

# BOSS LLM: ADAPTATION VIA NO-REGRET LEARNING

**Yu Feng\*, Avishree Khare\*, Nghia Nguyen\*, & Sikata Sengupta\***

Department of Computer Science

University of Pennsylvania

Philadelphia, PA 19104, USA

{akhare, fengyu1, nghianhh, sikata}@seas.upenn.edu

## ABSTRACT

The diversity of Large Language Models (LLMs) calls for more effective strategies to combine their strengths across various tasks. In this work, we learn an adaptive mixture of multiple expert models as Boss LLM. By extending the multi-objective optimization with exponential weights (MOEW) algorithm (18), Boss LLM selects the most suitable model for a given prompt that could potentially span multiple categories with provable low regret for every category and expert model. Empirical results demonstrate that Boss LLM not only effectively adapts its mixture based upon the categories of a given prompt and improves upon the expert models, but also exhibits generalization properties.<sup>1</sup>

## 1 INTRODUCTION

Rapid progress and development in large language models (LLMs) have led to an explosion of available models, each developed by different providers and has unique strengths and weaknesses across different various domains and tasks (16; 9; 17; 11). This growing variety reflects the diverse nature of real-world queries, which often span multiple categories and require different areas of expertise, e.g., factuality, creativity, coding, or mathematical reasoning. Even a single prompt or query could feasibly be associated with many categories. Given such diversity, relying exclusively on a single LLM for all tasks often leads to suboptimal performance, as no individual model consistently outperforms others across all scenarios. While leveraging multiple specialized models appears intuitive here, the challenge remains in effectively orchestrating their use, specifically in deciding which model should be deployed given a particular prompt.

A popular solution to this challenge is the Mixture-of-Experts (MoE) architecture (14; 3; 2), which partitions a model into specialized “experts” and selectively activates them based on the input. This approach enables efficient inference through conditional computation. However, MoE models typically require substantial training data to learn effective gating mechanisms, as the router must be trained extensively to activate appropriate experts based on input patterns. Moreover, MoE architectures tend to be less interpretable, making it difficult to understand or debug the routing decisions. Alternatively, post-hoc multi-agent methods coordinate independently trained LLMs after their training is complete (21; 5; 20). These systems leverage existing models but rely on heuristic or empirical rules for inter-LLM interactions to select and combine outputs. Both MoE and multi-agent approaches lack formal performance guarantees, and their heuristic nature can lead to unpredictable results, especially in critical applications.

In this work, as illustrated in Figure 1, we propose a fundamentally different approach that systematically and optimally integrates diverse LLMs by introducing a boss LLM capable of selecting which model to respond with or use in a context-dependent manner. We focus on learning an adaptive mixture of expert LLMs, our Boss LLM, such that it achieves provably low regret across every category and expert model. Specifically, we extend the multi-objective optimization with exponential weights (MOEW) framework (18), originally developed for action-independent subsequences, to the LLM setting. We demonstrate that the Boss LLM effectively adjusts its model mixture based

\*Equal contribution. Author names are listed in alphabetical order.

<sup>1</sup>We provide a flexible package, boss-llm, that allows researchers to apply this framework to their own choice of models and tasks. Interested parties may contact us for access.

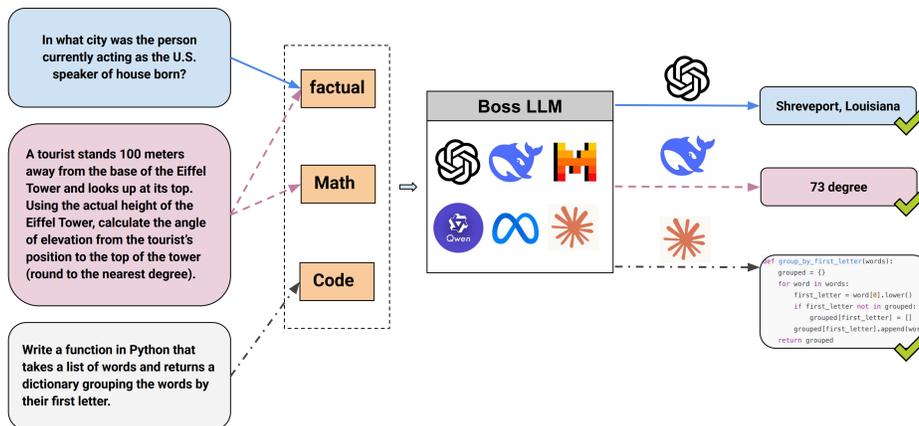


Figure 1: Overview of Boss LLM. We propose a novel approach that optimally integrates diverse LLMs by introducing a boss LLM, which dynamically selects the most suitable model to respond based on context.

on the prompt’s categories, consistently outperforming individual expert models. Additionally, we empirically evaluate its generalization capabilities and provide insights for further theoretical and empirical exploration.

### 1.1 RELATED WORK

In an online adversarial setting, one might seek regret guarantees that are not just with respect to each fixed expert model, but rather conditional on group (or in this context category) membership. In the setting where events are action-independent (as is the case for categories of prompts), (18) provides a nice closed-form algorithm (Multi-Objective Optimization with Exponential Weights). There are also several other lines of work around conditional regret guarantees and multi-objective optimization (10). In the space of Uncertainty Quantification and LLMs for factual correctness (with theoretical multi-group and multi-valid guarantees applied to LLM applications), (13) via calibration and conformal prediction.

Building on the theme of interactive updating processes, recent works on self-improvement for LLMs (25; 15; 19) explore iterative frameworks where models generate multiple outputs, verify their own results, and distill knowledge by selecting and refining the reranked or filtered results. One of the most closely related papers is (27), which formulates the problem of self-improvement as self-play with an LLM, allowing the policy to iteratively play against itself using no-regret learning. Our method adheres to the same philosophy but incorporates more controlled and precise updates for distilling knowledge from multiple LLMs.

The idea of LLM-merging has also recently emerged as it empirically has shown some signs of success regarding its ability to generalize to new tasks well from specialized fine-tuned existing models (Zhang et al.; 23). In a similar vein, mixture-of-experts models that use sparsity to selectively activate individual expert models have rapidly grown in interest (14). Our work offers a potential next step toward a more efficient and interpretable mixture-of-experts.

Meanwhile, there has also been quite a bit of work around the idea of Multi-LLM collaborations. Recent advancements in multi-LLM agent frameworks have resulted in the development of various platforms designed to enable collaboration among multiple agents (20). AutoGen (21) introduces a framework for building LLM-based applications where agents communicate with each other to accomplish tasks. CAMEL (12) is designed to enable autonomous collaboration among chat-based language models using role-playing and inception prompting. These frameworks complement recent efforts (22; 5; 7; 8) to improve factuality in LMs through multi-agent cooperation or debate, primarily using cross-model agents. Unlike the majority of related studies that rely on interactions between agents, our proposed method eliminates the need for inter-LLM interactions. By leveraging

the closed-form algorithm, a meta-agent is capable of effectively guiding the behavior of subagents without requiring direct communication among them.

## 2 MODEL SETUP

In this section, we present how a typical use case of LLMs fits into the Multi-Objective Optimization framework and how one can use the Exponential Weights algorithm to get a non-trivial guarantee on the performance of a multi-LLM setup compared with using the best LLM in hindsight. We consider the use case of answering a stream of prompts by the user, given access to multiple LLMs.

Formally, we have  $K$  expert LLMs, which we denote as  $\{LM_i\}_{i \in [K]}$ . We have to answer a stream of  $T$  prompts  $x_t \in \mathcal{X}$ ,  $t = 1 \dots T$  and each prompt belongs to one or more categories in the set of  $N$  categories  $\mathcal{C}$ . For example, the question ‘‘What is the Singular Value Decomposition of a matrix?’’ can belong to both the ‘‘factual’’ and ‘‘math’’ categories. We use the shorthand  $c(x_t)$  to denote the set of categories assigned to the prompt  $x_t$ , where  $c : \mathcal{X} \rightarrow 2^N$ . We define a collection of events  $\mathcal{E}$  as the collection of functions  $E_i(t, x_t) = I(c_i \in c(x_t))$  for all  $i \in [N]$ . In words,  $E_i(t, x_t)$  is the event that the query  $x_t$  belongs to the category  $c_i$ . We denote the space of responses  $\mathcal{Y}$ . We may think of a designated  $y_t$  or set of responses for each  $x_t$  corresponding to a desired or correct response.

Next, we formally define the sequential interaction setting between our multi-LLM method and the users of a chat system as a sequential (repeated) game between a learner and an adversary in Def 2.1.

**Definition 2.1** (Sequential Interaction Setting). For  $t = 1, \dots, T$

- Adversary selects categories  $C \in 2^{\mathcal{C}}$  and corresponding prompt  $x_t \in \mathcal{X}$ , which has the answer  $y_t \in \mathcal{Y}$
- Learner selects  $p_t \in \Delta[K]$  as its distribution over  $K$  expert LLMs.
- Adversary selects cost  $l_t \in [0, 1]^K$  with the knowledge of  $p_t$ .
- Learner experiences the expected loss  $\langle p_t, l_t \rangle$  over the random choice of the experts.

Here, the adversary represents how the set of prompts is selected and how the cost is calculated after the learner answers the prompts. The reason we chose this setting is that the categories of the prompts asked by the user and how satisfied the user is with the response of the system may not follow any distributional assumptions (e.g., identically and independently distributed prompts). Therefore, this setting allows us to have guarantees over potentially arbitrary category choices and values of the loss.

Next, we define the concept of regret in a repeated game, where  $i \in [K]$  is a fixed action.

**Definition 2.2** (Regret over a sequence compared to taking a fix action  $i$  (18)).

$$Reg(\pi^T, i) = \sum_{t=1}^T E(t, x_t) (\langle p_t, l_t \rangle - l_t^i)$$

Intuitively, regret quantifies how much worse playing a repeated game according to the set of distributions  $\{p_t\}_{t=1}^T$  compared with choosing a fixed action  $i$ . For the LLM setting, this helps us quantify how using our method compares with just choosing one LLM.

We can extend regret to also condition on questions belonging to a category (or even set of categories) with the following definition.

**Definition 2.3** (Action-Independent Subsequence Regret (18)).

$$Reg(\pi^T, E, i) = \sum_{t=1}^T E(t, x_t) (\langle p_t, l_t \rangle - l_t^i)$$

For our method, we calculate the probability of choosing the experts using the Multi-Objective Optimization with Exponential Weights (MOEW) algorithm. The complete algorithm is presented in Algorithm 1.

---

**Algorithm 1** Multi-Objective Optimization with Exponential Weights (MOEW) with LLMs. Adapted from (18, Algorithm 13, Pg. 50)

---

**Require:** Set of prompts  $\{x_t\}_{t=1}^T$  and categories  $\{c_t\}_{t=1}^T$ , Step size  $\eta$

**for**  $t = 1, \dots, T$  **do**

Adversary selects  $x_t, c_t$

Learner updates the distribution vector  $p^t \in [K]$  as

$$p_i^t = \frac{\sum_{k=1}^N E_k(t, x^t) \exp\left(\frac{\eta}{2} \sum_{t'=1}^{t-1} E_k(t', x^{t'}) (\langle p^{t'}, l^{t'} \rangle - l_i^{t'})\right)}{\sum_{j=1}^K \sum_{k=1}^N E_k(t, x^t) \exp\left(\frac{\eta}{2} \sum_{t'=1}^{t-1} E_k(t', x^{t'}) (\langle p^{t'}, l^{t'} \rangle - l_j^{t'})\right)}$$

Learner chooses among the  $K$  LLMs according to  $p^t$

Adversary selects  $l_t$  based on  $p^t$

**end for**

---

Using this algorithm, we get the following guarantee on the performance of our method compared to a fixed best LLM.

**Theorem 2.4.** *Algorithm 1 obtains action-independent subsequence regret of  $4\sqrt{T(\ln(N) + \ln(K))}$  where  $N$  is the number of categories,  $K$  is the number of LLMs, and  $T$  is the time horizon. This theorem (and proof) is from (18, Thm. 18, Pg. 50).*

$$\max_{k \in [N], i \in [K]} \text{Reg}(\pi^T, E_k, i) \leq 4\sqrt{T(\ln(N) + \ln(K))}.$$

*Proof.* We now provide a proof sketch. See (18, Thm. 18, Pg. 50) for full details. Note that we use slightly different notation for our adversary’s cost vectors and loss increment functions. In order to control the maximum regret over all categories and all experts, we formulate a multi-objective optimization game between an adversary and learner. We consider action-independent subsequences that only depend on the rounds in which each category is active for the prompt selected at this round. To do this, consider the following setup (related to the sequential interaction defined in Definition 2.1): there are  $N$  subsequence selection functions corresponding to the categories and  $K$  actions corresponding to each expert we could select. Then we can define an  $NK$ -dimensional multi-objective problem where at each round:

1. The action space for the learner is  $\mathcal{A}^t = \Delta[K]$  and the action space of the adversary is  $\mathcal{B}^t = [0, 1]^K$
2. The loss increment functions for every subsequence indicator function  $E \in \mathcal{E}$  and action  $i \in [K]$  is:  $d_{E,i}^t(p^t, l^t) = \mathbb{E}_{j \sim p^t} [E(t, j, x^t)(l_j^t - l_i^t)]$ . Notice in this case since we are considering action-independent subsequences we can simply write  $E(t, x^t)$  instead of  $E(t, j, x^t)$ .
3. We design a zero-sum game of the form  $u^t(a, b) = \sum_{i=1}^{d=NK} p_i^t d_{E,i}^t(a, b)$  and solve for the minimax equilibrium strategy of this game

In the case of action-independent subsequences, we can solve for a closed-form solution for the minimax strategies of these zero-sum games at each round. By reasoning about the Adversary Moves First value of the game (since it is zero-sum), we know that the value of the game for  $u^t$  is 0. Plugging in the closed form strategy of  $p^t$ , we are able to show that for any possible loss vector  $l$  chosen by the expert,  $u^t(p^t, l) \leq 0$  and therefore this strategy is able to achieve the value of the game. The regret we incur stems from the application of Exponential Weights to our  $d = NK$  dimensional objective optimization game (see (18, Thm. 15, Pg. 46)).  $\square$

This theorem provides us with the guarantee that by choosing the distribution over the  $K$  LLMs according to Algorithm 1 gives us an upper bound on the error we have after  $T$  rounds compared with playing the best action in hindsight after all  $T$  rounds. Because we do not know what the best action in hindsight is as we are playing the game (we can only calculate that at the end of  $T$  rounds), choosing the distribution according to the MOEW algorithm is the best we can realistically guarantee.

### 3 EXPERIMENTAL SETUP

In this section, we present how we set up the experiments to test how our algorithm performs with real-world open-source LLMs. First, we show how to use an LLM to simulate the adversary in Section 3.1. Next, we describe the datasets we use to simulate the prompts and categories that can come from a stream of queries from users in Section 3.2. We show how the LLMs we use in our experiments (both the  $K$  expert LLMs and the adversary LLM) behave for a prompt from the user in Section 3.3. Finally, we describe what we consider a generalization scenario for the learned mixture to unseen prompts in Section 3.4.

#### 3.1 ADVERSARY SIMULATION

The adversary, given access to transcript up until the current time step will select a prompt (alongside its corresponding category and true answer). To select a loss vector, our adversary is an LLM trained via zero-shot learning to select prompts, categories to follow this sequential interaction protocol. For a given prompt  $x_t$  and select correct answer  $y_t$ , the adversary will score the responses provided by the other expert models based upon some examples and overlap with the select correct answer. The adversary will also be provided with access to  $p_t$  before selecting its loss vector, though in practice it will likely not affect the adversary’s loss assignment in this setting. We consider utilizing a variety of different models for the adversary LLM, and, for our experiments, use Llama-3.1-8B-Instruct (6). For category selection, our main adversary is selecting the worst category performance of the learner in hindsight (weighted by proximity to current time step). We also include  $\epsilon$ -greedy to encourage some exploration of category selection.

For loss scoring, an example of a prompt to the Adversary LLM is given in Appendix Listing 1.

#### 3.2 CATEGORIES AND DATASETS

The main categories we consider in this project are factuality, math, and code. While we don’t explicitly include prompts / questions that correspond to multiple categories at once, our framework can be extended to such cases (for example, questions that require retrieving a fact and using that to answer a math problem). In the future, we also hope to expand our evaluation to incorporate additional categories from Llama3’s human evaluation (noting that three are already part of our current setup).

We next list the datasets we use for the factuality, math, and code categories, respectively.

**HotpotQA.** (24) provide “HotpotQA, which is a new dataset with 113k Wikipedia-based question-answer pairs with four key features: (1) the questions require finding and reasoning over multiple supporting documents to answer; (2) the questions are diverse and not constrained to any pre-existing knowledge bases or knowledge schemas; (3) it provides sentence-level supporting facts required for reasoning, allowing QA systems to reason with strong supervision and explain the predictions; (4) it offers a new type of factoid comparison questions to test QA systems’ ability to extract relevant facts and perform necessary comparison.” An example of a question from HotpotQA, with the answer and expert responses, is given in Appendix Listing 2.

**GSM8K.** (4) provide “GSM8K (Grade School Math 8K) which is a dataset of 8.5K high quality linguistically diverse grade school math word problems. The dataset was created to support the task of question answering on basic mathematical problems that require multi-step reasoning. These problems take between 2 and 8 steps to solve. Solutions primarily involve performing a sequence of elementary calculations using basic arithmetic operations to reach the final answer. An example of a question from GSM8K, with the answer and expert responses, is given in Appendix Listing 3.

**Code dataset.** The final dataset “contains problem descriptions and code in python language. This dataset is taken from sahil2801/code.instructions.120k, which adds a prompt column in alpaca style.” (1). An example of a question from this dataset, with the answer and expert responses, is given in Appendix Listing 4.

### 3.3 EXPERTS

We utilize the following expert models via Hugging Face for our experiments: 1) mistralai/Mistral-7B-Instruct-v0.3<sup>2</sup>; 2) Qwen/Qwen2.5-7B-Instruct<sup>3</sup>; 3) Qwen/Qwen2.5-0.5B-Instruct<sup>4</sup>; 4) meta-llama/Llama-3.2-1B-Instruct<sup>5</sup>. We present an example question from HotpotQA along with the corresponding response from each of the five possible experts in Appendix A.1.

### 3.4 GENERALIZATION

In this section, we explore a simple notion of generalization. We study how the topic-based weights over experts learned over train-time generalize to new unseen prompts when we are no longer updating weights based on losses. More precisely, we explore how the final mixture of weights selected over topics (or combinations of topics) presented by our algorithm fares against the base experts on new examples when our algorithm is no longer receiving the loss and updating the mixture weights. We do this by implementing our algorithm with the existing history, but only providing 0-loss vector feedback at test-time. We still document the “true” loss vectors selected by our adversary for the purposes of evaluation. We do this to simply check whether these mixtures have any properties that make them more successful on new unseen examples across topics empirically. Theoretical guarantees for the mixture weights are an interesting open question.

## 4 RESULTS

In this section, we design the experiments to test the empirical performance of our BossLLM framework on real-world open-source LLMs. First, we analyze the regret of the learner relative to all the base experts across all categories as shown in Figure 2 and Figure 3. We also provide the subsequence regret for every topic and expert in table 1. We also then consider the question of generalization (Figure 3). That is, we consider how the mixture model of the learner fares on unseen prompts compared to the expert model and across topics, when it is no longer receiving immediate information about the loss of all responses.

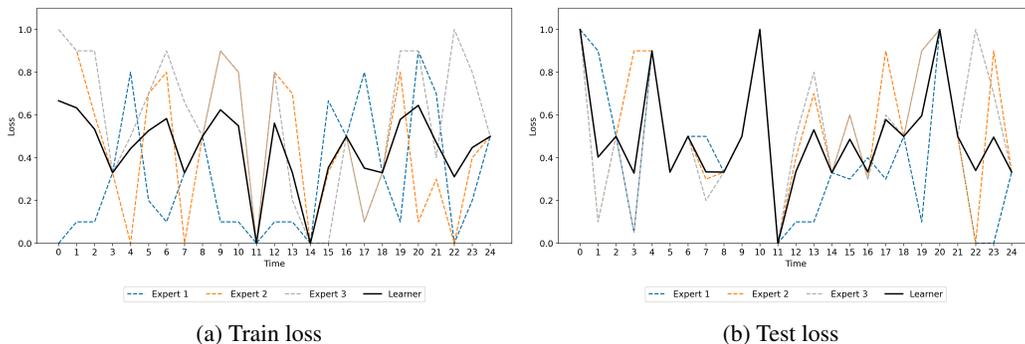


Figure 2: Loss over  $T=25$  rounds incurred by 3 base experts and the learner’s selection. The 3 base experts are Mistral-7B-Instruct-v0.3, Qwen2.5-0.5B-Instruct and Llama-3.2-1B-Instruct respectively. The adversary is simulated by Meta-Llama-3-8B-Instruct.

<sup>2</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

<sup>3</sup><https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

<sup>4</sup><https://huggingface.co/Qwen/Qwen2.5-0.5B-Instruct>

<sup>5</sup><https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct>

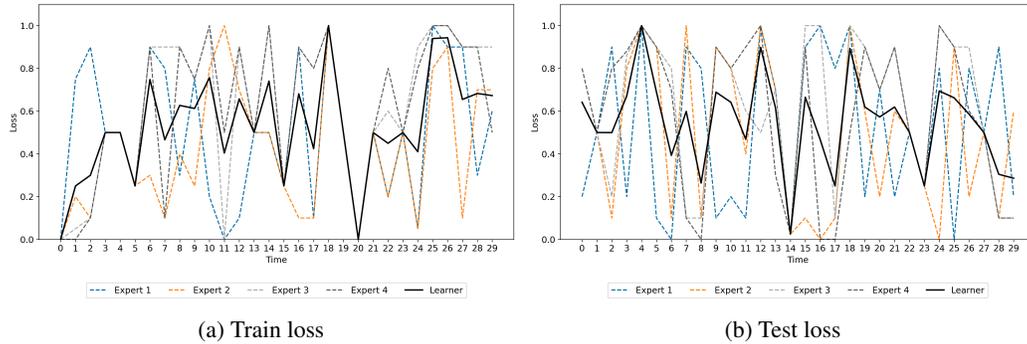


Figure 3: Loss over  $T = 30$  rounds incurred by 4 base experts and the learner’s selection. The 4 base experts are Mistral-7B-Instruct-v0.3, Qwen/Qwen2.5-7B-Instruct, Qwen2.5-0.5B-Instruct and Llama-3.2-1B-Instruct respectively. The adversary is simulated by Meta-Llama-3-8B-Instruct.

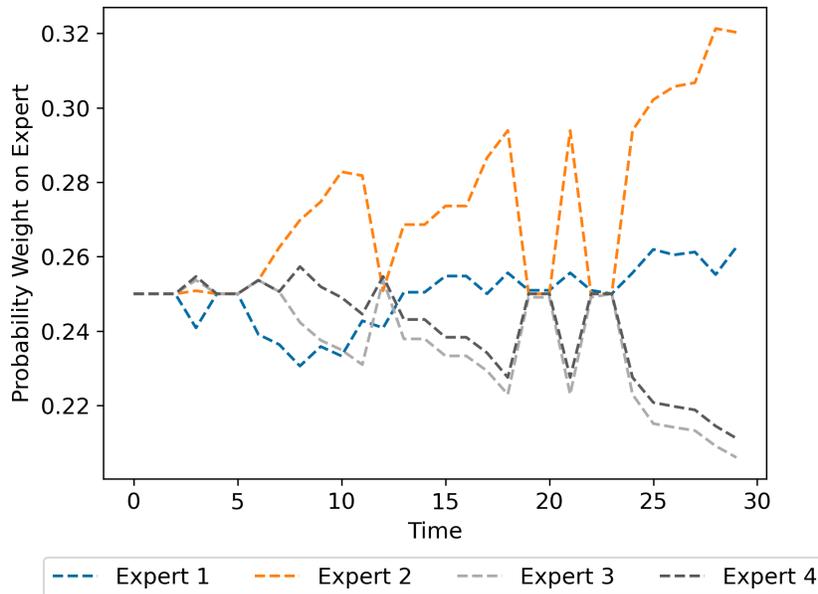


Figure 4: Learner’s Distribution over Experts from Figure 3, across  $T = 30$  rounds. The 4 base experts are Mistral-7B-Instruct-v0.3, Qwen2.5-7B-Instruct, Qwen2.5-0.5B-Instruct and Llama-3.2-1B-Instruct respectively. The adversary is simulated by Meta-Llama-3-8B-Instruct.

Category/Model	Mistral-7B-Instruct-v0.3	Qwen2.5-0.5B-Instruct	Llama-3.2-1B-Instruct
Factual	0.323	-0.077	-0.177
Math	1.502	0.065	-1.662
Code	1.718	-0.282	-1.582

Table 1: Subsequence Regret by Expert (from Figure 2) and Category for  $T = 25$  and  $l_t \in [0, 1]$ . A negative regret value means that our learner performs better than the model in the corresponding column at the category in the corresponding row and vice versa.

Category/Model	Mistral-7B-Instruct-v0.3	Qwen2.5-7B-Instruct	Qwen2.5-0.5B-Instruct	Llama-3.2-1B-Instruct
Factual	0.963	3.163	-2.787	-1.987
Math	0.000	0.000	0.000	0.000
Code	0.307	0.257	-0.193	-0.343

Table 2: Subsequence Regret by Expert (from Figure 3) and Category for  $T = 30$  and  $l_t \in [0, 1]$ .

In Figure 2 we find that when we include one larger model and two smaller models as experts, over time, more weight is put on the expert that is larger (first) model  $[0.37, 0.33, 0.29]$  which we may intuitively expect to be the expert performing better. We also broadly see in this figure that the learner is performing competitively with the best fixed expert over the training set. In the test set (regarding generalization) we also see that the learner still performs somewhat competitively. In Figure 3, however, when we include several larger models, we find that on categories where no expert performs particularly better than the others (loss wise over time) that our Boss LLM broadly remains with a uniform distribution over the base experts. In Table 1, we can see for different categories, there are different experts with respect to which we have maximum regret. All of these are quite small (given our time horizon and regret bound). However, in Figure 3 when there are two larger expert models and two smaller expert models, one can see that more weight goes towards the first two models (as seen in Figure 4,  $[0.26, 0.32, 0.21, 0.21]$ ) and the learner performs relatively competitively with the best fixed expert in hindsight.

We include plots of the evolution of probability mass put on each expert across time as well, as can be seen in Figure 4, to study how the learner is updating sequentially across categories. Recall that in this example, since the first two experts are larger in terms of number of parameters than the subsequent expert models, one can see that over time the learner is putting more weight on the first two models than the latter two. Moreover, we can see from Table 2 that for the categories sampled (here mainly Factual and Code), the subsequence regret values are for the most part quite low. In the future, we intend to evaluate several other ablations on our experiments. First, we want to consider the role of enabling the adversary to evaluate the loss of all responses combined vs. individual evaluation of each expert’s response separately. We also want to evaluate how the adversary generates losses when looking at the transcript vs just the current round expert responses.

Of course, our adversary is not necessarily always playing its best response against the learner’s selected distribution, and our losses are not necessarily always reflective of the underlying properties of the quality of the responses generated by the expert. In the next section, we propose some directions for future exploration to address this point.

## 5 CONCLUSION AND FUTURE WORK

In this work, we apply the MOEW algorithm of (18) to multiple expert LLM models to enable adaptation based upon various possibly intersecting categories of prompts. We empirically verify the theoretical results by showcasing how our Boss LLM models is competitive with the best of the expert models across topics (especially in the limit as the horizon becomes larger). A strength of our approach that it enables the Boss LLM to select weights according to the various categories that the prompt is a member of in a clean fashion without needing to handle disjoint categories and weights. Moreover, it has provable guarantees on how much worse than the base class of models it can do across all possible topics, which enables one to select parameters according to their specific objectives and goals.

A limitation of the current approach is our reliance on a large language model to simulate the adversary. That is, first, our adversary may not select the worst-case action against the learner in certain rounds. Moreover, we rely on the general success of few-shot learning to teach the adversary how to score the responses generated by the expert model LLMs. Lastly, we are also in the full information setting, which can become costly to generate for large  $K$ .

We also assume in our current experiments that each prompt belongs to a single category. Future work could explore more complex prompts that may fall under multiple categories. Another limitation of our current pipeline is that we do not currently generate a final answer for each prompt question. We leave it for future work to aggregate expert responses to a single, high-quality output for evaluating final task performance.

We propose a few possible extensions to address some of the limitations above. First, we could consider trying to calibrate the loss function of the adversary if we are provided with some labeled human feedback of possible responses to given prompts. We could also consider the bandit information setting rather than the full-feedback setting as this may be more realistic and practical. Lastly, we could consider extending this into other settings by considering no-swap regret versions of such algorithms. We could also consider generalizations of the following approach via Reinforcement Learning with Human Feedback (RLHF).

## REFERENCES

- [1] (2024). Dataset card for python\_code\_instructions.18k\_alpaca. *Hugging Face*.
- [2] Artetxe, M., Bhosale, S., Goyal, N., Mihaylov, T., Ott, M., Shleifer, S., Lin, X. V., Du, J., Iyer, S., Pasunuru, R., Anantharaman, G., Li, X., Chen, S., Akin, H., Baines, M., Martin, L., Zhou, X., Koura, P. S., O’Horo, B., Wang, J., Zettlemoyer, L., Diab, M., Kozareva, Z., and Stoyanov, V. (2022). Efficient large scale language modeling with mixtures of experts.
- [3] Cai, W., Jiang, J., Wang, F., Tang, J., Kim, S., and Huang, J. (2025). A survey on mixture of experts in large language models. *IEEE Transactions on Knowledge and Data Engineering*, page 1–20.
- [4] Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. (2021). Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- [5] Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. (2024). Improving factuality and reasoning in language models through multiagent debate.
- [6] Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. (2024). The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- [7] Feng, S., Shi, W., Wang, Y., Ding, W., Balachandran, V., and Tsvetkov, Y. (2024a). Don’t hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration. *arXiv preprint arXiv:2402.00367*.
- [8] Feng, Y., Htut, P. M., Qi, Z., Xiao, W., Mager, M., Pappas, N., Halder, K., Li, Y., Benajiba, Y., and Roth, D. (2024b). Diverseagententropy: Quantifying black-box llm uncertainty through diverse perspectives and multi-agent interaction.
- [9] Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes, D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan, E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J., Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J., Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton, J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia, J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li, K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhota, K., Rantala-Yearly, L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat, L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N., Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan, P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S., Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang, S., Raparthi, S., Shen, S., Wan, S., Bhosale, S., Zhang, S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan, S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speckbacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang, X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Goldschlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z., Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey, A., Shajnfeld, A., Gangidi, A., Victoria, A., Goldstand, A., Menon, A., Sharma, A., Boesenberg,

- A., Baevski, A., Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A., Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poulton, A., Ryan, A., Ramchandani, A., Dong, A., Franco, A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel, A., Bharambe, A., Eisenman, A., Yazdan, A., James, B., Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola, B. D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock, B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B., Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C., Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C., Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty, D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine, D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang, D., Le, D., Holland, D., Dowling, E., Jamil, E., Montgomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T., Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun, F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Caggioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz, G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov, G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H., Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H., Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan, H., Damaj, I., Molybog, I., Tufanov, I., Leontiadis, I., Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli, J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J., Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J., Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J., McPhie, J., Torres, J., Ginsburg, J., Wang, J., Wu, K., U, K. H., Saxena, K., Khandelwal, K., Zand, K., Matosich, K., Veeraraghavan, K., Michelena, K., Li, K., Jagadeesh, K., Huang, K., Chawla, K., Huang, K., Chen, L., Garg, L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L., Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M., Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso, M., Groshev, M., Naumov, M., Lathi, M., Keneally, M., Liu, M., Seltzer, M. L., Valko, M., Restrepo, M., Patel, M., Vyatskov, M., Samvelyan, M., Clark, M., Macey, M., Wang, M., Hermoso, M. J., Metanat, M., Rastegari, M., Bansal, M., Santhanam, N., Parks, N., White, N., Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta, N., Laptev, N. P., Dong, N., Cheng, N., Chernoguz, O., Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P., Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P., Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P., Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R., Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan, R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta, S., Siby, S., Bondu, S. J., Datta, S., Chugh, S., Hunt, S., Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma, S., Yamamoto, S., Ramaswamy, S., Lindsay, S., Lindsay, S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S., Zhang, S., Zhang, S., Wang, S., Agarwal, S., Sajuyigbe, S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield, S., Govindaprasad, S., Gupta, S., Deng, S., Cho, S., Virk, S., Subramanian, S., Choudhury, S., Goldman, S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson, T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked, T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V., Kumar, V. S., Mangla, V., Ionescu, V., Poenaru, V., Mihalescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W., Bouaziz, W., Constable, W., Tang, X., Wu, X., Wang, X., Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y., Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu, Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Y., He, Y., Rait, Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao, Z., and Ma, Z. (2024). The llama 3 herd of models.
- [10] Haghtalab, N., Jordan, M., and Zhao, E. (2024). A unifying perspective on multi-calibration: Game dynamics for multi-objective learning. *Advances in Neural Information Processing Systems*, 36.
- [11] Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. (2023). Mistral 7b.
- [12] Li, G., Hammoud, H. A. A. K., Itani, H., Khizbullin, D., and Ghanem, B. (2023). CAMEL: Communicative agents for "mind" exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- [13] Liu, T. and Wu, Z. S. (2024). Multi-group uncertainty quantification for long-form text generation. *arXiv preprint arXiv:2407.21057*.
- [14] Lo, K. M., Huang, Z., Qiu, Z., Wang, Z., and Fu, J. (2024). A closer look into mixture-of-experts in large language models. *arXiv preprint arXiv:2406.18219*.
- [15] Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhume, S., Yang, Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck, S., Yazdanbakhsh,

- A., and Clark, P. (2023). Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- [16] OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. (2024). Gpt-4 technical report.
- [17] Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and Qiu, Z. (2025). Qwen2.5 technical report.
- [18] Roth, A., editor (2024). *Learning in Games*. Philadelphia, PA.
- [19] Song, Y., Zhang, H., Eisenach, C., Kakade, S., Foster, D., and Ghai, U. (2024). Mind the gap: Examining the self-improvement capabilities of large language models. *arXiv preprint arXiv:2412.02674*.
- [20] Wang, J., Wang, J., Athiwaratkun, B., Zhang, C., and Zou, J. (2024). Mixture-of-agents enhances large language model capabilities.
- [21] Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang, L., Zhang, X., Zhang, S., Liu, J., Awadallah, A. H., White, R. W., Burger, D., and Wang, C. (2023). Autogen: Enabling next-gen llm applications via multi-agent conversation.

- [22] Xiong, K., Ding, X., Cao, Y., Liu, T., and Qin, B. (2023). Examining inter-consistency of large language models collaboration: An in-depth analysis via debate. In Bouamor, H., Pino, J., and Bali, K., editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7572–7590, Singapore. Association for Computational Linguistics.
- [23] Yang, E., Shen, L., Guo, G., Wang, X., Cao, X., Zhang, J., and Tao, D. (2024). Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities. *arXiv preprint arXiv:2408.07666*.
- [24] Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. (2018). HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- [25] Zelikman, E., Wu, Y., Mu, J., and Goodman, N. (2022). STar: Bootstrapping reasoning with reasoning. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K., editors, *Advances in Neural Information Processing Systems*.
- [Zhang et al.] Zhang, Y., Ding, Y., Wu, J., and Yang, Y. Llm merging competition technical report for neurips 2024: Efficiently building large language models through merging. In *LLM Merging Competition at NeurIPS 2024*.
- [27] Zhang, Y., Yu, D., Peng, B., Song, L., Tian, Y., Huo, M., Jiang, N., Mi, H., and Yu, D. (2024). Iterative nash policy optimization: Aligning llms with general preferences via no-regret learning. *arXiv preprint arXiv:2407.00617*.

## A APPENDIX

### A.1 A SAMPLE QUESTION WITH CORRESPONDING RESPONSES FROM EXPERTS

**Question from hotpotQA:**

prompt: Which actress does this Netflix’s project by Nick Bernardone and Tina Fey star in the title role?

category: factual

**Answer from hotpotQA:**

answer: Ellie Kemper

**Adversary response:**

To calculate the loss vector, I will evaluate the responses of the experts and the learner’s distribution over the experts. Since the category is factual, I will compare the responses to the true answer, Ellie Kemper. learner’s distribution over the experts is not provided, but I will assume it is uniform, i.e., the learner is equally uncertain about all experts. Expert 0 provides a correct response, mentioning Ellie Kemper as the star of the Netflix project *Girlboss*. Expert 1 provides an incorrect response, stating Emily Blunt as the star. Expert 2 provides an unclear response, stating that they couldn’t find any information about a Netflix project starring Tina Fey and Nick Jonas. on this information, I will calculate the loss vector as follows: \* For Expert 0, the loss is 0, as the response is correct. \* For Expert 1, the loss is 1, as the response is incorrect. \* For Expert 2, the loss is 1, as the response is unclear and does not provide the correct answer. vector for the learner:  $[0, 1, 1]$

Model	Response	Adversary’s Loss
Mistral-7B-Instruct-v0.3	The actress who stars in the title role of the Netflix project by Nick Bernardone and Tina Fey is Ellie Kemper. The project is called <i>Girlboss</i> , based on the autobiographical book of the same name by Sophia Amoruso. Ellie Kemper plays the character Sophia Marlowe, a fictionalized version of Amoruso.”	0
Qwen2.5-0.5B-Instruct	The actress who stars in the title role of the Netflix project by Nick Bernardone and Tina Fey is Emily Blunt.	1
Llama-3.2-1B-Instruct	I couldn’t find any information about a Netflix project starring Tina Fey and Nick Jonas. It’s possible that you may be thinking of a different project or incorrect information.	1

Table 3: Expert Responses to Prompt from HotpotQA and losses from the adversary

### A.2 ADVERSARY SIMULATION

**System:** You are a judge evaluating a learner on K tasks. The learner has access to N experts who provide responses to the tasks. Your task is to provide a loss vector over experts corresponding to a given question/prompt, the responses of the experts, and the category of the question/prompt. The loss vector should be in the range [0,1] but the values need not sum to 1. For example, if the true answer is closest to the response of the first expert and farthest from the responses of the other three experts, the loss vector could be [0.1, 0.9, 0.9, 0.9] if there are four experts and learner is equally sure about all four experts. On the other hand, if all the four experts are correct, the loss vector would be [0, 0, 0, 0]. The learner's response is provided as it's distribution over the experts. Your response should start with some reasoning about how the loss vector is calculated and should strictly end with the following: 'Loss vector for the learner: [<YOUR LOSS VECTOR OF LENGTH N>]'.

**User:** Here are your interactions with the learner:  
Start of round t=3:

*Category chosen:* factual

*Question from the category to the learner:* Command and Control was a book that discusses the 1980 incident in what state?

*True Answer to the question:* Arkansas

*Response from expert 0 to the question:* The book \"Command and Control\" by Eric Schlosser primarily focuses on the 1989 Titan II missile accident that occurred in Damascus, Arkansas.

*Response from expert 1 to the question:* New Mexico

*Learner's distribution over experts:* [0.49999062 0.50000938]

*Loss vector for the learner:* ?.

Let's calculate the loss vector of length 2 for round t=3:

**Response:**

Listing 1: An example of the adversary's prompt to score expert responses

**Question:** Command and Control was a book that discusses the 1980 incident in what state?

**Answer:** Arkansas

**Response (Expert 0):** The book \"Command and Control\" by Eric Schlosser primarily focuses on the 1989 Titan II missile accident that occurred in Damascus, Arkansas.

**Response (Expert 1):** New Mexico

—

Listing 2: An example from HotpotQA

**Question:** On Monday, Matt worked for 450 minutes in his office. On Tuesday, he worked half the number of minutes he worked on Monday. On Wednesday, he worked for 300 minutes. How many more minutes did he work on Wednesday than on Tuesday?

**Answer:** On Tuesday Matt worked 450 minutes / 2 =  $\ll 450/2=225 \gg$  225 minutes.  
On Wednesday Matt worked 300 minutes - 225 minutes =  $\ll 300-225=75 \gg$  75 minutes more.  
#### 75

**Response (Expert 0):** Matt worked for 450 minutes on Monday. On Tuesday, he worked half of that, which is  $450/2 = 225$  minutes. On Wednesday, he worked for 300 minutes. So, he worked  $300 - 225 = 75$  more minutes on Wednesday than on Tuesday.

**Response (Expert 1):** To find the difference in minutes worked on Wednesday and Tuesday, we need to subtract the minutes worked on Tuesday from the minutes worked on Wednesday.  
 $\backslash \backslash$  Wednesday: 300 minutes  
 $\backslash \backslash$  Tuesday: 450 minutes  
 $\backslash \backslash$   $300 - 450 = -150$  minutes  
 $\backslash \backslash$  So, Matt worked 150 minutes less on Wednesday than on Tuesday.

—

Listing 3: An example from GSM8k

**Question:** Implement a Python program to check if a number is divisible by another number. Input:7, 2

**Answer:**

```
def is_divisible(n, m):
    if (n % m == 0):
        return True
    else:
        return False
print(is_divisible(7, 2)) # Outputs False
```

**Response (Expert 0):**

```
```python
def is_divisible(num, divisor):
    """
    Check if the first number is divisible by the second number.
    Parameters:
        num (int): The number to check.
        divisor (int): The number to divide by.
    Returns:
        bool: True if num is divisible by divisor, False otherwise.
    """
    return num % divisor == 0

# Example usage:
num = 7
divisor = 2
print(is_divisible(num, divisor)) # Output: True
```
```

**Response (Expert 1):**

```
```python
def check_divisibility(num1, num2):
    if num2 == 0:
        print("Number 2 is not divisible by number 1")
    elif num1 % num2 == 0:
        print(f"{num1} is divisible by {num2}")
    else:
        print(f"{num1} is not divisible by {num2}")

# Test the function
check_divisibility(7, 2)
```
```

Listing 4: An example from the Code dataset