

How to Alleviate Catastrophic Forgetting in LLMs Finetuning? Hierarchical Layer-Wise and Element-Wise Regularization

Shezheng Song^{*1} Hao Xu^{*1} Jun Ma¹ Shasha Li¹ Long Peng¹ Qian Wan² Xiaodong Liu¹ Jie Yu¹

Abstract

Large Language Models (LLMs) exhibit strong general language capabilities. However, fine-tuning these models on domain-specific tasks often leads to catastrophic forgetting, where the model overwrites or loses essential knowledge acquired during pretraining. This phenomenon significantly limits the broader applicability of LLMs. To address this challenge, we propose a novel approach to compute the element-wise importance of model parameters crucial for preserving general knowledge during fine-tuning. Our method utilizes a dual-objective optimization strategy: (1) regularization loss based on element-wise parameter importance, which constrains the updates to parameters crucial for general knowledge; (2) cross-entropy loss to adapt to domain-specific tasks. Additionally, we introduce layer-wise coefficients to account for the varying contributions of different layers, dynamically balancing the dual-objective optimization. Extensive experiments on scientific, medical, and physical tasks using GPT-J and LLaMA-3 demonstrate that our approach mitigates catastrophic forgetting while enhancing model adaptability. Compared to previous methods, our solution is approximately 20 times faster and requires only 10%–15% of the storage, highlighting the practical efficiency. The code will be released.

1. Introduction

Large Language Models (LLMs) are pretrained on massive and diverse datasets, equipping them with remarkable general capabilities (Wang & Komatsuzaki, 2021; Touvron et al., 2023b; OpenAI, 2024). This pretraining process allows LLMs to serve as versatile tools for a wide range of natural language processing tasks. However, in domains such as medical and scientific fields, LLMs often struggle

^{*}Equal contribution ¹NUDT ²CCNU. Correspondence to: Jun Ma <majun@nudt.edu.cn>.

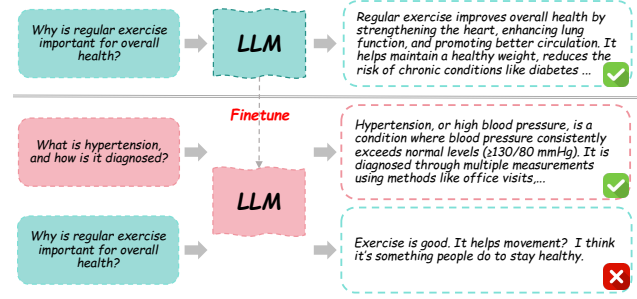


Figure 1: Illustration of catastrophic forgetting: the fine-tuned LLM fails to answer previously known questions.

to perform effectively, necessitating fine-tuning domain-specific data. While fine-tuning could enhance the model task-specific performance, it also introduces a critical challenge: **catastrophic forgetting** (Kirkpatrick et al., 2016; Kemker et al., 2018; Shao & Feng, 2022; Ren et al., 2024). As shown in Figure 1, catastrophic forgetting refers to the phenomenon where a model, during the process of fine-tuning, loses or overwrites knowledge learned during pretraining. This issue poses a severe limitation on the broader applicability of LLMs, as it undermines their versatility and reusability across domains. The fixed data composition and format in the fine-tuning data may impair the general knowledge previously learned by the model. This results in a loss of logical reasoning abilities and related general knowledge, which affects the model performance on domain-specific tasks. On the other hand, it may also lead to a decline in the ability to answer general tasks, including questions it was previously capable of answering.

Addressing catastrophic forgetting is therefore a crucial requirement for maximizing the utility of LLMs. A successful solution needs to achieve a delicate balance: **retaining the essential general knowledge** when learning new domain-specific expertise. This balance is critical when fine-tuning LLMs for specialized tasks, as both domain adaptation and generalizability are necessary for practical applications. EW-CLoRA (Xiang et al., 2024) focuses on the issue of catastrophic forgetting in LLM fine-tuning and uses the Fisher matrix to measure the importance of parameters for general capabilities. However, it requires gradients computed with

labels from the model distribution, necessitating an additional backpropagation pass for online computation. Thus, its computational cost is very high. For GPT-J-6B, calculating the Fisher matrix takes 22 hours on an A800 and requires 23GB of storage, and these requirements increase for larger LLMs. Besides, rsLoRA (Kalajdziewski, 2023) aims to stabilize learning by introducing a rank-stabilized scaling factor, but it does not effectively protect general capabilities as expected.

To address catastrophic forgetting, we calculate parameter importance from two dimensions—element-wise and layer-wise—to constrain the updates of parameters crucial for general capabilities during fine-tuning. Firstly, our approach calculates the path integral during parameter updates as the **element-wise parameter importance for regularization** on the general capabilities of the LLM. This helps preserve parameters critical for general knowledge, minimizing significant modifications to it. Our method could avoid the computation and storage of the Fisher matrix, enabling faster and more storage-efficient computation of parameter importance. Specifically, we define domain ν as the general knowledge, representing the general capabilities of LLMs, and domain μ as the knowledge learned during fine-tuning for specific tasks. Our approach leverages a dual-objective optimization strategy that combines two losses: regularization loss, which reduces updates to parameters critical for domain ν to preserve general knowledge; and cross-entropy (CE) loss, which encourages alignment of domain μ parameters to enhance domain-specific learning. Through the constraint of a dual-objective loss, we aim to maintain general capabilities while performing domain-specific fine-tuning.

Besides, we propose a **layer-wise coefficient** to adjust the weight of regularization loss. In LLMs, different layers contribute differently to generalization ability and domain-specific ability. The impact of each layer on the learning process is not uniform; some layers capture high-level abstract features, while others focus on lower-level details. Traditional approaches often treat the importance of each layer as equal, which overlooks the varying degrees of influence that different layers have on the model learning and generalization ability. Thus, we propose layer-wise coefficients to dynamically adjust the balance between task learning and the retention of general knowledge in each layer, allowing some layers to prioritize task learning, while others preserve general knowledge. We use the L2 norm of the computed element-wise importance of each layer weight to capture their contribution to both objectives.

Through extensive experiments on scientific, physical, and medical tasks using LLMs (GPT-J and LLaMA-3), we demonstrate that our framework achieves state-of-the-art performance, mitigating catastrophic forgetting while enhancing LLM adaptability. To maintain general capabilities,

it is essential to identify and quantify the importance of various parameters that contribute to these capabilities. The computation of parameter importance is typically time-consuming, and storing the associated weights requires substantial memory resources. Our experimental results demonstrate that our method is nearly **20 times faster** and requires only **10%~15% of the storage memory** compared to the previous method, demonstrating the practicality of our approach. Our contributions are as follows:

- We introduce a framework that first records parameter importance on general data, and then applies regularization constraints during fine-tuning on domain-specific data to effectively address catastrophic forgetting in large language models (LLMs).
- We propose the element-wise and layer-wise importance metrics to dynamically adjust parameter updates, preserving critical general knowledge while allowing domain-specific expertise to be learned effectively.
- Our method achieves state-of-the-art performance across multiple datasets using mainstream backbone LLMs. It significantly reduces computational time (20x faster) and storage (10%~15%) for parameter importance estimation compared to prior methods.

2. Related Work

2.1. Continual Learning

Traditionally, continual learning (Wickramasinghe et al., 2024; Hadsell et al., 2020; Wickramasinghe et al., 2024; Vijayan & Sridhar, 2021) refers to developing learning algorithms to accumulate knowledge on non-stationary data. In general, continuous learning could be categorized into the following methods: **Regularization-based methods**. Synaptic Intelligence (SI) (Zenke et al., 2017) dynamically estimates the importance of each parameter in an online fashion, penalizing significant changes to parameters that are important for previously learned tasks during training on new tasks. This method adjusts the learning rate for parameters, ensuring that important parameters are not excessively modified. Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2016) grounded in a Bayesian perspective, estimates the importance of parameters by calculating the Fisher Information Matrix. During new task training, EWC introduces a regularization term that restricts the updates to important parameters, thereby preventing catastrophic forgetting. From a probabilistic viewpoint, EWC derives an importance matrix that quantifies the significance of network parameters for previous tasks. **Architecture-based methods**. Researches (Wu et al., 2024; Wang et al., 2023; 2022; Chen et al., 2022) learn new tasks by adapting the structure of existing models. For instance, Wang et al. inserts trainable task-specific

prompts to the input layer to expand the domain ability. **Replay-based methods.** Researchers (Jin et al., 2022; Liu et al., 2021; Qin et al., 2022; Bai et al., 2022) retain a subset of previously encountered data, which are reintegrated into the training process of the new tasks. **Distillation-based methods.** Researches (Li & Hoiem, 2017; Cao et al., 2021; Shao & Feng, 2022; Gu et al., 2022; Qin & Joty, 2022) learn new tasks under the guidance of a teacher model. For instance, Learning without Forgetting (LwF) (Li & Hoiem, 2017) transfers knowledge from old tasks to new tasks, allowing the model to retain performance on the previous task while learning new ones.

2.2. Catastrophic Forgetting in LLM and LoRA

With the rapid advancement of large language models (LLMs) (Touvron et al., 2023a;b), directly using pretrained models for domain-specific tasks has become prohibitively expensive. As a result, fine-tuning has become the preferred approach, typically divided into full parameter tuning and parameter-efficient fine-tuning (PEFT) methods, such as LoRA (Low-Rank Adaptation) (Hu et al., 2021; Wang et al., 2024). Full parameter fine-tuning (Lv et al., 2023) updates all model parameters to improve task adaptability but often causes catastrophic forgetting. PEFT methods like LoRA, by updating only a small subset of parameters through low-rank matrices, reduce computational costs and mitigate forgetting, though some still occur.

To further reduce catastrophic forgetting, researchers have proposed combining EWC with LoRA in a method known as EWCLoRA (Xiang et al., 2024). This method leverages EWC to calculate the Fisher Information Matrix for parameter importance and uses low-rank matrices of LoRA to limit the scope of parameter updates. However, the calculation of the Fisher matrix introduces significant computational and memory overhead. Additionally, an interpolation-based LoRA (I-LoRA) method is introduced by Ren et al.. I-LoRA constructs a dual-memory experience replay framework, utilizing LoRA parameter interpolation to simulate the weight interpolation process. However, it requires maintaining an additional set of LoRA parameters throughout the process, increasing space cost.

3. Preliminary

LoRA is a lightweight and parameter-efficient fine-tuning method that introduces low-rank decomposition into the weight matrix θ of a pretrained model. Only the newly added low-rank matrices B and A are optimized, while the main weight θ_0 remains frozen. The parameter at time t during fine-tuning can be expressed as $\theta_t = \theta_0 + \delta_t$; $\delta_t = B_t A_t$, where $\theta_0 \in \mathbb{R}^{d \times d}$ are pretrained weights; $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times d}$ are the low-rank matrices with $r \ll d$.

The optimization objective of LoRA is given by:

$$\mathcal{L}_{\text{LoRA}} = \mathcal{L}(y, f(x; \theta(t))) \quad (1)$$

where \mathcal{L} is the task-specific loss function.

Although LoRA achieves parameter efficiency and training effectiveness, it suffers from catastrophic forgetting, where fine-tuning specific tasks hurts the general ability.

4. Hierarchical Importance Regularization

Inspired by Synaptic Intelligence (SI) (Zenke et al., 2017), we propose a framework to constrain LLMs from making significant changes to their general capabilities during fine-tuning, thus addressing catastrophic forgetting in LoRA tuning. As shown in Figure 2, the framework is to compute the importance of each parameter during the training of the initial general task (e.g. ν) and constrain their updates when fine-tuning on subsequent task (e.g. μ). Specifically, the importance scores measure how much each parameter contributes to reducing the loss in the ν task, and these scores are used to guide the fine-tuning process for the new μ task. This ensures that the critical parameters for ν task are modified to a lesser extent when learning μ task.

4.1. General Element-Wise Importance Recording

In the general task ν , LoRA fine-tuning is performed by minimizing the task-specific loss $\mathcal{L}_{\text{task}}$. The training process is characterized by a trajectory $\theta(t)$ in parameter space. The task-specific loss $\mathcal{L}_{\text{task}}$ is generally computed using cross-entropy loss.

$$L_\nu = \mathcal{L}_{\text{task}}^\nu(y_\nu, f(x_\nu; \theta(t))) = - \sum_{i=1}^N y_k \log(p_k) \quad (2)$$

where y_k is the true label (target) for the i -th example, p_k is the predicted probability of the model for the correct label, and N is the total number of examples.

We define the contribution of parameter i to the reduction of the loss function as ω_i . The larger the value of ω_i , the more important the parameter i is for maintaining the performance of the task ν . The change in the loss function from time t_0 to time t_1 can be defined as the sum of the contributions of all parameters:

$$L(\theta_{t_1}) - L(\theta_{t_0}) = - \sum_i w_i \quad (3)$$

In accordance with the typical behavior of the loss value, which generally decreases, we introduced a negative sign on the right-hand side of Equation (3) to ensure that the value of ω_i remains positive.

During the training process of task ν , the total change in the loss function can be obtained by performing a path

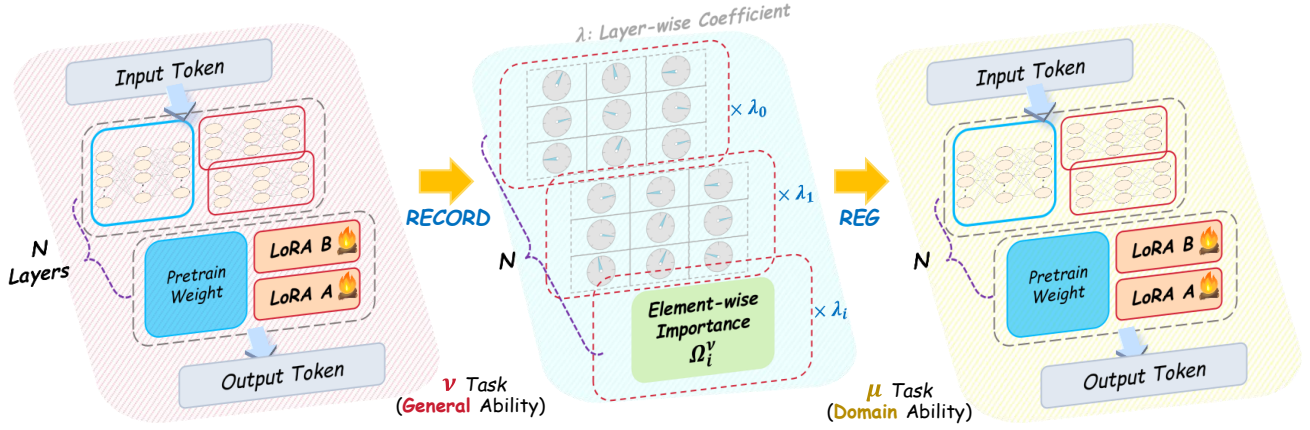


Figure 2: Adaptive constraint combining element-wise and layer-wise importance to preserve general capabilities from the ν task while learning domain-specific abilities for the μ task. RECORD means the general importance recording in Section 4.1. REG means the regularization in Section 4.2 and Section 4.3.

integral of the gradient of the loss function with respect to the parameters, that is, the path integral from the initial parameter value θ_{t_0} to the final parameter value θ_{t_1} :

$$L(\theta_{t_1}) - L(\theta_{t_0}) = \int_{\theta_{t_0}}^{\theta_{t_1}} g(\theta(t)) d\theta(t) \quad (4)$$

where g represents the gradient of the loss function with respect to the parameters. By expanding $d\theta(t)$ in Equation (4), we can derive the following expression:

$$\begin{aligned} L(\theta_{t_1}) - L(\theta_{t_0}) &= \int_{t_0}^{t_1} g(\theta(t)) \theta'(t) dt \\ &= \sum_i \int_{t_0}^{t_1} g(\theta_i(t)) \theta'_i(t) dt \end{aligned} \quad (5)$$

In accordance with Equation (3) and Equation (5), it is deduced that the defined quantity ω_i corresponds precisely to the negative of the path integral of the gradient g_i .

$$w_i = - \int_{t_0}^{t_1} g(\theta_i(t)) \theta'_i(t) dt \quad (6)$$

This indicates that we can represent ω_i using the product of $g(\theta_i(t)) = \frac{\partial L}{\partial \theta_i}$ and $\theta'_i(t) = \frac{\partial \theta_i}{\partial t}$ (Zenke et al., 2017).

Considering that LoRA utilizes low-rank matrix approximation for fine-tuning, the parameter updates and gradients need to be adjusted accordingly.

The parameters updating process of low-rank matrices B and A at time $t + 1$ are defined as:

$$\begin{aligned} B(t+1) &= B(t) - \eta g^B(t) \\ A(t+1) &= A(t) - \eta g^A(t) \end{aligned} \quad (7)$$

where η is the learning rate, $g^A(t)$ and $g^B(t)$ are the gradients of the loss functions with respect to A and B. Based on Equation (7), we derive the following expression:

$$\begin{aligned} B(t+1)A(t+1) &= B(t)A(t) - \eta g^B(t)A(t) \\ &\quad - \eta B(t)g^A(t) + \eta^2 g^B(t)g^A(t) \end{aligned} \quad (8)$$

According to the definition of LoRA, the parameters at time $t + 1$ and time t are respectively defined as:

$$\begin{aligned} \theta(t+1) &= \theta_0 + B(t+1)A(t+1) \\ \theta(t) &= \theta_0 + B(t)A(t) \end{aligned} \quad (9)$$

Based on Equation (8), we derive the change of the parameters, which is expressed in terms of $g^A(t)$ and $g^B(t)$:

$$\begin{aligned} \Delta\theta &= \theta(t+1) - \theta(t) \\ &= B(t+1)A(t+1) - B(t)A(t) \\ &= -\eta(g^B(t)A(t) + B(t)g^A(t) - \eta g^B(t)g^A(t)) \end{aligned} \quad (10)$$

According to the definition of batch gradient descent, the change in parameters is the negative product of the gradient and the learning rate. If we regard LoRA as a special form of full fine-tuning, we can assume that there exists a gradient $\tilde{g}(t)$ that completes the parameter update process (Wang et al., 2024).

Based on Equation (10) and the definition of $\tilde{g}(t)$, we obtain the parameter change and hypothetical gradient at time t .

$$\begin{aligned} \tilde{\theta}'(t) &= B(t+1)A(t+1) - B(t)A(t) \\ \tilde{g}(t) &= g^B(t)A(t) + B(t)g^A(t) - \eta g^B(t)g^A(t) \end{aligned} \quad (11)$$

In this way, we obtain the value of ω_i for the LoRA scenario.

$$w_i = - \int_{t_0}^{t_1} \tilde{g}_i(t) \tilde{\theta}'_i(t) dt \quad (12)$$

To quantify the importance of each parameter, we calculate an importance score Ω_i^ν based on its contribution to the change in loss during training of task ν . Specifically, the importance of a parameter is computed as:

$$\Omega_i^\nu = \sum_{\nu} \frac{\omega_i^\nu}{(\Delta_i^\nu)^2 + \xi} \quad (13)$$

where $\Delta_i^\nu = \theta_i(t^\nu) - \theta_i(t^0)$ is whole change of the i -th parameter θ_i during task ν , $\theta_i(t^\nu)$ is the final parameter after fine-tuning on task ν . In the context of LoRA fine-tuning, the Δ_i^ν is defined as $(B(t^\nu)A(t^\nu))_i$. This relationship stems from the fact that, at the initialization of LoRA at time 0, the B matrix is set to zero. The term in the denominator $(\Delta_i^\nu)^2$ ensures that the regularization term carries the same units as the loss L . ξ is a small positive constant to prevent division by zero. This formulation assigns higher scores to parameters that have a significant impact on loss reduction while accounting for their magnitude to avoid bias toward large updates.

4.2. Element-Wise Regularization in Domain Tuning

After fine-tuning the ν task, we extend the optimization objective to include both task-specific and regularization losses during μ finetuning. The task-specific loss $\mathcal{L}_{\text{task}}^\mu$ drives the adaptation to the μ task. To preserve knowledge from the ν task, the regularization loss penalizes deviations from the important parameter values recorded in the ν task. The regularization loss $\mathcal{L}_{\text{reg},l}^\nu$ of the l -th layer is defined as:

$$\mathcal{L}_{\text{reg},l}^\nu = \sum_i^n \sum_{\nu < t < \mu} \Omega_i^\nu (\theta_i^t - \theta_i^\nu)^2 \quad (14)$$

where Ω_i^ν represents the importance of the i -th parameter in the ν task, and θ_i^ν is the reference parameter after ν task fine-tuning. This loss ensures that parameters with high importance scores remain close to their ν task values while allowing less important parameters more flexibility for adaptation. During training, ω_i values are updated continuously, while the cumulative importance Ω_i^ν is updated only at the end of task ν . After updating Ω_i^ν , the ω_i is reset to zero.

4.3. Layer-Wise Coefficient Regularization

We compute the importance of each layer based on its contribution to the parameters learned in the ν task. This layer-specific importance metric allows the model to dynamically adjust the regularization across different layers. The layer-wise weighted regularization is defined as :

$$\mathcal{L}_{\text{reg}}^\nu = \sum_l \text{softmax}(\|\Omega_l^\nu\|_2) \mathcal{L}_{\text{reg},l}^\nu \quad (15)$$

where $\|\Omega_l^\nu\|_2$ denotes the L2 norm of the parameter importance matrix Ω_l^ν for the l -th layer, which reflects the

significance of the parameters learned in the ν task. The total loss for the μ task is defined as:

$$\mathcal{L}^\mu = \mathcal{L}_{\text{task}}^\mu + \varphi \mathcal{L}_{\text{reg}}^\nu \quad (16)$$

The use of this adaptive regularization $\mathcal{L}_{\text{reg}}^\nu$ helps mitigate catastrophic forgetting by maintaining the integrity of essential features learned in prior tasks. φ is the hyperparameter that controls the weight of the domain ($\mathcal{L}_{\text{task}}$) and general (\mathcal{L}_{reg}) ability of LLM.

5. Experiments

5.1. Backbone LLMs and Baseline Methods

Following the previous work (Xiang et al., 2024), two mainstream LLMs are used for the evaluation of our method: (1) *GPT-J* (Wang & Komatsuzaki, 2021) is a GPT-2-like causal language model trained on the Pile dataset. It is suitable for various understanding and generation tasks. (2) *LLaMA-3* (Dubey et al., 2024) is the third-generation open-source LLM. It is designed with enhanced efficiency and scalability, offering state-of-the-art performance across various benchmarks. These models vary in architecture and parameter count, enabling a robust evaluation of our method.

We compare our method with the following approaches: (1) *Base*: the model without any tuning. (2) *LoRA*(μ) (Hu et al., 2021): the method is fine-tuned using only data from the μ task (domain-specific task). (3) *LoRA*($\nu + \mu$): the method is first fine-tuned using data from the ν task (general task), and then fine-tuned using data from the μ task (domain-specific task). (4) *EWCLoRA* (Xiang et al., 2024): a method using the EWC method, where the Fisher matrix is computed and regularization constraints are applied to preserve the important parameters while updating for the new task. (5) *rsLoRA*: an enhanced LoRA method that modifies the scaling factor to prevent gradient collapse, enabling better fine-tuning performance with higher-rank adapters while maintaining the same inference cost.

5.2. Tasks, Metrics, and Hyperparameters

ν Task (General Ability): The ν task focuses on learning which parameters are important for general tasks. Following previous work (Xiang et al., 2024), we take Pile (Gao et al., 2020) as the evaluation datasets for LLM general ability. LoRA is applied to fine-tune the model on the ν task, and parameter importance for Synaptic Intelligence (SI) is recorded during this stage.

μ Task (Domain Ability): The μ task evaluates the ability to adapt to specific tasks while mitigating catastrophic forgetting of general knowledge. We select three representative tasks: (1) *Medical task*: MedMCQA dataset (Pal et al., 2022). (2) *Scientific task*: SciQ dataset (Welbl et al., 2017). (3) *Physics task*: PiQA dataset (Bisk et al., 2020).

Table 1: General and domain ability of LLMs. (Acc \uparrow : Accuracy of domain ability, PPL \downarrow : Perplexity of general ability.)

	LLaMA-3						GPT-J					
	SciQ		PiQA		MedMCQA		SciQ		PiQA		MedMCQA	
	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow
Base	4.94	95.10	4.94	48.53	4.94	18.50	3.28	91.60	3.28	49.13	3.28	21.30
LoRA(μ)	5.05	96.20	5.43	48.75	5.04	53.69	3.43	96.50	3.54	50.16	3.49	38.35
LoRA($\nu + \mu$)	5.31	96.10	5.58	46.91	5.15	53.12	3.39	96.20	3.52	49.51	3.37	33.66
rsLoRA	5.28	96.50	5.71	47.50	5.24	51.92	3.50	96.20	3.65	49.62	3.35	35.69
EWC-L	4.88	96.30	4.98	48.45	4.79	56.39	3.38	96.10	3.47	49.40	3.38	36.48
Ours	4.64	97.10	4.90	51.14	4.64	55.80	3.35	96.80	3.40	50.49	3.34	36.10

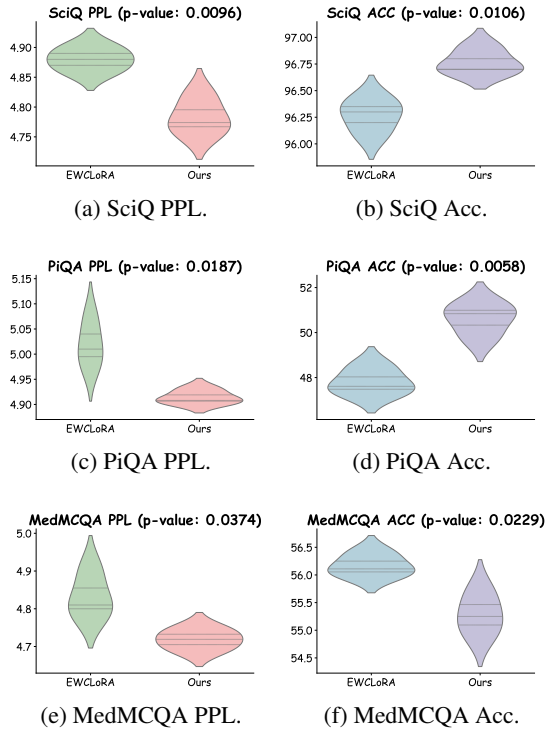


Figure 3: Independent samples t-test of EWCLoRA and our method on LLaMA-3: violin plots of perplexity (PPL) and accuracy (Acc) across datasets

The LLMs selected for our experiments are GPT-J-6B and LLaMA 3.2-3B. The batch size is set to 20, and the learning rate is set to $8e-4$. The rank for LoRA fine-tuning is set to 8, with the LoRA alpha value set to 32. Both the ν and μ tasks are trained for 5 epochs.

6. Results and Analysis

6.1. Comparison of General and Domain Capabilities

As shown in Table 1, our method achieves better preservation of general ability (as reflected by the lowest PPL) while maintaining domain-specific accuracy comparable to, or even better than, previous methods. This demonstrates that our approach effectively balances domain accuracy and general perplexity.

Figure 3 presents a comparison between the results of EWCLoRA and our method through independent samples t-tests. The six subplots show the Perplexity (PPL) and Accuracy (Acc) across SciQ, PiQA, and MedMCQA datasets. The p-values for perplexity on SciQ, PiQA, and MedMCQA, and for accuracy on SciQ and PiQA are below 0.05, indicating statistically significant differences and demonstrating the superiority of our method over EWCLoRA.

Figure 4 shows the loss curves in the learning process of GPT-J and LLaMA-3 across three datasets. The total loss is the weighted sum of the task loss $\mathcal{L}_{\text{task}}$ and general loss \mathcal{L}_{reg} . As observed, the task loss continuously decreases, while the \mathcal{L}_{reg} exhibits an initial increase followed by a decrease. As defined in Equation (14), \mathcal{L}_{reg} measures the difference between the model parameters θ_ν after learning on task ν and the model parameters θ_μ learned on the current task μ . Initially, when learning on task μ , the model parameters are not yet updated, so the general loss is zero. As the task loss updates the parameters, the model starts to deviate from θ_ν , causing the general loss to rise. This mechanism enforces the model to learn in a way that minimizes both general and task losses simultaneously.

6.2. Complexity Comparison

We compare our HLoRA method with the previous SOTA method, EWCLoRA, from two aspects: the time required for importance calculation and the storage memory needed. As shown in Figure 5, our method is nearly **20 times faster** and

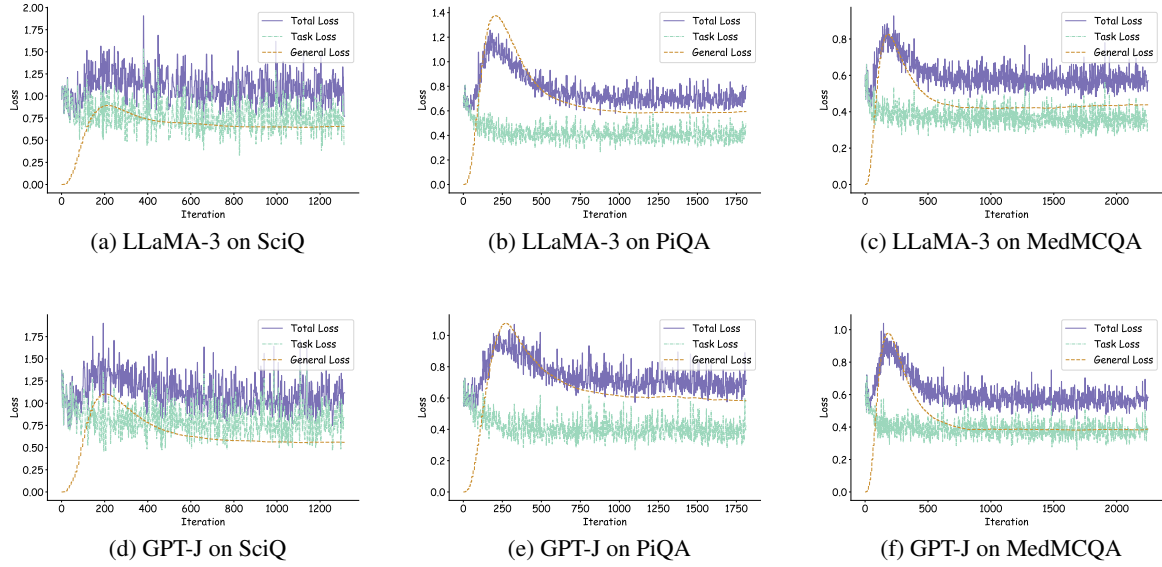


Figure 4: Loss curves on three datasets: balancing task learning and generalization. The total loss consists of task loss ($\mathcal{L}_{\text{task}}$) and a scaled version of general loss (\mathcal{L}_{reg}), where task loss controls the model learning on new domain data, and general loss helps maintain the model generalization ability.

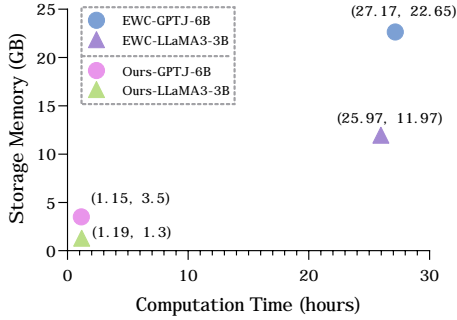


Figure 5: Comparison of computation time and storage for importance Ω_l^V between previous method and ours.

requires only **10%~15% of the storage memory** compared to EWCLoRA, demonstrating the practicality of ours.

Time Complexity: The experiments were conducted on an A800 GPU to evaluate the time complexity of our method in comparison with EWCLoRA. For EWCLoRA, the Fisher matrix computation followed the approach described in the original paper, using 20,000 randomly sampled data points from the Pile dataset with a maximum batch size of 8. In contrast, for our method, the time measurement was based on 5 training epochs, a setting determined through empirical evaluation to achieve optimal performance. The experimental results show that for GPT-J-6B and LLaMA-3-3B, EWCLoRA requires **27.17** and **25.97** hours, respectively, to

compute the importance matrix, while our HLoRA method only takes **1.15** and **1.19** hours.

Storage Memory: EWCLoRA necessitates the computation and storage of the Fisher matrix based on the Pile dataset before calculating the parameter importance. According to the original paper, the Fisher matrix for GPT-J-6B occupies approximately **22.65** GB of memory. Similarly, for LLaMA-3-3B, the Fisher matrix occupies **11.97** GB of memory, calculated based on the Fisher computation method described in the original work. In contrast, the storage memory of our method is only **3.5** GB and **1.3** GB, offering a significant advantage in terms of memory efficiency. This demonstrates that EWCLoRA incurs substantial storage overhead, whereas our method avoids such requirements, providing a more space-efficient solution.

6.3. Regularization Coefficient Analysis

Figure 6 demonstrate the effect of the regularization coefficient φ in Equation (16) on PPL and accuracy across three tasks. As φ increases, PPL gradually decreases, indicating a stronger emphasis on preserving general ability. Higher values of φ correspond to better general ability retention. However, as shown in Figure 6b, increasing φ negatively impacts the average accuracy on PiQA. Thus, e^{-3} is selected as the optimal value for the regularization coefficient to balance task performance and general ability (lower PPL).

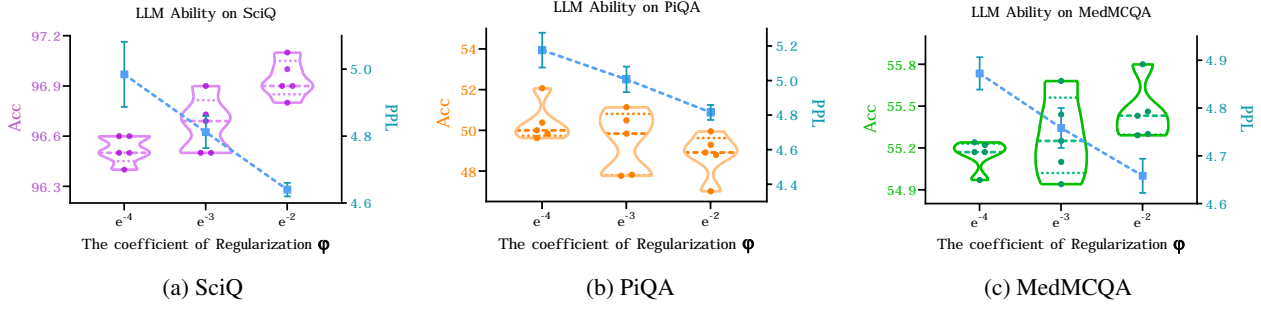


Figure 6: The influence of regularization coefficient ϕ on LLaMA-3 across datasets. (Acc \uparrow : Accuracy, PPL \downarrow : Perplexity.)

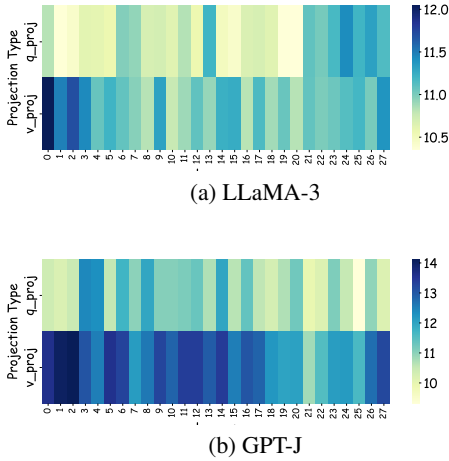


Figure 7: Log-scaled heatmap of L2 norms of parameter importance Ω_l^ν for q_proj and v_proj after LoRA fine-tuning on ν task across layers.

6.4. Parameters Importance Visualization

Figure 7 highlights the importance in Equation (13) of q_proj and v_proj layers for general capabilities during the LoRA fine-tuning process. The heatmap illustrates that the v_proj layers, particularly in the first four and the last layer, are crucial for preserving general knowledge. In contrast, the importance of the q_proj layers is relatively weaker across the model. The L2 norms have been log-transformed to facilitate the comparison of the relative significance of these parameters across layers.

6.5. Ablation Study

As shown in Table 2, to investigate the role of different components in our proposed HLoRA, we conduct ablation studies by selectively removing certain structures and observing the resulting impact. Specifically, we exclude two sets of components: (1) *layer*: eliminating the differentiation of importance among layers, and (2) *layer, element*:

Table 2: Ablation experiments. (layer: layer-wise weighted regularization, element: element-wise regularization.)

	SciQ		PiQA		MedMCQA	
	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow	PPL \downarrow	Acc \uparrow
LLaMA-3						
Ours	4.64	97.10	4.90	51.14	4.64	55.80
- layer	4.75	96.80	4.96	49.70	4.74	54.41
- layer, element	5.31	96.10	5.58	46.91	5.15	53.12
GPT-J						
Ours	3.35	96.80	3.40	50.49	3.34	36.10
- layer	3.36	96.30	3.41	49.95	3.35	35.62
- layer, element	3.39	96.20	3.52	49.51	3.37	33.66

removing both layer-wise and element-wise importance, i.e., training the ν task first and then training the μ task without imposing any regularization constraints throughout the process. Upon removing the two components, the performance of methods based on two backbone LLMs declines across three datasets, thereby highlighting the effectiveness of the layer-wise and element-wise importance introduced.

7. Conclusion

This paper addresses the critical issue of catastrophic forgetting in large language models (LLMs) during domain-specific fine-tuning. We propose a novel fine-tuning framework that preserves general capabilities while enabling efficient adaptation to new domains, minimizing knowledge loss in tasks outside the fine-tuned domain. Additionally, we introduce a layer-wise coefficient to adjust the balance between regularization loss and cross-entropy loss dynamically. This adjustment accounts for the varying contributions of different layers to both generalization and domain-specific learning. Extensive experiments in scientific, physical, and medical tasks show that our framework effectively mitigates catastrophic forgetting while maintaining performance in domain-specific tasks.

Impact Statements

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Bai, G., He, S., Liu, K., and Zhao, J. Incremental intent detection for medical domain with contrast replay networks. In *ACL (Findings)*, pp. 3549–3556. Association for Computational Linguistics, 2022.
- Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Cao, Y., Wei, H.-R., Chen, B., and Wan, X. Continual learning for neural machine translation. *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*, pp. 3964 – 3974, 2021.
- Chen, C., Yin, Y., Shang, L., Jiang, X., Qin, Y., Wang, F., Wang, Z., Chen, X., Liu, Z., and Liu, Q. bert2bert: Towards reusable pretrained language models. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 1:2134 – 2148, 2022.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Gu, S., Hu, B., and Feng, Y. Continual learning of neural machine translation within low forgetting risk regions. *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, pp. 1707 – 1718, 2022.
- Hadsell, R., Rao, D., Rusu, A. A., and Pascanu, R. Embracing change: Continual learning in deep neural networks. *Trends in Cognitive Sciences*, 24(12):1028 – 1040, 2020.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Jin, X., Zhang, D., Zhu, H., Xiao, W., Li, S.-W., Wei, X., Arnold, A., and Ren, X. Lifelong pretraining: Continually adapting language models to emerging corpora. *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*, pp. 4764 – 4780, 2022.
- Kalajdziewski, D. A rank stabilization scaling factor for fine-tuning with lora. *arXiv preprint arXiv:2312.03732*, 2023.
- Kemker, R., McClure, M., Abitino, A., Hayes, T., and Kanan, C. Measuring catastrophic forgetting in neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N. C., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114:3521 – 3526, 2016.
- Li, Z. and Hoiem, D. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Liu, Q., Cao, P., Liu, C., Chen, J., Cai, X., Yang, F., He, S., Liu, K., and Zhao, J. Domain-lifelong learning for dialogue state tracking via knowledge preservation networks. *Conference on Empirical Methods in Natural Language Processing, Proceedings*, pp. 2301 – 2311, 2021.
- Lv, K., Yang, Y., Liu, T., Gao, Q., Guo, Q., and Qiu, X. Full parameter fine-tuning for large language models with limited resources. *arXiv preprint arXiv:2306.09782*, 2023.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2024.
- Pal, A., Umapathi, L. K., and Sankarasubbu, M. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pp. 248–260. PMLR, 2022.
- Qin, C. and Joty, S. Lfpt5: A unified framework for lifelong few-shot language learning based on prompt tuning of t5. *ICLR 2022 - 10th International Conference on Learning Representations*, 2022.
- Qin, Y., Zhang, J., Lin, Y., Liu, Z., Li, P., Sun, M., and Zhou, J. Elle: Efficient lifelong pre-training for emerging data. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pp. 2789 – 2810, 2022.
- Ren, W., Li, X., Wang, L., Zhao, T., and Qin, W. Analyzing and reducing catastrophic forgetting in parameter efficient tuning. *arXiv preprint arXiv:2402.18865*, 2024.

-
- Shao, C. and Feng, Y. Overcoming catastrophic forgetting beyond continual learning: Balanced training for neural machine translation. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 1:2023 – 2036, 2022.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Vijayan, M. and Sridhar, S. Continual learning for classification problems: A survey. *IFIP Advances in Information and Communication Technology*, 611 IFIPAICT:156 – 166, 2021.
- Wang, B. and Komatsuzaki, A. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- Wang, P., Panda, R., and Wang, Z. Data efficient neural scaling law via model reusing. *Proceedings of Machine Learning Research*, 202:36193 – 36204, 2023.
- Wang, Z., Zhang, Z., Lee, C.-Y., Zhang, H., Sun, R., Ren, X., Su, G., Perot, V., Dy, J., and Pfister, T. Learning to prompt for continual learning. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2022-June:139 – 149, 2022.
- Wang, Z., Liang, J., He, R., Wang, Z., and Tan, T. Lora-pro: Are low-rank adapters properly optimized? *arXiv preprint arXiv:2407.18242*, 2024.
- Welbl, J., Liu, N. F., and Gardner, M. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*, 2017.
- Wickramasinghe, B., Saha, G., and Roy, K. Continual learning: A review of techniques, challenges, and future directions. *IEEE Transactions on Artificial Intelligence*, 5(6):2526 – 2546, 2024.
- Wu, C., Gan, Y., Ge, Y., Lu, Z., Wang, J., Feng, Y., Shan, Y., and Luo, P. Llama pro: Progressive llama with block expansion. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 1:6518 – 6537, 2024.
- Xiang, J., Tao, T., Gu, Y., Shu, T., Wang, Z., Yang, Z., and Hu, Z. Language models meet world models: Embodied experiences enhance language models. *Advances in neural information processing systems*, 36, 2024.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In *International conference on machine learning*, pp. 3987–3995. PMLR, 2017.