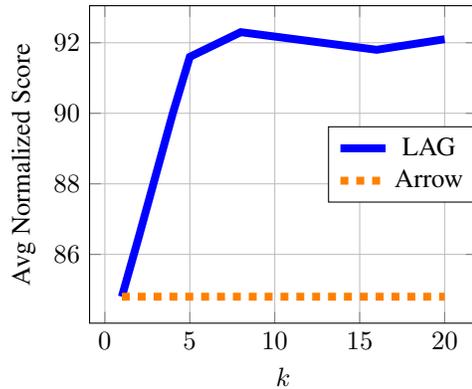


Figure 3: The average normalized performance of LAG goes up and plateaus as k increases. LAG is equivalent to Arrow routing at $k = 1$ and would be equivalent to SpectR if k could be set to the total number of adapters in the library.

rank-8 LoRAs in our library. Indeed, we found SpectR to be intractable for our evaluation using our library of 1000 adapters, exactly the issue LAG was developed to address. Therefore, we first compare LAG against Arrow routing and the instruction-tuned model without any adapters. We also report the performance of the Oracle model, using the ground-truth knowledge and task adapters applied to each sample. The Oracle provides a reference for the model’s achievable performance. For LAG, we use $k = 20$ for filtering, which reduces the number of adapters considered by SpectR by 98%. We later explore the impact of k on our results. We select a single adapter to apply per-token and per-layer for both Arrow and LAG.

Table 3 displays the result for each dataset using the associated KILT metric. The increase in performance when using the Oracle model provides insight into how much the knowledge and task adapters improve performance over the instruct model. The Oracle model performs much better than the base instruct model across all tasks except for QA, where the Oracle performs better on NQ but achieves the same performance on TriviaQA. This could indicate that these datasets were highly represented in the original LM training data or that the instruction-tuning is well-suited for QA in general. LAG consistently outperformed both baselines. The two exceptions were the chat task, where Arrow and LAG both do almost as well as the Oracle model, and on TriviaQA where the Oracle performance suggests the LoRA adapters provide little value. Table 4 provides the normalized performance per-task. LAG captures 92.1% of the Oracle performance on average, scoring 7.3 points higher than Arrow and improving the instruct model by 32.6 points.

Aggressive Filtering

We filtered the LoRAs using $k = 20$ with LAG because using SpectR on the entire adapter library is intractable. The value of k modulates the trade-off between Arrow efficiency and the superior accuracy of SpectR at a higher computational cost. With $k = 1$, LAG becomes equivalent to Arrow routing because the SpectR step is forced to choose the

given adapter. As k increases, there are better chances for Arrow to include the appropriate adapters for SpectR reranking. For example, Fleshman and Van Durme (2025) showed Arrow top- k accuracy increasing by more than $3\times$ when going from top-1 to top-4 in their experiments. If k is set to the total number of LoRAs in the library, then LAG is equivalent to SpectR, as no Arrow filtering occurs.

Figure 3 shows the average normalized performance of LAG across tasks as a function of k . The performance rises steeply and then flattens around $k = 5$, suggesting that more aggressive Arrow filtering improves efficiency with minimal sacrifice to downstream performance. In practice, the value of k can be chosen for the given compute budget, or future research could explore setting k dynamically as a form of test-time scaling (Snell et al. 2025).

Knowledge Access

Finally, we loosen our core objective and explore LAG in the case where the documents associated with the knowledge library are available during inference. This scenario allows for a comparison with alternative methods of incorporating knowledge into the LM, such as Retrieval-Augmented Generation (RAG) and Parametric RAG (PRAG) (Lewis et al. 2020; Su et al. 2025). Specifically, we use documents relevant to the query to either augment the prompt (RAG) or to retrieve the adapter trained on the selected document (PRAG). In keeping with Su et al. (2025), we use BM25 (Robertson and Zaragoza 2009) for document retrieval in both cases. We continue to use $k = 20$ for LAG to more easily compare with our previous results.

We include PRAG and RAG as individual baselines using each approach to augment the instruction-tuned model with additional knowledge. We then evaluate combinations using LAG, PRAG, or RAG for knowledge, with LAG selecting the task adapters in all cases. The individual dataset performances are shown in Table 5. All LAG combinations perform better than either PRAG or RAG alone, except on QA where only RAG + LAG outperforms. The best method for incorporating knowledge varied across tasks, with no statistically significant winner. Table 6 reports the normalized task scores where using LAG for both knowledge and tasks performed best on fact checking and entity linking, while PRAG-LAG did best on chat, and RAG-LAG on slot filling and QA. Notably, RAG-LAG achieved a 102.7 normalized score on slot filling, meaning it outperformed the Oracle model by 2.7%. The samples for slot filling are of the form ‘*entity* [SEP] *relationship*’, making BM25 especially effective since the entity and relationship can both be found in the document containing the necessary information. This contrasts with a task like entity linking, where samples can contain arbitrary content with a single entity marked for identification. The majority of the content in each sample can be unrelated to the correct response. LAG’s per-token selection of adapters provides more flexibility in such cases, and LAG did achieve a better score on entity linking in our evaluation.

Conclusion

We introduce LoRA-Augmented Generation (LAG), a completely unsupervised approach for filtering a large LoRA li-

	PRAG	RAG	LAG	P-LAG	R-LAG
WoW	12.7	12.6	17.7	17.9	17.8
FEV	68.4	78.6	89.6	89.5	89.2
AY2	24.6	23.4	62.4	61.8	57.2
WnCw	20.6	16.9	47.2	44.9	32.9
WnWi	19.4	21.1	61.2	62.4	47.3
zsRE	29.3	42.2	44.8	44.8	47.4
T-REx	25.4	41.4	46.1	46.4	56.4
NQ	24.7	32.6	30.9	30.8	35.9
TQA	48.1	48.2	46.5	47.0	49.0

Table 5: Using LAG to incorporate task-capabilities outperforms PRAG or RAG alone, with variation across datasets. P-LAG: PRAG + LAG, R-LAG: RAG + LAG.

	Chat	Fact	Link	Slot	QA	AVG
PRAG	70.2	75.9	35.3	50.9	80.5	62.6
RAG	69.9	87.2	33.2	79.9	89.9	72.0
LAG	98.2	99.4	89.2	87.5	86.0	92.1
P-LAG	99.0	99.3	88.2	88.0	86.4	92.2
R-LAG	98.6	98.9	80.0	102.7	94.7	95.0

Table 6: Normalized performance across tasks. Different LAG combinations perform better across the tasks. Combining RAG with LAG results in the best average performance.

brary and applying the most suitable adapters at test time on a per-token and per-layer basis. LAG addresses scalability issues with previous solutions while maintaining superior downstream task performance over alternative training and data-free approaches. We evaluated LAG on a set of multiple knowledge-intensive tasks, including fact checking, entity linking, slot filling, question answering, and chat. We demonstrated LAG’s ability to select from this diverse set of knowledge and task capabilities and efficiently incorporate the selected adapters into the LM. LAG significantly outperformed Arrow routing with access to the same adapter libraries. We demonstrated that more aggressive arrow filtering can heavily reduce necessary computation with minimal degradation in task performance. We performed additional experimentation using the Wikipedia articles corresponding to our knowledge library. This allowed us to evaluate methods such as Retrieval-Augmented Generation (RAG) and Parametric RAG, and to combine these methods with LAG. Using LAG to incorporate task capabilities significantly outperformed using RAG or PRAG alone, but there was not a statistically significant difference between the combined approaches. The performance varied by task, and we discussed how different tasks might have characteristics that make each of the various approaches for incorporating knowledge with LAG more suitable in certain circumstances.

Overall, LAG successfully unifies the existing approaches in the area of unsupervised adapter routing, incorporating the efficiency of Arrow and the discriminative power of SpectR to provide scalable multi-task performance.

References

- Alani, H.; Guo, Z.; and Barbosa, D. 2018. Robust named entity disambiguation with random walks. *Semant. Web*, 9(4): 459–479.
- Allen-Zhu, Z.; and Li, Y. 2024. Physics of language models: part 3.1, knowledge storage and extraction. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Bapna, A.; and Firat, O. 2019. Simple, Scalable Adaptation for Neural Machine Translation. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 1538–1548. Hong Kong, China: Association for Computational Linguistics.
- Barnett, S.; Kurniawan, S.; Thudumu, S.; Brannelly, Z.; and Abdelrazek, M. 2024. Seven Failure Points When Engineering a Retrieval Augmented Generation System. In *Proceedings of the IEEE/ACM 3rd International Conference on AI Engineering - Software Engineering for AI, CAIN '24*, 194–199. New York, NY, USA: Association for Computing Machinery. ISBN 9798400705915.
- Caccia, L.; Ansell, A.; Ponti, E.; Vulić, I.; and Sordoni, A. 2025. Training Plug-n-Play Knowledge Modules with Deep Context Distillation. arXiv:2503.08727.
- Caccia, L.; Ponti, E.; Su, Z.; Pereira, M.; Roux, N. L.; and Sordoni, A. 2023. Multi-Head Adapter Routing for Cross-Task Generalization. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Chronopoulou, A.; Peters, M.; Fraser, A.; and Dodge, J. 2023. AdapterSoup: Weight Averaging to Improve Generalization of Pretrained Language Models. In Vlachos, A.; and Augenstein, I., eds., *Findings of the Association for Computational Linguistics: EACL 2023*, 2054–2063. Dubrovnik, Croatia: Association for Computational Linguistics.
- Dinan, E.; Roller, S.; Shuster, K.; Fan, A.; Auli, M.; and Weston, J. 2019. Wizard of Wikipedia: Knowledge-Powered Conversational Agents. In *International Conference on Learning Representations*.
- Elsahar, H.; Vougiouklis, P.; Remaci, A.; Gravier, C.; Hare, J.; Laforest, F.; and Simperl, E. 2018. T-REx: A Large Scale Alignment of Natural Language with Knowledge Base Triples. In Calzolari, N.; Choukri, K.; Cieri, C.; Declerck, T.; Goggi, S.; Hasida, K.; Isahara, H.; Maegaard, B.; Mariani, J.; Mazo, H.; Moreno, A.; Odijk, J.; Piperidis, S.; and Tokunaga, T., eds., *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. Miyazaki, Japan: European Language Resources Association (ELRA).
- Fleshman, W.; Khan, A.; Marone, M.; and Durme, B. V. 2024. AdapterSwap: Continuous Training of LLMs with Data Removal and Access-Control Guarantees. In *Proceedings of Conference on Applied Machine Learning in Information Security (CAMLIS) 2024*.
- Fleshman, W.; and Van Durme, B. 2024. RE-Adapt: Reverse Engineered Adaptation of Large Language Models. arXiv:2405.15007.
- Fleshman, W.; and Van Durme, B. 2025. SpectR: Dynamically Composing LM Experts with Spectral Routing. arXiv:2504.03454.
- Gao, Y.; Xiong, Y.; Gao, X.; Jia, K.; Pan, J.; Bi, Y.; Dai, Y.; Sun, J.; Guo, Q.; Wang, M.; and Wang, H. 2024. Retrieval-Augmented Generation for Large Language Models: A Survey. arXiv:2312.10997.
- Geva, M.; Schuster, R.; Berant, J.; and Levy, O. 2021. Transformer Feed-Forward Layers Are Key-Value Memories. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 5484–5495. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Glass, M.; Rossiello, G.; Chowdhury, M. F. M.; Naik, A.; Cai, P.; and Gliozzo, A. 2022. Re2G: Retrieve, Rerank, Generate. In Carpuat, M.; de Marneffe, M.-C.; and Meza Ruiz, I. V., eds., *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2701–2715. Seattle, United States: Association for Computational Linguistics.
- Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; et al. 2024. The Llama 3 Herd of Models. arXiv:2407.21783.
- Hoffart, J.; Yosef, M. A.; Bordino, I.; Fürstenau, H.; Pinkal, M.; Spaniol, M.; Taneva, B.; Thater, S.; and Weikum, G. 2011. Robust Disambiguation of Named Entities in Text. In Barzilay, R.; and Johnson, M., eds., *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, 782–792. Edinburgh, Scotland, UK.: Association for Computational Linguistics.
- Houlsby, N.; Giurghi, A.; Jastrzebski, S.; Morrone, B.; De Laroussilhe, Q.; Gesmundo, A.; Attariyan, M.; and Gelly, S. 2019. Parameter-Efficient Transfer Learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning*.
- Hu, E. J.; yelong shen; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.
- Huang, C.; Liu, Q.; Lin, B. Y.; Pang, T.; Du, C.; and Lin, M. 2024. LoraHub: Efficient Cross-Task Generalization via Dynamic LoRA Composition. In *First Conference on Language Modeling*.
- Ilharco, G.; Ribeiro, M. T.; Wortsman, M.; Schmidt, L.; Hajishirzi, H.; and Farhadi, A. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Joshi, M.; Choi, E.; Weld, D.; and Zettlemoyer, L. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In Barzilay, R.; and Kan, M.-Y., eds., *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1601–1611. Vancouver, Canada: Association for Computational Linguistics.

- Karpukhin, V.; Oguz, B.; Min, S.; Lewis, P.; Wu, L.; Edunov, S.; Chen, D.; and Yih, W.-t. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 6769–6781. Online: Association for Computational Linguistics.
- Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; Toutanova, K.; Jones, L.; Kelcey, M.; Chang, M.-W.; Dai, A. M.; Uszkoreit, J.; Le, Q.; and Petrov, S. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics*, 7: 452–466.
- Lester, B.; Al-Rfou, R.; and Constant, N. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 3045–3059. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Levy, O.; Seo, M.; Choi, E.; and Zettlemoyer, L. 2017. Zero-Shot Relation Extraction via Reading Comprehension. In Levy, R.; and Specia, L., eds., *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, 333–342. Vancouver, Canada: Association for Computational Linguistics.
- Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; Yih, W.-t.; Rocktäschel, T.; Riedel, S.; and Kiela, D. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781713829546.
- Li, X. L.; and Liang, P. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *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)*, 4582–4597. Online: Association for Computational Linguistics.
- Liu, N. F.; Lin, K.; Hewitt, J.; Paranjape, A.; Bevilacqua, M.; Petroni, F.; and Liang, P. 2023. Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics*, 12: 157–173.
- Lu, Q.; Dou, D.; and Nguyen, T. H. 2021. Parameter-Efficient Domain Knowledge Integration from Multiple Sources for Biomedical Pre-trained Language Models. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Findings of the Association for Computational Linguistics: EMNLP 2021*, 3855–3865. Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Mangrulkar, S.; Gugger, S.; Debut, L.; Belkada, Y.; Paul, S.; and Bossan, B. 2022. PEFT: State-of-the-art Parameter-Efficient Fine-Tuning methods. <https://github.com/huggingface/peft>.
- Ostapenko, O.; Su, Z.; Ponti, E. M.; Charlin, L.; Le Roux, N.; Caccia, L.; and Sordoni, A. 2024. Towards modular LLMs by building and reusing a library of LoRAs. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.
- Petroni, F.; Piktus, A.; Fan, A.; Lewis, P.; Yazdani, M.; De Cao, N.; Thorne, J.; Jernite, Y.; Karpukhin, V.; Mailard, J.; Plachouras, V.; Rocktäschel, T.; and Riedel, S. 2021. KILT: a Benchmark for Knowledge Intensive Language Tasks. In Toutanova, K.; Rumshisky, A.; Zettlemoyer, L.; Hakkani-Tur, D.; Beltagy, I.; Bethard, S.; Cotterell, R.; Chakraborty, T.; and Zhou, Y., eds., *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2523–2544. Online: Association for Computational Linguistics.
- Pfeiffer, J.; Kamath, A.; Rücklé, A.; Cho, K.; and Gurevych, I. 2021. AdapterFusion: Non-Destructive Task Composition for Transfer Learning. In Merlo, P.; Tiedemann, J.; and Tsarfaty, R., eds., *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 487–503. Online: Association for Computational Linguistics.
- Ponti, E. M.; Sordoni, A.; Bengio, Y.; and Reddy, S. 2023. Combining Parameter-efficient Modules for Task-level Generalisation. In Vlachos, A.; and Augenstein, I., eds., *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 687–702. Dubrovnik, Croatia: Association for Computational Linguistics.
- Robertson, S.; and Zaragoza, H. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.*, 3(4): 333–389.
- Snell, C. V.; Lee, J.; Xu, K.; and Kumar, A. 2025. Scaling LLM Test-Time Compute Optimally Can be More Effective than Scaling Parameters for Reasoning. In *The Thirteenth International Conference on Learning Representations*.
- Su, W.; Tang, Y.; Ai, Q.; Yan, J.; Wang, C.; Wang, H.; Ye, Z.; Zhou, Y.; and Liu, Y. 2025. Parametric Retrieval Augmented Generation. arXiv:2501.15915.
- Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; and Mittal, A. 2011. FEVER: a Large-scale Dataset for Fact Extraction and VERification. In *NAACL-HLT*.
- Wang, X.; Mikaelyan, L.; Isazawa, T.; and Hensman, J. 2024. KBLaM: Knowledge Base augmented Language Model. ArXiv.
- Wang, Y.; Agarwal, S.; Mukherjee, S.; Liu, X.; Gao, J.; Awadallah, A. H.; and Gao, J. 2022. AdaMix: Mixture-of-Adaptations for Parameter-efficient Model Tuning. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 5744–5760. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Wei, J.; Bosma, M.; Zhao, V.; Guu, K.; Yu, A. W.; Lester, B.; Du, N.; Dai, A. M.; and Le, Q. V. 2022. Finetuned Language Models Are Zero-Shot Learners. In *Proceedings of*

the 10th International Conference on Learning Representations (ICLR 2022).

Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davison, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Scao, T. L.; Gugger, S.; Drame, M.; Lhoest, Q.; and Rush, A. M. 2020. HuggingFace’s Transformers: State-of-the-art Natural Language Processing. arXiv:1910.03771.

Yang, Z.; Band, N.; Li, S.; Candès, E.; and Hashimoto, T. 2025. Synthetic continued pretraining. In *International Conference on Learning Representations*.

Zadouri, T.; Üstün, A.; Ahmadian, A.; Ermis, B.; Locatelli, A.; and Hooker, S. 2024. Pushing Mixture of Experts to the Limit: Extremely Parameter Efficient MoE for Instruction Tuning. In *The Twelfth International Conference on Learning Representations*.

Zhang, Z.; Zeng, Z.; Lin, Y.; Wang, H.; Ye, D.; Xiao, C.; Han, X.; Liu, Z.; Li, P.; Sun, M.; and Zhou, J. 2023. Plug-and-Play Knowledge Injection for Pre-trained Language Models. In Rogers, A.; Boyd-Graber, J.; and Okazaki, N., eds., *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 10641–10658. Toronto, Canada: Association for Computational Linguistics.

Zhao, Z.; Gan, L.; Wang, G.; Zhou, W.; Yang, H.; Kuang, K.; and Wu, F. 2024. LoraRetriever: Input-Aware LoRA Retrieval and Composition for Mixed Tasks in the Wild. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Findings of the Association for Computational Linguistics: ACL 2024*, 4447–4462. Bangkok, Thailand: Association for Computational Linguistics.

Adapter Details

We fit LoRA adapters for our knowledge and task libraries using the *peft* package (Mangrulkar et al. 2022). We use common hyperparameter settings, purposely not optimizing for LAG, which is designed to work with externally trained adapters. We use a learning rate of $1e^{-4}$ and LoRA dropout of 0.05 for both libraries. Following Fleshman and Van Durme (2025), we use rank-8 LoRAs for our task adapters, targeting the `k_proj`, `q_proj`, `v_proj`, and `o_proj` attention layers of the instruction-tuned model. We train for a single epoch with a batch size of 8. We use rank-6 knowledge adapters targeting the `gate_proj`, `up_proj`, and `down_proj` layers, following work suggesting that transformers store knowledge in their feed-forward layers (Geva et al. 2021). These are fit using the pretrained version of the model to mitigate impacting the instruction-following capabilities (Fleshman and Van Durme 2024). We train these LoRAs for 4 epochs with a batch size of 1. For both sets of adapters, we use a LoRA α which is twice the rank. All training and inference was done using a single Nvidia A100 GPU.