

# FUNCTIONAL-LEVEL UNCERTAINTY QUANTIFICATION FOR CALIBRATED FINE-TUNING ON LLMs

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

From common-sense reasoning to domain-specific tasks, parameter-efficient fine tuning (PEFT) methods for large language models (LLMs) have showcased significant performance improvements on downstream tasks. However, fine-tuned LLMs often struggle with overconfidence in uncertain predictions, particularly due to sparse training data. This overconfidence reflects poor epistemic uncertainty calibration, which arises from limitations in the model’s ability to generalize with limited data. Existing PEFT uncertainty quantification methods for LLMs focus on the post fine-tuning stage and thus have limited capability in calibrating epistemic uncertainty. To address these limitations, we propose Functional-Level Uncertainty Quantification for Calibrated Fine-Tuning (UQ4CT), which captures and calibrates functional-level epistemic uncertainty during the fine-tuning stage via a mixture-of-expert framework. We show that UQ4CT reduces Expected Calibration Error (ECE) by more than 25% while maintaining high accuracy across 5 benchmarks. Furthermore, UQ4CT maintains superior ECE performance with high accuracy under distribution shift, showcasing improved generalizability.

## 1 INTRODUCTION

Large Language Models (LLMs) have revolutionized various domains as general task solvers (Chang et al., 2024). To adapt LLMs for specific downstream tasks or create instruction-following models, fine-tuning have become increasingly important (Houlsby et al., 2019; Hu et al., 2021a; Liu et al., 2022; Ding et al., 2022; 2023). This involves additional training on pre-trained LLMs using a smaller dataset (Zhong et al., 2021; Ren et al., 2022). Through fine-tuning, the model parameters are updated to better adapt to the domain-specific knowledge (Peng et al., 2023). To reduce the computational cost for fine-tuning, Hu et al. (2021a) proposed Low-Rank Adaptation (LoRA), which effectively reduces the parameters required for fine-tuning by introducing low-rank trainable matrices at each layer of the transformer architecture instead of fine-tuning the full model parameters. Li et al. (2024); Wu et al. (2024b) proposed LoRA Mixture-of-Experts (MoE) models which grants better performance while maintaining parameter efficiency.

However, previous studies have shown that fine-tuned LLMs are often overconfident with their predictions (Xiao et al., 2022c; He et al., 2023; Tian et al., 2023; OpenAI, 2023). This resembles poorly calibrated uncertainty (Zhou et al., 2022) due to the sparsity of fine-tuning data. Overconfidence is a crucial problem in safety-related decision making or in fields where data is very limited, such as experimental design, climate science and medical diagnosis (Singhal et al., 2022; Wu et al., 2023a; Lampinen et al., 2023; Li et al., 2022). Thus, methods that enhance uncertainty quantification of fine-tuned LLMs is urgently needed to assure trustworthy predictions for better application.

Established uncertainty quantification methods have been studied in conjunction with the LoRA structure. Monte-Carlo dropout (Gal & Ghahramani, 2016b) interprets dropout in neural networks as approximate Bayesian inference in deep Gaussian processes, allowing uncertainty estimates to be obtained from existing LoRA adapters without modifying them. Checkpoint ensemble (Chen et al., 2017) utilizes predictions from multiple LoRA checkpoints saved during a single fine-tuning process to calibrate uncertainty. Deep ensemble (Lakshminarayanan et al., 2017; Wang et al., 2023; Zhai et al., 2023a) combines the predictions from multiple LoRA adapters for better uncertainty calibration. Laplace-LoRA (Yang et al., 2024a) applies Bayesian inference via Laplace approximation to the LoRA parameters after fine-tuning, resulting in improved calibration and uncertainty estimates.

054 Although these methods have demonstrated improved uncertainty estimations, they all utilize a fixed  
 055 set of LoRA parameters fine-tuned over the entire downstream task dataset. The point estimates of  
 056 parameters have very limited capabilities in capturing epistemic uncertainty, while direct calibration  
 057 of epistemic uncertainty over the entire LoRA parameter space is an ideal but not practical approach.

058 Therefore, we propose Functional-Level Uncertainty Quantification for Calibrated Fine-Tuning  
 059 (UQ4CT) to calibrate the functional-level epistemic uncertainty via the LoRA MoE architecture.  
 060 We propose a functional perspective on LoRA parameters where we treat the LoRA experts as basis  
 061 functions and consider the more complex, prompt dependent functions as mixtures of those basis  
 062 functions. On top of learning the parameters in the LoRA experts, UQ4CT also trains a prompt-  
 063 dependent LoRA mixture to form a calibrated distribution over the functional space. The LoRA  
 064 experts capture different functional relationships in the fine-tuning data throughout training, and the  
 065 MoE routers dynamically select these functional bases conditioned on the input. The selection pro-  
 066 cess models the functional level epistemic uncertainty, and consequently captures the uncertainty in  
 067 the output space. We calibrate functional level epistemic uncertainty to align with predictive cor-  
 068 rectness during training time. This significantly improves uncertainty estimations of the model on  
 069 its predictions without compromising the accuracy. To summarize, our contributions include:

- 070 • A novel epistemic uncertainty quantification approach with Mixture-of-Experts architec-  
 071 ture during fine-tuning stage to model functional level epistemic uncertainty and align with  
 072 predictive correctness, which mitigates overconfidence issue and improves generalizability.
- 073 • A novel training calibration loss function incorporating predictive correctness to calibrate  
 074 the prompt-dependent LoRA mixture for better uncertainty estimation.
- 075 • More than 25% Expected Calibration Error reduction on 4 common-sense reasoning tasks  
 076 and 1 domain-specific question answering task, superior ECE performance under distri-  
 077 bution shift scenarios on 2 common-sense reasoning tasks and 4 domain-specific question  
 078 answering tasks without compromising accuracy.

## 080 2 PRELIMINARIES

### 081 2.1 LOW-RANK ADAPTATION (LORA)

082 LLMs have numerous large weight matrices to perform matrix multiplication, denoted as  $\mathbf{W}_0 \in$   
 083  $\mathbb{R}^{n_{\text{out}} \times n_{\text{in}}}$  that maps inputs  $\mathbf{x}$  to outputs  $\mathbf{h}$ . Hu et al. (2021a) proposes LoRA, which fixes  $\mathbf{W}_0$  and  
 084 introduces a low-rank perturbation  $\Delta\mathbf{W}$  to the weight matrix:

$$085 \mathbf{h} = \mathbf{W}_0\mathbf{x} + \Delta\mathbf{W}\mathbf{x} = \mathbf{W}_0\mathbf{x} + \mathbf{B}\mathbf{A}\mathbf{x}. \quad (1)$$

086 Here,  $\Delta\mathbf{W}$  is calculated as the product of two matrices,  $\mathbf{B} \in \mathbb{R}^{n_{\text{out}} \times n_{\text{lr}}}$  and  $\mathbf{A} \in \mathbb{R}^{n_{\text{lr}} \times n_{\text{in}}}$  where  $n_{\text{lr}}$  is  
 087 significantly smaller than  $n_{\text{in}}$  or  $n_{\text{out}}$ . For example, we use  $n_{\text{lr}} = 32$  while  $n_{\text{in}} = n_{\text{out}} = 4096$  for the  
 088 Llama2-7b model (Touvron et al., 2023c). Therefore, the total number of LoRA parameters for this  
 089  $\Delta\mathbf{W}$  is  $n_{\text{lr}}(n_{\text{in}} + n_{\text{out}})$ , which is far smaller than the parameter count of the full matrix,  $n_{\text{in}}n_{\text{out}}$ . One  
 090 of the key motivations of incorporating LoRA to fine-tune LLMs is the vast amount of memory cost  
 091 reduction compared with fine-tuning on the full model. For an LLM with 7 billion parameters, main-  
 092 taining the average gradient and average squared gradients for optimization multiplies the memory  
 093 required by a factor of 3 compared to simply loading model weights. LoRA greatly mitigates this  
 094 memory cost as the tripled memory consumption only applies to LoRA adapters.

### 095 2.2 MIXTURE OF EXPERTS (MOE)

096 LoRA Mixture-of-Experts (Li et al., 2024; Wu et al., 2024b) is an efficient approach to scale the  
 097 number of parameters while maintaining the same computational bounds. LoRA MoE utilizes the  
 098 top-k router to assign each token to the LoRA experts (Lepikhin et al., 2020). The router is a linear  
 099 layer that maps the input hidden state  $\mathbf{h}$  to a probability distribution of candidate experts.

100 Let  $\mathbf{h}_i^\ell \in \mathbb{R}^{1 \times d}$  ( $1 \leq i \leq s, 1 \leq \ell \leq L$ ) denote the output hidden state of the  $i$ -th token at the  $\ell$ -th  
 101 layer of the LLM, where  $L$  is the number of LLM layers and  $d$  is the hidden dimension. With  $\mathbf{W}_r^\ell$  as

the trainable router weight at layer  $\ell$ , the top-k gate router chooses  $k$  experts with highest probability given a hidden state  $\mathbf{h}_i^\ell$ :

$$R^\ell(\mathbf{h}_i^\ell) = \text{KeepTop-k}(\text{Softmax}(\mathbf{W}_r^\ell \cdot \mathbf{h}_i^\ell)). \quad (2)$$

Finally, we obtain the final MixLoRA prediction with:

$$\text{MixLoRA}(\mathbf{h}^\ell) = \sum_{k=1}^K R^\ell(\mathbf{h}^\ell)_k E_k^\ell(\mathbf{h}^\ell), \quad E_k^\ell(\mathbf{h}^\ell) = \mathbf{W}^\ell \cdot \mathbf{h}^\ell + \mathbf{B}_i^\ell \mathbf{A}_i^\ell \cdot \mathbf{h}^\ell \quad (3)$$

where  $\mathbf{W}$  is the pretrained weights of the feed-forward network (FFN) layer and  $\mathbf{B}_i^\ell \mathbf{A}_i^\ell$  is the  $i$ -th LoRA expert.

### 2.3 ALEATORIC AND EPISTEMIC UNCERTAINTIES

In machine learning models, uncertainty can be categorized into aleatoric (data-wise) and epistemic (model-wise) uncertainty (Hora, 1996; Hüllermeier & Waegeman, 2021). For LLMs, aleatoric uncertainty arises from the inherently ambiguous and context dependent nature of natural languages where a single phrase or sentence can have multiple valid interpretations in different contexts. Epistemic uncertainty is introduced by the model’s lack of knowledge due to limited learning capabilities, suboptimal modeling or sparse training data.

Epistemic uncertainty is highly related to several well-known limitations of generative models. For example, it has been observed that when an LLM is pretrained on a diverse range of text data, it is generally well-calibrated, i.e. the predicted probability of the next token generally aligns with what is observed in real text. However, after fine-tuning or alignment with human preferences, the calibration error deteriorates (Zhao et al., 2021; Achiam et al., 2023a). A related phenomenon is forgetting, where the performance of a fine-tuned LLM diminishes on tasks outside the scope of the target downstream task (Lin et al., 2023; Luo et al., 2023).

Motivated by these observations, we explore functional-level epistemic uncertainty in generative models and aim to develop metrics that assess model performance on specific problem instances to fine-tune the parameter mixture.

## 3 METHODOLOGY

The high level goal of UQ4CT is to balance the exploration and exploitation of different LoRA experts during fine-tuning. In particular, we incorporate the functional-level epistemic uncertainty (FEU) to calibrate the prompt-dependent parameter mixture with LoRA MoE.

Assume that our answer  $a$  is generated via a mixture of mechanisms or models  $M$ , conditioning on the input prompt  $x$ . Assume that  $e(a)$  is an embedding of  $a$  so that least squares distance is a natural distance on the space of  $e(a)$ . For an expressive enough model class  $\mathcal{M}$ , and a calibrated distribution  $P(M|x)$  over the model class, we can measure the deviation of the generated answer to the “ideal” one  $a_* = f_*(x)$  as:

$$\begin{aligned} & \mathbb{E}_{M \sim P(M|x)} \left[ \mathbb{E}_{a \sim P(a|M,x)} \left[ \|e(a) - e(a_*)\|^2 \right] \right] \\ &= \mathbb{E}_{P(M|x)} \underbrace{\left[ \mathbb{E}_{P(a|M,x)} \left[ \|e(a) - \mathbb{E}_{P(a|M,x)}[e(a)]\|^2 \right] \right]}_{\text{Aleatoric Uncertainty}} + \underbrace{\mathbb{E}_{P(M|x)} \left[ \|\mathbb{E}_{P(a|M,x)}[e(a)] - e(a_*)\|^2 \right]}_{\text{Epistemic Uncertainty}}. \end{aligned} \quad (4)$$

Note that we hereby quantify uncertainty as a function of the input prompt  $x$ , since the distribution of model  $M$  conditions on  $x$ . We hence name the task “functional-level uncertainty quantification”.

### 3.1 FUNCTIONAL-LEVEL EPISTEMIC UNCERTAINTY

Motivated by the decomposition in Eq. (4) for least squares loss, we may consider a general distance  $d$ , and define epistemic uncertainty that characterizes the variation caused by model training procedure. Specifically, we focus on the variation introduced in the model fine-tuning stage of LLMs.

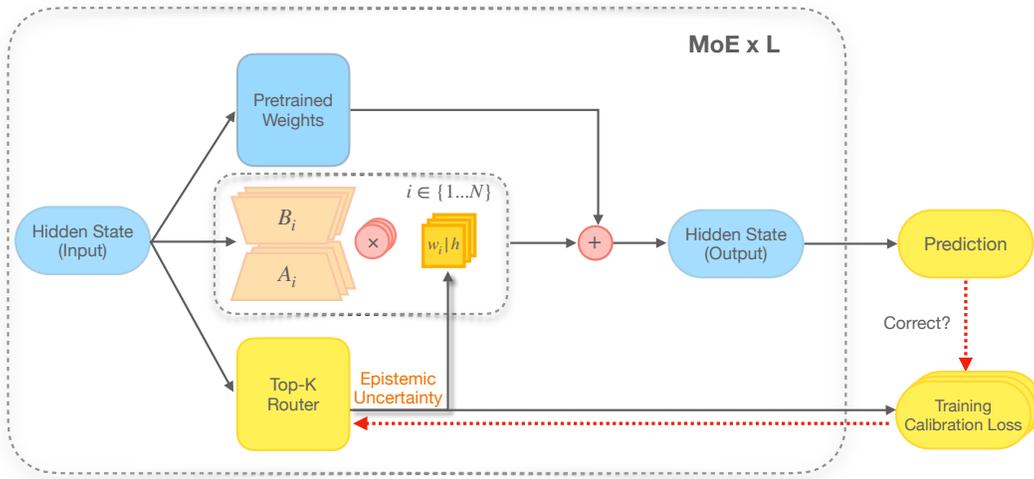


Figure 1: Mixture of Experts (MoE) architecture to capture and calibrate functional-level epistemic uncertainty. Experts  $B_{1 \dots N} A_{1 \dots N}$  capture functional relationships in the data throughout fine-tuning, the weights  $w_{1 \dots N} | h$  quantify the uncertainty in selecting these functional bases conditioned on the input hidden state  $h$ , which is the semantic representation of the input token  $x$ . In the UQ4CT workflow, we align this uncertainty with predictive correctness. When the router makes a correct prediction, the loss reinforces this decision, thereby increasing the router confidence in its selection, which aligns with a lower epistemic uncertainty. Conversely, when the router makes an incorrect prediction, the loss penalizes this selection, potentially causing the router to distribute its probabilities more broadly across experts, which is indicative of higher epistemic uncertainty.

Mathematically, given prompt  $x$ , we consider the following definition of epistemic uncertainty:

$$\text{Epistemic Uncertainty} = \mathbb{E}_{M \sim P(M|x)} \mathbb{E}_{M' \sim P(M'|x)} [\|\mathbb{E}_{P(a|M,x)}[e(a)] - \mathbb{E}_{P(a'|M',x)}[e(a')]\|^2]. \quad (5)$$

Here,  $a'$  represents the ground truth output sampled from the ideal MoE model  $M'$  conditioned on the prompt  $x$ . The epistemic uncertainty measures the least squares distance between  $e(a)$  and  $e(a')$  from current mixture distribution  $P(M|x)$  and the ideal mixture distribution  $P(M'|x)$ .

### 3.2 QUANTIFYING FEU WITH LORA MOE FRAMEWORK

We quantify the functional-level epistemic uncertainty (FEU) with the MoE architecture, which is represented by the embedding  $e(a)$  in Eq. (5). As shown in Figure 1, the LoRA experts  $B_{1 \dots N} A_{1 \dots N}$  capture important functional relationships in the data during fine-tuning. We treat these functions represented by the LoRA experts as basis functions  $f_{1 \dots N}$  and define them as follows:

$$f_1, f_2, \dots, f_N = \{B_1 A_1, B_2 A_2, \dots, B_N A_N\}. \quad (6)$$

Conditional on the input prompts, the more complex functional relationships that recursively map inputs to outputs are represented as linear combinations (mixture of experts) of the basis functions. Uncertainty quantification in this functional space reduces to quantifying the uncertainty of the weights over the basis functions. In particular, the weights  $w_{1 \dots N} | h$  from the top-k router dynamically selects these basis functions conditioned on the input hidden state  $h$ :

$$h = \sum_{i=1}^N (w_i | h) \cdot f_i = \sum_{i=1}^N (w_i | h) \cdot (B_i A_i). \quad (7)$$

The top-k weights that produce the final output hidden state quantify the functional level epistemic uncertainty in the function selection.

In the mixture of LoRA experts architecture, we follow the routing mechanisms of the MoE layers as in Eq. (2) and (3). Specifically, we employ top-2 gate routers, which chooses the 2 experts with

highest probability given a hidden state  $\mathbf{h}_i^\ell$ :

$$R^\ell(\mathbf{h}_i^\ell) = \text{KeepTop-2}(\text{Softmax}(\mathbf{W}_r^\ell \cdot \mathbf{h}_i^\ell)). \quad (8)$$

Given an input prompt  $x$  with length  $s$ , we model functional-level epistemic uncertainty (FEU) by aggregating  $R^\ell(\mathbf{h}_i^\ell)$  over both layer dimension and sequence dimension:

$$\text{FEU}(x) = \frac{1}{s} \sum_{i=1}^s \left[ \frac{1}{L} \sum_{\ell=1}^L R^\ell(\mathbf{h}_i^\ell) \right] \quad (9)$$

### 3.3 TRAINING CALIBRATION LOSS

We then calibrate the FEU model of the epistemic uncertainty against predictive accuracy. Specifically for the MoE top-k routers, we design the following calibration loss for training:

$$\mathcal{L}_{\text{cal}} = \|\mathbb{1}\{\text{MixLoRA}(x) = y^*\} - \text{FEU}(x)\|^2, \quad (10)$$

where the first term is an indicator function of whether the model prediction matches the ground truth  $y^*$  given the prompt  $x$ . Here, the indicator function resembles  $\mathbb{E}_{P(a'|M',x)}[e(a')]$  in Equation 5, where the ground truth  $y^*$  is  $a'$  and the indicator function maps the predictive correctness to a confidence space  $e \in [0, 1]$ . We employ a one-hot definition of the ground truth confidence. When the prediction from current mixture model matches the ground truth, the ground truth confidence is 1. Otherwise, when the predictions do not match, the ground truth confidence is 0.

As shown in Figure 1, this term effectively promotes expert exploitation for correct predictions and expert exploration for incorrect predictions by directly calibrating the functional level epistemic uncertainty to align with the predictive correctness. Ideally, when the  $N$  LoRA experts together capture all the functional relationships across the data distribution with cross entropy during fine-tuning, our proposed loss  $\mathcal{L}_{\text{cal}}$  also finds proper mixture of LoRA experts conditioned on the input  $x$  by conditionally promoting expert exploitation and exploration. This allows the model to select correct functional relationships regarding  $x$  to generate an output that better matches the data distribution, which grants calibrated uncertainty estimations.

Load balancing is a common technique to ensure even exploitation across experts with the MoE architecture (Fedus et al., 2022). We follow the load balancing loss  $\mathcal{L}_b$  proposed by (Li et al., 2024) and define our loss function as:

$$\mathcal{L} = \text{CE} + \alpha \cdot \mathcal{L}_b + \beta \cdot \mathcal{L}_{\text{cal}}, \quad (11)$$

where  $\text{CE}$  represents cross entropy loss,  $\alpha$  and  $\beta$  are the hyperparameters of two auxiliary terms. Details about  $\mathcal{L}_b$  can be found in Appendix A.1.

## 4 RELATED WORK

### 4.1 PARAMETER-EFFICIENT FINE TUNING FOR LLMs

Large Language Models (Brown et al., 2020; Chowdhery et al., 2022; Hoffmann et al., 2022; Touvron et al., 2023a;d) have shown impressive abilities in handling various natural language processing tasks. Building on these advances, instruction fine-tuning (Chung et al., 2022; Iyer et al., 2022; Zheng et al., 2024) has enhanced LLMs' capacity to comprehend and follow human instructions, forming the core of modern conversational AI systems (Wu et al., 2023b; Achiam et al., 2023b). However, as LLMs increase in size, the fine-tuning process demands much more time and memory.

To address these challenges, several strategies have been proposed, including parameter-efficient fine-tuning (PEFT) (Mangrulkar et al., 2022), model distillation (Liu et al., 2023; Xiao et al., 2023), quantization (Frantar et al., 2022; Xiao et al., 2022a), and pruning (Frantar & Alistarh, 2023; Ma et al., 2023). Among these, LoRA (Hu et al., 2021b), which leverages low-rank matrix decomposition of linear layer weights, is a widely adopted PEFT technique that boosts model performance without adding computational costs during inference. For example, VeRA (Kopiczko et al., 2023) introduces learnable scaling vectors to modify shared pairs of frozen random matrices across layers, while FedPara (Hyeon-Woo et al., 2021) focuses on low-rank Hadamard products for federated learning

270 settings. Tied-Lora(Renduchintala et al., 2023) applies weight tying to further minimize the number  
271 of trainable parameters. AdaLoRA(Zhang et al., 2023) uses Singular Value Decomposition (SVD) to  
272 prune less important singular values for efficient updates, and DoRA(Liu et al., 2024) separates pre-  
273 trained weights into magnitude and direction components, applying LoRA to update the directional  
274 component during fine-tuning, thus reducing the number of parameters to be trained.

## 275 276 277 4.2 MIXTURE OF EXPERTS

278  
279 The Mixture-of-Experts concept(Jacobs et al., 1991), introduced as early as 1991, presented a novel  
280 supervised learning framework where multiple networks (experts) specialize in handling distinct  
281 subsets of training data. Modern MoE variants adapt this by modifying the traditional feed-forward  
282 sub-layer within transformer blocks, incorporating sparsely activated LoRA experts, which allows  
283 for significant expansion in model width without a proportional increase in computational overhead.

284 Different MoE architectures have since emerged, distinguished by their expert sampling and rout-  
285 ing strategies. For example, LLaVA-MoLE(Chen et al., 2024) improves token routing to domain-  
286 specific experts within transformer layers, reducing data conflicts and consistently outperforming  
287 standard LoRA baselines. Other MoE-based approaches include MoRAL(Yang et al., 2024b),  
288 which focuses on efficiently adapting LLMs to new domains and tasks for lifelong learning, and  
289 LoRAMoE(Dou et al., 2024), which incorporates LoRAs via a router network to mitigate the is-  
290 sue of world knowledge forgetting. PESC(Wu et al., 2024a) transforms dense models into sparse  
291 ones through an MoE structure, lowering computational and GPU memory requirements. MoE-  
292 LoRA(Luo et al., 2024) introduces a new parameter-efficient MoE method using Layer-wise Ex-  
293 pert Allocation (MoLA) for transformer models, while MoCLE(Gou et al., 2023) activates task-  
294 specific model parameters based on instruction clusters. MixLoRA (Li et al., 2024) implements a  
295 high-throughput framework for LoRA MoE training and inference process, constructing LoRAs as  
296 stochastic experts to reduce computational overhead while expanding model capacity.

297 Despite the performance improvements these architecture advancements have brought, the overcon-  
298 fidence problem of fine-tuned models is lacking attention (Xiao et al., 2022c; He et al., 2023; Tian  
299 et al., 2023; OpenAI, 2023). Enhancing the uncertainty estimation capabilities of these models is  
300 fundamental toward more reliable, interpretable and trustworthy applications of LLMs.

## 301 302 303 4.3 UNCERTAINTY QUANTIFICATION IN LLMs

304  
305 Uncertainty quantification has garnered substantial attention in various tasks and domains within  
306 neural networks(Gal & Ghahramani, 2015; Gal & Ghahramani, 2016a; Malinin & Gales, 2018;  
307 Ovadia et al., 2019; Malinin et al., 2021; Lin et al., 2022; Kuhn et al., 2023; Lin et al., 2023). This  
308 focus extends to LLMs, where the precise quantification of prediction uncertainty has become a  
309 critical area of research(Xiao et al., 2022b; Lin et al., 2022; Mielke et al., 2022; Chen & Mueller,  
310 2023; Duan et al., 2023; Huang et al., 2023). LLMs, particularly in generative tasks, pose unique  
311 challenges, especially when it comes to measuring the uncertainty of their outputs(Liu et al., 2019;  
312 Malinin & Gales, 2021; Kuhn et al., 2023; Lin et al., 2023). The distinction between aleatoric  
313 and epistemic uncertainty was recently examined in the context of LLMs(Hou et al., 2023), though  
314 this was approached by ensembling model inputs rather than model instances and did not address  
315 fine-tuning tasks specifically.

316 Existing works have investigated the application of ensembling in fine-tuning LLMs for uncertainty  
317 quantification. Gleave & Irving (2022); Sun et al. (2022) focus on uncertainty estimation in full  
318 model fine-tuning, while this approach inherently incurs significant memory overhead. Wang et al.  
319 (2023); Zhai et al. (2023b); Balabanov & Linander (2024) explore the use of LoRA ensembles  
320 for uncertainty estimation in LLMs. Yang et al. (2024a) applies a post-hoc Laplace approxima-  
321 tion Mackay (1992) to model LoRA parameters after fine-tuning. BatchEnsemble(Wen et al., 2020),  
322 introduces component-specific rank-1 matrices as multiplicative modifications to a base model.  
323 Though this method has been applied to LLMs, it has been used in the pre-training phase rather  
than fine-tuning(Tran et al., 2022). None of these methods provide calibrations on epistemic uncer-  
tainty, which is crucial to mitigate overconfidence in the fine-tuning stage given the sparse dataset.

## 5 EXPERIMENTS

### 5.1 DATASETS

We include 5 multiple-choice question answering benchmarks to evaluate UQ4CT: OpenBookQA (OBQA, Mihaylov et al. (2018)), ARC-Easy (ARC-E) and ARC-Challenge (ARC-C) from AI2 Reasoning Challenge (Clark et al., 2018), BOOLQ (Clark et al., 2019) and ClimateQA, an expert-annotated domain specific benchmark for climate science. We also use computer science, law, medication and engineering subsets from MMLU dataset(Hendrycks et al., 2020) to evaluate performance under distribution shift. We fine-tune on the publicly available training split and test on the validation split from these benchmarks to evaluate model performance. [We report the average model performance over 3 random runs and the standard deviations in the subscript.](#)

### 5.2 EXPERIMENT SETUP

We implement UQ4CT with PyTorch (Paszke et al., 2019), extending the MixLoRA repository in (Li et al., 2024) and compare the average performance in three random runs and report the mean and standard deviation with following baselines. We use the publicly available LLaMA-2-7B-hf (Touvron et al., 2023c) as our base model. In particular, we apply MixLoRA to query, key, value and output layers, together with the feed-forward networks in LLaMA-2-7B-hf (gate layer, down layer and up layer). Details are provided in Appendix A.4

- **LoRA**(Hu et al., 2021a). We use standard LoRA fine-tuning as lower performance bound.
- **Monte Carlo (MC) Dropout**(Gal & Ghahramani, 2016b) keeps dropout on at both training and testing time. By performing multiple forward passes, MC dropout randomly shuts down a portion of model nodes, producing ensemble-alike predictions. To combine LoRA fine-tuning with MC dropout, we apply dropout on the input of the LoRA adapter, following the implementation of Mangrulkar et al. (2022).
- **Deep Ensemble**(Lakshminarayanan et al., 2017) averages the predictions from each ensemble member which have been trained with varying random initialization. We combine deep ensemble with LoRA by fine-tuning 3 randomly initialized LoRAs together and ensembling their output as final predictions.
- **Laplace-LoRA (LA)**(Yang et al., 2024a) applies a post-hoc Laplace approximation on fine-tuned LoRA parameters for better uncertainty estimation.
- **MixLoRA**(Li et al., 2024) incorporates LoRAs via a router network to reduce computational overhead while expanding model capacity. We add this as a baseline to resemble plain LoRA MoE model performance.

**Evaluation Metrics.** We measure the prediction accuracy on the validation set for all 5 tasks. For uncertainty calibration, we incorporate expected calibration error (ECE, Guo et al. (2017)) with 15 bins, which measures the alignment between predicted probabilities and empirical accuracy. We also investigate model performance under distribution shift to ensure the model has predictable behavior when given data from other domains as this is a crucial component for real-world applications. Specifically, we test models fine-tuned on OBQA dataset with 4 domain-specific MMLU subtask ensembles focusing on different professionalities and ARC-E/C datasets to approximate larger and smaller distribution shifts. Metric details are provided in Appendix A.5.

### 5.3 RESULTS

We assess the prediction accuracy and uncertainty calibration of models under both in-distribution and distribution shift scenarios. The in-distribution scenario examines the alignment of the fine-tuned model on the target downstream task, while the distribution shift scenario evaluates the generalizability of the model on novel tasks beyond the fine-tuned domain. These two scenarios combined provides a comprehensive assessment of model robustness in real-world applications where it is essential for the model to excel on its primary task while maintaining the ability to effectively handle unforeseen or out-of-distribution inputs.

Table 1: Performance comparison of different methods fine-tuned with LLaMA2-7B across 4 common sense reasoning tasks and a domain-specific task. UQ4CT shows substantial ECE improvements while maintaining high accuracy.

Metrics	Methods	BoolQ	ARC-E	ARC-C	OBQA	ClimateQA
ACC $\uparrow$	LoRA	69.5 <sub>1.93</sub>	74.8 <sub>1.39</sub>	53.8 <sub>0.6</sub>	72.1 <sub>0.87</sub>	59.9 <sub>2.13</sub>
	MC Drop	66.8 <sub>3.66</sub>	76.8 <sub>1.30</sub>	50.9 <sub>2.01</sub>	74.8 <sub>1.34</sub>	58.2 <sub>2.11</sub>
	Ensemble	66.2 <sub>3.7</sub>	71.2 <sub>1.0</sub>	47.5 <sub>0.57</sub>	75.5 <sub>1.4</sub>	59.6 <sub>6.9</sub>
	LA	68.7 <sub>1.32</sub>	74.6 <sub>2.11</sub>	51.4 <sub>0.83</sub>	70.8 <sub>1.24</sub>	55.2 <sub>3.29</sub>
	MixLoRA	71.5 <sub>1.05</sub>	77.7 <sub>2.27</sub>	54.3 <sub>1.07</sub>	75.5 <sub>2.91</sub>	61.6 <sub>1.76</sub>
	UQ4CT	73.5 <sub>0.52</sub>	76.6 <sub>1.30</sub>	52.8 <sub>1.77</sub>	77.3 <sub>1.36</sub>	63.3 <sub>1.74</sub>
ECE $\downarrow$	LoRA	11.9 <sub>0.78</sub>	11.9 <sub>2.04</sub>	19.4 <sub>4.75</sub>	10.2 <sub>1.07</sub>	14.3 <sub>0.64</sub>
	MC Drop	12.2 <sub>0.85</sub>	11.9 <sub>1.99</sub>	19.8 <sub>4.85</sub>	10.9 <sub>0.24</sub>	14.3 <sub>0.56</sub>
	Ensemble	7.28 <sub>2.3</sub>	9.1 <sub>1.49</sub>	10.23 <sub>1.39</sub>	8.83 <sub>2.35</sub>	13.5 <sub>3.29</sub>
	LA	17.1 <sub>1.72</sub>	16.6 <sub>3.7</sub>	18.1 <sub>0.5</sub>	17.2 <sub>1.2</sub>	12.6 <sub>1.9</sub>
	MixLoRA	7.88 <sub>2.09</sub>	9.09 <sub>0.81</sub>	10.74 <sub>1.07</sub>	12.9 <sub>1.99</sub>	12.5 <sub>1.32</sub>
	UQ4CT	<b>2.3</b> <sub>0.82</sub>	<b>6.0</b> <sub>0.2</sub>	<b>6.1</b> <sub>1.11</sub>	<b>5.0</b> <sub>1.15</sub>	<b>8.1</b> <sub>0.52</sub>

### 5.3.1 IN-DISTRIBUTION PERFORMANCE

We first evaluate UQ4CT and baseline models fine-tuned on the 4 common sense reasoning tasks and the climate question answering task under the in-distribution scenario, where models are trained and evaluated on different splits of the same dataset. Note that one of the key advantages of UQ4CT is that uncertainty calibration happens during the fine-tuning stage with little computational overhead, while other UQ methods require repetitive sampling or other post-hoc complexities.

As shown in Table 1, UQ4CT demonstrates notable improvements in uncertainty calibration across a variety of tasks. Across all benchmarks, UQ4CT maintains competitive accuracy (ACC) compared to the baseline methods. For example, on the BoolQ and ClimateQA datasets, UQ4CT achieves accuracy rates of 73.5% and 63.3%, respectively. This empirically demonstrates that UQ4CT is capable of maintaining high accuracy with uncertainty calibration, which assures the gain in UQ performance does not compromise accuracy.

The most substantial performance improvement is observed in the reduction of Expected Calibration Error (ECE). UQ4CT consistently outperforms other methods, reducing ECE by more than 25% on average across the evaluated benchmarks. Unlike other methods where the ECE performance is worsened on the more challenging ARC-C benchmark, UQ4CT achieves an ECE score of 6.1, showcasing the effectiveness of the calibration.

In addition to experiments on LLaMA-2-7B in the main text, we also present additional experiments on fine-tuning Mistral-7B in Appendix A.2 for more comprehensive evaluation of our method. For both LLaMA-2-7B and Mistral-7B models, UQ4CT consistently shows substantial improvements in uncertainty calibration across various tasks. The improvements are critical in applications where the model’s confidence must align with its predictive accuracy given limited data, particularly in safety-critical and domain-specific tasks.

### 5.3.2 PERFORMANCE UNDER DISTRIBUTION SHIFT

Due to the sparse nature of the fine-tuning data, real world deployment of LLMs often requires the model to be robust to out-of-distribution knowledge (Ouyang et al., 2022; Touvron et al., 2023b;c). Therefore, we evaluate the performance of UQ4CT along with other baseline models fine-tuned on the OBQA dataset under smaller and larger distribution shift scenarios. Similar to the dataset setup in (Yang et al., 2024a), we use ARC-C and ARC-E dataset to simulate smaller distribution shift because the ARC dataset has similar domain focus on general science reasoning, but is generally more challenging and covers a broader range of scientific topics than OBQA. For larger distribution shift, we ensemble the domain-specific MMLU subtasks into 4 benchmarks focusing on different professionalities: Computer Science (CS), Engineering (Eng), Law and Health. These tasks have very broad coverage of the domain task at various knowledge levels ranging from elementary school

Table 2: Performance comparison of different methods fine-tuned on OBQA dataset with LLaMA2-7B across 2 smaller distribution shift (DS) tasks and 4 larger distribution shift tasks. UQ4CT shows substantial ECE improvements while maintaining high accuracy.

Metrics	Methods	ID	Smaller DS		Larger DS			
		OBQA	ARC-C	ARC-E	CS	Eng	Law	Health
ACC $\uparrow$	LoRA	72.1 <sub>1.87</sub>	58.6 <sub>1.93</sub>	66.5 <sub>3.38</sub>	35.5 <sub>2.35</sub>	30.8 <sub>1.72</sub>	34.9 <sub>1.41</sub>	39.1 <sub>1.52</sub>
	MC Drop	74.8 <sub>1.34</sub>	58.7 <sub>2.07</sub>	66.6 <sub>3.30</sub>	36.0 <sub>1.69</sub>	30.3 <sub>2.25</sub>	35.1 <sub>1.86</sub>	39.1 <sub>1.35</sub>
	Ensemble	75.5 <sub>1.4</sub>	57.7 <sub>0.78</sub>	69.1 <sub>0.48</sub>	36.7 <sub>2.18</sub>	30.3 <sub>1.13</sub>	35.3 <sub>1.02</sub>	39.9 <sub>1.89</sub>
	LA	70.8 <sub>1.24</sub>	58.7 <sub>0.58</sub>	67.9 <sub>0.41</sub>	33.7 <sub>1.22</sub>	29.6 <sub>1.32</sub>	35.4 <sub>0.75</sub>	38.5 <sub>1.61</sub>
	MixLoRA	75.5 <sub>2.91</sub>	58.5 <sub>1.44</sub>	69.2 <sub>1.02</sub>	35.2 <sub>2.92</sub>	30.3 <sub>0.98</sub>	35.9 <sub>0.43</sub>	40.6 <sub>1.13</sub>
	UQ4CT	77.3 <sub>1.36</sub>	58.8 <sub>1.06</sub>	65.8 <sub>1.31</sub>	36.2 <sub>1.24</sub>	34.1 <sub>2.31</sub>	35.8 <sub>1.01</sub>	40.0 <sub>1.24</sub>
ECE $\downarrow$	LoRA	10.2 <sub>1.07</sub>	16.7 <sub>2.28</sub>	13.3 <sub>2.48</sub>	29.7 <sub>2.69</sub>	32.3 <sub>1.85</sub>	29.2 <sub>3.08</sub>	31.0 <sub>2.13</sub>
	MC Drop	10.9 <sub>0.24</sub>	16.7 <sub>2.20</sub>	13.2 <sub>2.21</sub>	23.2 <sub>2.32</sub>	31.6 <sub>1.64</sub>	28.0 <sub>2.93</sub>	25.9 <sub>2.27</sub>
	Ensemble	8.83 <sub>2.35</sub>	15.1 <sub>1.09</sub>	11.1 <sub>0.99</sub>	22.4 <sub>1.32</sub>	28.5 <sub>2.13</sub>	29.0 <sub>1.37</sub>	24.5 <sub>0.39</sub>
	LA	17.2 <sub>1.2</sub>	16.2 <sub>0.5</sub>	24.4 <sub>0.42</sub>	28.6 <sub>2.61</sub>	30.5 <sub>1.43</sub>	29.5 <sub>1.83</sub>	30.7 <sub>1.2</sub>
	MixLoRA	12.9 <sub>1.99</sub>	19.0 <sub>1.88</sub>	14.5 <sub>2.57</sub>	26.4 <sub>3.25</sub>	33.7 <sub>1.87</sub>	30.3 <sub>2.27</sub>	28.3 <sub>1.07</sub>
	UQ4CT	5.0 <sub>1.15</sub>	8.9 <sub>3.46</sub>	6.5 <sub>1.85</sub>	19.6 <sub>2.90</sub>	23.1 <sub>1.17</sub>	25.9 <sub>3.43</sub>	21.9 <sub>3.49</sub>

to professionals. This domain-specificity demonstrates larger distribution shift from OBQA, which is a general common sense reasoning task. Details of the ensemble is provided in Appendix. A.6.

The distribution shift evaluations are provided in Table 2. UQ4CT provides substantial improvements in terms of ECE while maintains similar accuracy for both smaller and larger distribution shift scenarios. For smaller distribution shifts, UQ4CT shows comparable ECE performance as the in-distribution scenario. For the more challenging larger distribution shifts, UQ4CT still achieves the best ECE performance among all baseline models. Note that UQ4CT also achieves competitive prediction accuracy across all domain-specific tasks. This empirically shows that our proposed alignment of the functional epistemic uncertainty with predictive correctness improves generalizability and mitigating the overconfidence problem on the fine-tuned model.

#### 5.4 ABLATION STUDY

In this section, we conduct ablation studies to investigate the effectiveness of our designed calibration loss,  $\mathcal{L}_{cal}$ . We first evaluate the incremental weighting performance of the calibration term, which investigates the effectiveness of  $\mathcal{L}_{cal}$  at the early stage of fine-tuning. Then, we perform a sensitivity test, where we explore the overall performance impact of  $\mathcal{L}_{cal}$ . We also conduct an ablation study on the impact of active LoRA experts in Appendix A.3.

##### 5.4.1 INCREMENTAL WEIGHTING ON CALIBRATION TERM

Due to the random initialization of LoRA experts, the predictions during early fine-tuning stage are likely to be incorrect as the model has little knowledge on the functional relationships regarding the data. Thus, it is intuitive to incrementally increase the weight parameter  $\beta$  over the calibration term  $\mathcal{L}_{cal}$  in the training loss for the LoRA experts to learn before calibration. We conduct this study by incrementally increase  $\beta$  from 0 to 1 within 50 gradient steps during the early stage of fine-tuning:

$$\beta = \min \left\{ 1, \frac{\text{current\_grad\_step}}{50} \right\}. \quad (12)$$

We choose 50 gradient steps from our observation that training loss generally stabilizes after 50 gradient steps, indicating the LoRA experts have learned some functional relationships from data.

As shown in Table 3, the incremental loss has significantly worse ECE performance across all tasks. This demonstrates the advantage of uncertainty calibration even in the early stage. In the beginning, the lack of functional relationships on the training data in LoRA experts lead to high epistemic uncertainty. Thus, UQ4CT encourages exploration over all LoRA experts while UQ4CT Incremental lacks it due to the small weighting in the beginning.

Table 3: Performance comparison of UQ4CT with and without incremental weighting. Incremental weighting has worse ECE performance while maintains similar accuracy.

Metrics	Methods	BoolQ	ARC-E	ARC-C	OBQA	ClimateQA
ACC $\uparrow$	UQ4CT	73.5 <sub>0.52</sub>	76.6 <sub>1.30</sub>	52.8 <sub>1.77</sub>	77.3 <sub>1.36</sub>	63.3 <sub>1.74</sub>
	UQ4CT-Incremental	72.0 <sub>0.19</sub>	75.4 <sub>0.81</sub>	54.6 <sub>0.95</sub>	77.6 <sub>0.43</sub>	60.2 <sub>3.17</sub>
ECE $\downarrow$	UQ4CT	<b>2.3<sub>0.82</sub></b>	<b>6.0<sub>0.2</sub></b>	<b>6.1<sub>1.11</sub></b>	<b>5.0<sub>1.15</sub></b>	<b>8.1<sub>0.52</sub></b>
	UQ4CT-Incremental	4.0 <sub>0.16</sub>	9.8 <sub>1.51</sub>	13.8 <sub>2.08</sub>	10.3 <sub>1.73</sub>	12.2 <sub>0.88</sub>

Table 4: Performance of UQ4CT with varying  $\beta$  value on OBQA dataset. Prediction accuracy and uncertainty alignment increases with  $\beta$ , highlighting the effectiveness of the calibration term.

$\beta$	ACC $\uparrow$	ECE $\downarrow$
0	75.5 <sub>2.91</sub>	12.9 <sub>1.99</sub>
0.2	76.0 <sub>0.6</sub>	7.4 <sub>7</sub> <sub>0.78</sub>
0.5	75.9 <sub>0.31</sub>	7.8 <sub>2</sub> <sub>0.93</sub>
0.8	76.6 <sub>0.4</sub>	5.9 <sub>4</sub> <sub>1.16</sub>
1	<b>77.3<sub>1.36</sub></b>	<b>5.0<sub>1.15</sub></b>

#### 5.4.2 SENSITIVITY TEST ON CALIBRATION TERM

To further understand the effectiveness of the calibration loss, we perform a sensitivity test of the coefficient  $\beta$  in Equation 11. This evaluates how our proposed calibration of parameter mixtures affect the overall model prediction and uncertainty quantification capabilities. We evaluate  $\beta$  values of 0, 0.2, 0.5, 0.8 and 1, where  $\beta = 0$  resembles the original MixLoRA method.

Results in Table 4 demonstrate the effectiveness of the calibration loss. When  $\beta = 0$ , the model is optimized without calibration on parameter mixtures, resulting in high ECE value. Even with small  $\beta = 0.2$  or  $\beta = 0.5$ , the ECE scores drastically improved compared to no calibration setting. Finally, when  $\beta = 1$ , the calibration term effectively optimizes the conditional parameter mixtures to generate outputs that fit data distribution well, resulting in lower ECE scores and higher accuracies.

## 6 DISCUSSION & CONCLUSION

In this work, we propose Functional-Level Uncertainty Quantification for Calibrated Fine-Tuning (UQ4CT), which addresses the overconfidence issues commonly encountered during fine-tuning of large language models. We present a functional perspective on quantifying epistemic uncertainty in LLMs and utilize it for uncertainty-calibrated fine-tuning. By incorporating functional-level epistemic uncertainty quantification with a mixture-of-experts framework, our proposed uncertainty-calibrated training loss effectively addresses the challenge of overconfidence in fine-tuned LLMs by significantly improving uncertainty calibration while maintaining high accuracy. Our evaluations demonstrate that UQ4CT reduces the Expected Calibration Error by more than 25% without compromising accuracy across a variety of downstream tasks, including common-sense and domain-specific reasoning, under in-distribution and out-of-distribution scenarios.

The limitation of UQ4CT lies in its dependency on predictive correctness. For general language modeling tasks such as chat completion, there lacks a clear metric on whether the response is correct or not. This limits the application of UQ4CT as naively token matching is a poor indicator of semantic correctness due to the ambiguous nature of language. For future work, we are exploring ways to adapt UQ4CT on open-ended problems that lacks a definitive optimization objective.

### REPRODUCIBILITY STATEMENT

We share our experimental details in Appendix A.4, and also provide the code and model weights for running experiments in the supplementary materials to reproduce our model performance results.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. *arXiv preprint arXiv:2303.08774*, 2023a.
- 545  
546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
547 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
548 report. *arXiv preprint arXiv: 2303.08774*, 2023b.
- 549  
550 Oleksandr Balabanov and Hampus Linander. Uncertainty quantification in fine-tuned llms using  
551 lora ensembles. *arXiv preprint arXiv:2402.12264*, 2024.
- 552  
553 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
554 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
555 wal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh,  
556 Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
557 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
558 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot  
559 learners. *ArXiv*, abs/2005.14165, 2020. URL [https://api.semanticscholar.org/  
560 CorpusID:218971783](https://api.semanticscholar.org/CorpusID:218971783).
- 561  
562 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan  
563 Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM  
564 Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.
- 565  
566 Hugh Chen, Scott Lundberg, and Su-In Lee. Checkpoint ensembles: Ensemble methods from a  
567 single training process. *arXiv preprint arXiv:1710.03282*, 2017.
- 568  
569 Jiuhai Chen and Jonas Mueller. Quantifying uncertainty in answers from any language model and  
570 enhancing their trustworthiness, 2023.
- 571  
572 Shaoxiang Chen, Zequn Jie, and Lin Ma. Llava-mole: Sparse mixture of lora experts for mitigating  
573 data conflicts in instruction finetuning mllms. *arXiv preprint arXiv: 2401.16160*, 2024.
- 574  
575 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
576 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
577 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay,  
578 Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson,  
579 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,  
580 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier  
581 García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David  
582 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shiv-  
583 ani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie  
584 Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee,  
585 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Ja-  
586 son Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel.  
587 Palm: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24, 2022. URL  
588 <https://api.semanticscholar.org/CorpusID:247951931>.
- 589  
590 Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi  
591 Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun  
592 Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gau-  
593 rav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav  
594 Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and  
595 Jason Wei. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416, 2022. URL  
596 <https://api.semanticscholar.org/CorpusID:253018554>.
- 597  
598 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
599 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint  
600 arXiv:1905.10044*, 2019.

- 594 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
595 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
596 *arXiv preprint arXiv:1803.05457*, 2018.  
597
- 598 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
599 Chen, Chi-Min Chan, Weize Chen, et al. Delta tuning: A comprehensive study of parameter  
600 efficient methods for pre-trained language models. *arXiv preprint arXiv:2203.06904*, 2022.
- 601 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
602 Chen, Chi-Min Chan, Weize Chen, Jing Yi, Weilin Zhao, Xiaozhi Wang, Zhiyuan Liu, Hai-Tao  
603 Zheng, Jianfei Chen, Yang Liu, Jie Tang, Juanzi Li, and Maosong Sun. Parameter-efficient fine-  
604 tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235,  
605 2023.  
606
- 607 Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi,  
608 Xiao Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing  
609 Huang. Loramoe: Alleviate world knowledge forgetting in large language models via moe-style  
610 plugin. *arXiv*, 2024.
- 611 Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura,  
612 and Kaidi Xu. Shifting attention to relevance: Towards the uncertainty estimation of large lan-  
613 guage models, 2023.
- 614 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter  
615 models with simple and efficient sparsity. *JMLR*, 2022.  
616
- 617 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in  
618 one-shot. *ArXiv*, abs/2301.00774, 2023. URL [https://api.semanticscholar.org/  
619 CorpusID:255372747](https://api.semanticscholar.org/CorpusID:255372747).
- 620 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training  
621 quantization for generative pre-trained transformers. *ArXiv*, abs/2210.17323, 2022. URL  
622 <https://api.semanticscholar.org/CorpusID:253237200>.
- 623
- 624 Yarin Gal and Zoubin Ghahramani. Bayesian Convolutional Neural Networks with Bernoulli Ap-  
625 proximate Variational Inference. *arXiv e-prints*, art. arXiv:1506.02158, June 2015.  
626
- 627 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model un-  
628 certainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings  
629 of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Ma-  
630 chine Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016a. PMLR.  
631 URL <https://proceedings.mlr.press/v48/gall16.html>.
- 632 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model  
633 uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059.  
634 PMLR, 2016b.
- 635
- 636 Adam Gleave and Geoffrey Irving. Uncertainty estimation for language reward models, 2022.
- 637 Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James T  
638 Kwok, and Yu Zhang. Mixture of cluster-conditional lora experts for vision-language instruction  
639 tuning. *arXiv preprint arXiv: 2312.12379*, 2023.  
640
- 641 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
642 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- 643 Guande He, Jianfei Chen, and Jun Zhu. Preserving pre-trained features helps calibrate fine-tuned  
644 language models. In *ICLR*, 2023.  
645
- 646 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
647 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint  
arXiv:2009.03300*, 2020.

- 648 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
649 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
650 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aure-  
651 lia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and  
652 L. Sifre. Training compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022. URL  
653 <https://api.semanticscholar.org/CorpusID:247778764>.
- 654 Stephen C Hora. Aleatory and epistemic uncertainty in probability elicitation with an example from  
655 hazardous waste management. *Reliability Engineering & System Safety*, 54(2-3):217–223, 1996.  
656
- 657 Bairu Hou, Yujian Liu, Kaizhi Qian, Jacob Andreas, Shiyu Chang, and Yang Zhang. Decomposing  
658 uncertainty for large language models through input clarification ensembling, 2023.
- 659 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-  
660 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp.  
661 In *International Conference on Machine Learning*, pp. 2790–2799. PMLR, 2019.  
662
- 663 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
664 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*  
665 *arXiv:2106.09685*, 2021a.
- 666 J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and  
667 Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685,  
668 2021b. URL <https://api.semanticscholar.org/CorpusID:235458009>.
- 669 Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu,  
670 and Lei Ma. Look Before You Leap: An Exploratory Study of Uncertainty Measurement for  
671 Large Language Models, October 2023. URL <http://arxiv.org/abs/2307.10236>.  
672 arXiv:2307.10236 [cs].  
673
- 674 Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning:  
675 An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.  
676
- 677 Nam Hyeon-Woo, Moon Ye-Bin, and Tae-Hyun Oh. Fedpara: Low-rank hadamard product for  
678 communication-efficient federated learning. *arXiv preprint arXiv: 2108.06098*, 2021.
- 679 Srinivas Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu,  
680 Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O’Horo, Gabriel  
681 Pereyra, Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Veselin Stoy-  
682 anov. Opt-impl: Scaling language model instruction meta learning through the lens of gener-  
683 alization. *ArXiv*, abs/2212.12017, 2022. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:255096269)  
684 [CorpusID:255096269](https://api.semanticscholar.org/CorpusID:255096269).
- 685 Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures  
686 of local experts. *Neural Computation*, 1991.  
687
- 688 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
689 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
690 Mistral 7b. *arXiv preprint arXiv: 2310.06825*, 2023.
- 691 Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. Vera: Vector-based random  
692 matrix adaptation. *arXiv preprint arXiv: 2310.11454*, 2023.  
693
- 694 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic Uncertainty: Linguistic Invariances for  
695 Uncertainty Estimation in Natural Language Generation, April 2023. URL [http://arxiv.](http://arxiv.org/abs/2302.09664)  
696 [org/abs/2302.09664](http://arxiv.org/abs/2302.09664). arXiv:2302.09664 [cs].
- 697 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
698 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,  
699 30, 2017.  
700
- 701 Andrew Kyle Lampinen, Stephanie C. Y. Chan, Ishita Dasgupta, Andrew J. Nam, and Jane X. Wang.  
Passive learning of active causal strategies in agents and language models, May 2023.

- 702 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
703 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional  
704 computation and automatic sharding. In *ICLR*, 2020.
- 705
- 706 Dengchun Li, Yingzi Ma, Naizheng Wang, Zhiyuan Cheng, Lei Duan, Jie Zuo, Cal Yang, and  
707 Mingjie Tang. Mixlor: Enhancing large language models fine-tuning with lora based mixture of  
708 experts. *arXiv preprint arXiv:2404.15159*, 2024.
- 709
- 710 Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton Wang, Linxi Fan, Tao Chen, De-An Huang,  
711 Ekin Akyürek, Anima Anandkumar, et al. Pre-trained language models for interactive decision-  
712 making. *Advances in Neural Information Processing Systems*, 35:31199–31212, 2022.
- 713
- 714 Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching Models to Express Their Uncertainty in  
715 Words. *arXiv e-prints*, art. arXiv:2205.14334, May 2022. doi: 10.48550/arXiv.2205.14334.
- 716
- 717 Yong Lin, Lu Tan, Hangyu Lin, Zeming Zheng, Renjie Pi, Jipeng Zhang, Shizhe Diao, Haoxiang  
718 Wang, Han Zhao, Yuan Yao, et al. Speciality vs generality: An empirical study on catastrophic  
719 forgetting in fine-tuning foundation models. *arXiv preprint arXiv:2309.06256*, 2023.
- 720
- 721 Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with Confidence: Uncertainty Quanti-  
722 fication for Black-box Large Language Models. *arXiv e-prints*, art. arXiv:2305.19187, May 2023.  
723 doi: 10.48550/arXiv.2305.19187.
- 724
- 725 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and  
726 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context  
727 learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- 728
- 729 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
730 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv  
731 preprint arXiv: 2402.09353*, 2024.
- 732
- 733 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
734 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
735 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 736
- 737 Zechun Liu, Barlas Oğuz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang  
738 Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware  
739 training for large language models. *ArXiv*, abs/2305.17888, 2023. URL [https://api.  
740 semanticscholar.org/CorpusID:258959117](https://api.semanticscholar.org/CorpusID:258959117).
- 741
- 742 Tongxu Luo, Jiahe Lei, Fangyu Lei, Weihao Liu, Shizhu He, Jun Zhao, and Kang Liu. Moelora:  
743 Contrastive learning guided mixture of experts on parameter-efficient fine-tuning for large lan-  
744 guage models. *arXiv preprint arXiv: 2402.12851*, 2024.
- 745
- 746 Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study  
747 of catastrophic forgetting in large language models during continual fine-tuning. *arXiv preprint  
748 arXiv:2308.08747*, 2023.
- 749
- 750 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large  
751 language models. *ArXiv*, abs/2305.11627, 2023. URL [https://api.semanticscholar.  
752 org/CorpusID:258823276](https://api.semanticscholar.org/CorpusID:258823276).
- 753
- 754 D. J. C. Mackay. Information-based objective functions for active data selection. *Neural Computa-  
755 tion*, 4(2):550–604, 1992.
- 756
- 757 Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks, 2018.
- 758
- 759 Andrey Malinin and Mark Gales. Uncertainty Estimation in Autoregressive Structured Prediction,  
760 February 2021. URL <http://arxiv.org/abs/2002.07650>. arXiv:2002.07650 [cs, stat].
- 761
- 762 Andrey Malinin, Liudmila Prokhorenkova, and Aleksei Ustimenko. Uncertainty in gradient boosting  
763 via ensembles, 2021.

- 756 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul.  
757 Peft: State-of-the-art parameter-efficient fine-tuning methods. [https://github.com/  
758 huggingface/peft](https://github.com/huggingface/peft), 2022.
- 759 Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational  
760 agents’ overconfidence through linguistic calibration. *Transactions of the Association for  
761 Computational Linguistics*, 10:857–872, 2022. doi: 10.1162/tacl.a.00494. URL [https:  
762 //aclanthology.org/2022.tacl-1.50](https://aclanthology.org/2022.tacl-1.50).
- 764 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
765 electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
- 766 OpenAI. GPT-4 technical report, 2023.
- 768 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
769 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
770 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:  
771 27730–27744, 2022.
- 772 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua V.  
773 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?  
774 evaluating predictive uncertainty under dataset shift, 2019.
- 776 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
777 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-  
778 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 779 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning  
780 with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- 782 Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan,  
783 and Peter J Liu. Out-of-distribution detection and selective generation for conditional language  
784 models. In *The Eleventh International Conference on Learning Representations*, 2022.
- 785 Adithya Renduchintala, Tugrul Konuk, and Oleksii Kuchaiev. Tied-lora: Enhancing parameter effi-  
786 ciency of lora with weight tying. *arXiv preprint arXiv: 2311.09578*, 2023.
- 788 Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan  
789 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul  
790 Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera  
791 y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad  
792 Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam,  
793 and Vivek Natarajan. Large Language Models Encode Clinical Knowledge. *arXiv*, 2022.
- 794 Meiqi Sun, Wilson Yan, Pieter Abbeel, and Igor Mordatch. Quantifying uncertainty in foundation  
795 models via ensembles. In *NeurIPS 2022 Workshop on Robustness in Sequence Modeling*, 2022.  
796 URL <https://openreview.net/forum?id=LpBlkATV24M>.
- 798 Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea  
799 Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting calibrated  
800 confidence scores from language models fine-tuned with human feedback. *arXiv preprint  
801 arXiv:2305.14975*, 2023.
- 802 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
803 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-  
804 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation lan-  
805 guage models. *ArXiv*, abs/2302.13971, 2023a. URL [https://api.semanticscholar.  
806 org/CorpusID:257219404](https://api.semanticscholar.org/CorpusID:257219404).
- 808 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
809 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023b.

- 810 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
811 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
812 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023c.
- 813
- 814 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
815 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas  
816 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,  
817 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S.  
818 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian  
819 Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut  
820 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,  
821 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta,  
822 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh  
823 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov,  
824 Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert  
825 Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat  
826 models. *ArXiv*, abs/2307.09288, 2023d. URL [https://api.semanticscholar.org/  
CorpusID:259950998](https://api.semanticscholar.org/CorpusID:259950998).
- 827
- 828 Dustin Tran, Jeremiah Liu, Michael W. Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han,  
829 Zi Wang, Zeldia Mariet, Huiyi Hu, Neil Band, Tim G. J. Rudner, Karan Singhal, Zachary Nado,  
830 Joost van Amersfoort, Andreas Kirsch, Rodolphe Jenatton, Nithum Thain, Honglin Yuan, Kelly  
831 Buchanan, Kevin Murphy, D. Sculley, Yarin Gal, Zoubin Ghahramani, Jasper Snoek, and Balaji  
832 Lakshminarayanan. Plex: Towards reliability using pretrained large model extensions, 2022.
- 833
- 834 Xi Wang, Laurence Aitchison, and Maja Rudolph. Lora ensembles for large language model fine-  
835 tuning. *arXiv preprint arXiv:2310.00035*, 2023.
- 836
- 837 Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: An alternative approach to efficient  
838 ensemble and lifelong learning, 2020.
- 839
- 840 Haoyuan Wu, Haisheng Zheng, and Bei Yu. Parameter-efficient sparsity crafting from dense to  
841 mixture-of-experts for instruction tuning on general tasks. *arXiv preprint arXiv: 2401.02731*,  
842 2024a.
- 843
- 844 Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prab-  
845 hanjan Kambadur, David Rosenberg, and Gideon Mann. BloombergGPT: A Large Language  
846 Model for Finance, May 2023a.
- 847
- 848 Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief  
849 overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal  
850 of Automatica Sinica*, 10(5):1122–1136, 2023b.
- 851
- 852 Xun Wu, Shaohan Huang, and Furu Wei. Mixture of lora experts. *arXiv preprint arXiv:2404.13628*,  
853 2024b.
- 854
- 855 Guangxuan Xiao, Ji Lin, Mickael Seznec, Julien Demouth, and Song Han. Smoothquant: Accurate  
856 and efficient post-training quantization for large language models. *ArXiv*, abs/2211.10438, 2022a.  
857 URL <https://api.semanticscholar.org/CorpusID:253708271>.
- 858
- 859 Guangxuan Xiao, Ji Lin, and Song Han. Offsite-tuning: Transfer learning without full  
860 model. *ArXiv*, abs/2302.04870, 2023. URL [https://api.semanticscholar.org/  
CorpusID:256697131](https://api.semanticscholar.org/CorpusID:256697131).
- 861
- 862 Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-  
863 Philippe Morency. Uncertainty quantification with pre-trained language models: A large-  
scale empirical analysis. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Find-  
ings of the Association for Computational Linguistics: EMNLP 2022*, pp. 7273–7284, Abu  
 Dhabi, United Arab Emirates, December 2022b. Association for Computational Linguistics.  
 doi: 10.18653/v1/2022.findings-emnlp.538. URL [https://aclanthology.org/2022.  
findings-emnlp.538](https://aclanthology.org/2022.findings-emnlp.538).

864 Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-  
865 Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scale  
866 empirical analysis. In *EMNLP*, 2022c.

867 Adam X. Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation  
868 for large language models, 2024a.

869  
870 Shu Yang, Muhammad Asif Ali, Cheng-Long Wang, Lijie Hu, and Di Wang. Moral: Moe augmented  
871 lora for llms’ lifelong learning. *arXiv preprint arXiv: 2402.11260*, 2024b.

872  
873 Yuanzhao Zhai, Han Zhang, Yu Lei, Yue Yu, Kele Xu, Dawei Feng, Bo Ding, and Huaimin Wang.  
874 Uncertainty-penalized reinforcement learning from human feedback with diverse reward lora en-  
875 sembles. *arXiv preprint arXiv:2401.00243*, 2023a.

876  
877 Yuanzhao Zhai, Han Zhang, Yu Lei, Yue Yu, Kele Xu, Dawei Feng, Bo Ding, and Huaimin Wang.  
878 Uncertainty-penalized reinforcement learning from human feedback with diverse reward lora en-  
879 sembles, 2023b.

880  
881 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and  
882 Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *ICLR*, 2023.

883  
884 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving  
885 few-shot performance of language models. In *International conference on machine learning*, pp.  
886 12697–12706. PMLR, 2021.

887  
888 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
889 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
890 chatbot arena. *NeurIPS*, 36, 2024.

891  
892 Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein. Adapting language models for zero-shot  
893 learning by meta-tuning on dataset and prompt collections. *arXiv preprint arXiv:2104.04670*,  
894 2021.

895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Table 5: Performance comparison of different methods fine-tuned with Mistral-7B across 4 common sense reasoning tasks and a domain-specific task. UQ4CT shows significant ECE improvements while maintaining high accuracy.

Metrics	Methods	BoolQ	ARC-E	ARC-C	OBQA	ClimateQA
ACC $\uparrow$	LoRA	70.3 <sub>0.62</sub>	84.8 <sub>0.47</sub>	70.2 <sub>0.84</sub>	82.8 <sub>0.62</sub>	72.5 <sub>1.6</sub>
	MC Drop	69.6 <sub>1.07</sub>	84.6 <sub>0.91</sub>	69.6 <sub>0.76</sub>	82.6 <sub>0.71</sub>	72.5 <sub>1.6</sub>
	Ensemble	71.8 <sub>1.29</sub>	84.2 <sub>0.66</sub>	71.0 <sub>1.41</sub>	82.5 <sub>0.6</sub>	72.9 <sub>2.88</sub>
	LA	70.7 <sub>1.82</sub>	82.4 <sub>2.05</sub>	68.5 <sub>3.31</sub>	82.5 <sub>0.77</sub>	71.6 <sub>1.56</sub>
	MixLoRA	73.1 <sub>0.38</sub>	85.5 <sub>1.27</sub>	71.2 <sub>1.75</sub>	83.3 <sub>1.14</sub>	72.0 <sub>1.69</sub>
	UQ4CT	73.6 <sub>0.28</sub>	85.9 <sub>0.82</sub>	74.4 <sub>0.82</sub>	83.7 <sub>1.22</sub>	73.2 <sub>1.29</sub>
ECE $\downarrow$	LoRA	10.17 <sub>0.24</sub>	9.46 <sub>1.62</sub>	18.42 <sub>1.91</sub>	13.3 <sub>0.25</sub>	13.72 <sub>2.62</sub>
	MC Drop	10.62 <sub>0.51</sub>	8.91 <sub>1.35</sub>	18.38 <sub>1.66</sub>	13.3 <sub>0.31</sub>	13.72 <sub>2.61</sub>
	Ensemble	8.72 <sub>1.13</sub>	8.72 <sub>1.49</sub>	17.0 <sub>0.97</sub>	9.14 <sub>2.82</sub>	12.86 <sub>1.78</sub>
	LA	5.33 <sub>2.16</sub>	20.3 <sub>5.7</sub>	21.27 <sub>4.15</sub>	<b>6.41<sub>3.22</sub></b>	14.64 <sub>2.21</sub>
	MixLoRA	8.81 <sub>1.03</sub>	8.16 <sub>0.99</sub>	15.51 <sub>3.86</sub>	10.53 <sub>1.73</sub>	14.05 <sub>3.09</sub>
	UQ4CT	<b>3.07<sub>0.83</sub></b>	<b>5.7<sub>0.69</sub></b>	<b>7.04<sub>0.58</sub></b>	<u>7.92<sub>1.14</sub></u>	<b>11.4<sub>1.14</sub></b>

## A APPENDIX

### A.1 LOAD BALANCING LOSS

We follow the load balancing loss in (Li et al., 2024). Given  $N$  experts indexed by  $i = 1$  to  $N$  and a batch  $B$  with  $T$  tokens, the auxiliary loss is computed as:

$$\mathcal{L}_{aux} = a \cdot N \cdot \sum_{i=1}^N \mathcal{F}_i \cdot \mathcal{P}_i, \quad (13)$$

where

$$\mathcal{F}_i = \frac{1}{T} \sum_{x \in B} \mathbb{1}\{\text{argmax}_k \mathcal{R}(x)_k = i\}, \mathcal{P}_i = \frac{1}{T} \sum_{x \in B} \mathcal{R}(x)_i. \quad (14)$$

Here,  $\mathcal{R}(\cdot)$  is the top-k router,  $\mathcal{F}_i$  is the fraction of tokens dispatched to expert  $i$  and  $\mathcal{P}_i$  is the fraction of the router probability allocated for expert  $i$ . The final loss is multiplied by the expert count  $N$  to keep the loss constant as the number of experts varies, and the constant term  $a$  is set to  $10^{-2}$  as a multiplicative coefficient, which is large enough to ensure load balancing while remaining small enough not to overwhelm the primary objective.

### A.2 EXPERIMENTAL RESULTS WITH MISTRAL-7B

In this section, we present the results using Mistral-7B (Jiang et al., 2023), a different decoder-based LLM backbone. Table 5 shows the results of fine-tuning Mistral-7B on 4 common-sense reasoning tasks and one domain-specific climate question-answering task.

For each of the tasks, UQ4CT effectively calibrates the parameter mixtures, leading to the best ECE performance in 4 out of 5 tasks. This indicates the robustness of UQ4CT across different LLMs.

### A.3 DECIDING NUMBER OF ACTIVE EXPERTS

One important aspect of the LoRA MoE architecture is how many experts to activate. Here, we investigate the performance impact of different number of active LoRA experts. We evaluate the model performance with 1 to 5 active experts with 8 in total.

As shown in Table 6, 2 active experts give the optimal performance in terms of accuracy and ECE scores. One expert alone cannot capture complicated functional relationships, while more than 2 experts could potentially introduce redundant functional bases to the model, which deviates the output distribution more from data distribution, thus worsening predictive and calibration performance. Additionally, more active experts lead to a more flattened distribution across experts, which hardens the alignment of parameter mixtures during fine-tuning.

Table 6: Performance comparison of UQ4CT with varying number of experts on OBQA dataset. Top-2 expert selection strategy grants best accuracy and calibration.

Top-K	ACC $\uparrow$	ECE $\downarrow$
Top-1	74.8 <sub>0.62</sub>	7.69 <sub>1.96</sub>
Top-2	<b>77.3<sub>1.36</sub></b>	<b>5.0<sub>1.15</sub></b>
Top-3	75.2 <sub>0.8</sub>	5.8 <sub>0.81</sub>
Top-4	75.8 <sub>0.53</sub>	7.67 <sub>0.46</sub>
Top-5	75.3 <sub>0.5</sub>	6.29 <sub>0.61</sub>

#### A.4 TRAINING DETAILS

We train our model with total of 8 LoRA experts, and select 2 experts with the highest probability. For each expert, we use  $rank = 16$  and  $alpha = 32$ . We use batch size of 16 to train our model. For climate task, we set the learning rate to  $5e - 4$  and dropout rate to 0.1 to incorporate the small dataset size. For other tasks, we use  $2e - 4$  as our learning rate with dropout 0.05. We use AdamW as our optimizer and a cutoff length of 512 for prompts during training.

The experimental setup for single LoRA based models is similar with LoRA ranks set to 80 to accommodate the MoE model size. For the ensemble baseline, we use an ensemble size of 8 with  $rank = 16$ . For Laplace-LoRA, we follow the Laplace hyperparameters in this Github Repository.

#### A.5 EXPECTED CALIBRATION ERROR

Expected calibration error (ECE) is a commonly used metric to asses uncertainty quantification performance. ECE measures the alignment between prediction accuracy and model confidence through regrouping the predicted probabilities into  $m$  bins. This method then computes the weighted average of the difference between average accuracy and confidence in each bin:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (15)$$

where  $|B_m|$  is the number of evaluated datapoints in bin  $m$ , acc and conf is calculated as following:

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbf{1}(\hat{y}_i = y_i), \quad (16)$$

$$\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} P(\hat{y}_i). \quad (17)$$

#### A.6 MMLU DISTRIBUTION SHIFT DATASET COMPOSITION

- **Computer Science (CS):**
  - College Computer Science
  - Computer Security
  - High School Computer Science
  - Machine Learning
- **Engineering (Eng):**
  - Electrical Engineering
- **Law:**
  - International Law
  - Jurisprudence
  - Professional Law
- **Health:**

1026            – Anatomy  
1027            – Clinical Knowledge  
1028            – College Medicine  
1029            – Human Aging  
1030            – Nutrition  
1031            – Professional Medicine  
1032            – Virology  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079