

# GR-LoRA: Gradient-Recycling Low-Rank Adaptation for Class-Incremental Learning

Yipeng Lin<sup>\*1</sup> Fengqiang Wan<sup>\*1</sup> Yang Yang<sup>1</sup>

## Abstract

Pre-trained models with parameter-efficient fine-tuning have shown strong effectiveness in Class-Incremental Learning (CIL), which seeks to balance model plasticity and stability. In this context, orthogonality constraints can significantly enhance model stability, yet their reliance on subspace inevitably compromises model plasticity over long tasks. To address this, we propose Gradient-Recycling Low-Rank Adaptation (GR-LoRA), which reconciles stability and plasticity by recycling the gradients discarded in orthogonal projection. Specifically, GR-LoRA recycles post-decomposition non-orthogonal gradient components into task-specific lightweight modules and selects optimal module via entropy to improve plasticity, while incorporating local and global mismatch suppression to preserve stability by synthesizing out-of-distribution representations across all tasks. Theoretical analysis confirms that this recycling strategy preserves stability and improves plasticity. Experimental results from multiple CIL benchmarks verify the effectiveness and general applicability of GR-LoRA.

## 1. Introduction

Class-Incremental Learning (CIL) aims to enable models to progressively acquire knowledge from sequential data streams in real-world scenarios (De Lange et al., 2021; Parisi et al., 2019; Zhou et al., 2022). The primary challenge in CIL is stability-plasticity dilemma (Mermillod et al., 2013; CHEN et al., 2023), where stability characterizes the preservation of prior knowledge, and plasticity reflects the capacity to adapt to new concepts. Large-scale pre-trained models (PTMs) (Steiner et al., 2022; Radford et al., 2021),

<sup>\*</sup>Equal contribution <sup>1</sup>Nanjing University of Science and Technology, Nanjing, China. Correspondence to: Yang Yang <yyang@njust.edu.cn>.

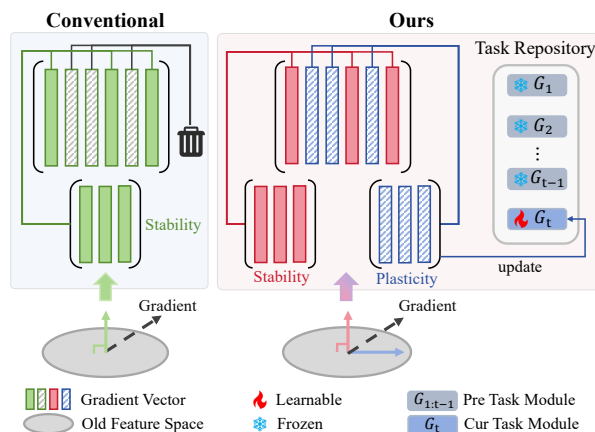


Figure 1. Left: Projection-based method enforces orthogonality to past task subspaces, discarding non-orthogonal gradient components. Right: Our method preserves the orthogonal update in the shared model, while routing the non-orthogonal component into a task repository implemented as task-specific modules.

together with parameter-efficient fine-tuning (PEFT) methods (Houlsby et al., 2019; Han et al., 2024) facilitate this dilemma by leveraging strong representations while updating only a minimal set of model parameters.

However, existing PTM-based CIL methods still grapple with distinct limitations. In the initial exploration, selection-based methods (Wang et al., 2022b; Smith et al., 2023; Wang et al., 2022a) dynamically select task-specific modules from a shared pool using input similarity, which can lead to erroneous selection due to overconfident matching (Pan, 2019). To circumvent unreliable selection, prototype-rectification methods (Zhou et al., 2024; Wu et al., 2025) adopt an expandable parameter learning strategy on the same model, update historical prototypes using current-task samples and rely on the rectified prototypes at inference time, which inevitably accumulates errors over long tasks. To avoid accumulated errors of prototype rectification, orthogonality-based methods (Liang & Li, 2024; Liu & Chang, 2025) project current-task gradients onto subspaces orthogonal to past tasks to prevent interference, but explicitly discarding the non-orthogonal components, as illustrated in Figure 1 (left). Consequently, as tasks accumulate, this constraint progressively contracts the feasible optimization space (Liang

& Li, 2024). As established in Theorem 3.1, this contraction results in reduced plasticity due to the elimination of non-orthogonal gradient components during within-task optimization. This observation motivates an investigation into whether the non-orthogonal gradient components can be explicitly repurposed to balance stability and plasticity.

Motivated by this, we propose Gradient Recycling Low-Rank Adaptation (GR-LoRA), which redirects the non-orthogonal gradient components into task-specific module, as illustrated in Figure 1(right). Specifically, adhering to orthogonality principles, GR-LoRA decomposes task gradients into orthogonal and non-orthogonal components, encodes the latter into task-specific LoRA (Hu et al., 2022) modules and selects optimal module via entropy to improve plasticity without compromising stability. Theorem 3.2 demonstrates that preserving these components in auxiliary parameter suffices to recover the original optimization trajectory. To further improve the reliability of this selection for preserving stability, we incorporate two suppression strategies, consisting of Local Mismatch Suppression (LMS) and Global Mismatch Suppression (GMS). LMS suppresses spurious activations by routing current-task samples through mismatched LoRA modules, forcing alignment with current categories. Subsequently, GMS leverages prototypes to synthesize out-of-distribution representations via inter-class similarity compensation. These synthetic features approximate the global distribution, enabling the classification head to effectively suppress mismatched signals across the entire tasks. Extensive experiments across diverse benchmark datasets demonstrate that GR-LoRA outperforms existing state-of-the-art (SOTA) approaches.

## 2. Related Work

### 2.1. PTM-Based Class-Incremental Learning

The advent of large-scale PTM (Steiner et al., 2022; Radford et al., 2021) has provided robust and generalizable representations for downstream tasks. However, fully fine-tuning these models is computationally prohibitive and prone to catastrophic forgetting (French, 1999; French & Ferrara, 2020). Consequently, PEFT (Han et al., 2024; Houlsby et al., 2019) has established itself as the mainstream approach for CIL. By freezing the pre-trained backbone and optimizing only a minimal set of parameters, PEFT effectively preserves historical knowledge while accommodating new tasks. Based on the structural integration of trainable parameters, existing approaches generally fall into three categories: prompt-based, adapter-based, and LoRA-based methods. Prompt-based methods insert learnable tokens to encode knowledge via key-query retrieval mechanisms, such as L2P (Wang et al., 2022b) selects dynamic prompts from a pool to instruct the frozen backbone; DualPrompt (Wang et al., 2022a) utilizes general and expert prompts to separate

task-invariant knowledge from task-specific ones; CODA-Prompt (Smith et al., 2023) employs a decomposed attention scheme to assemble prompts as weighted linear combinations of pool components. Adapter-based methods, such as EASE (Zhou et al., 2024) mitigates feature degradation by integrating multiple adapter predictions with semantic-guided prototype synthesis; SSIAT (Tan et al., 2024) aligns new and old task features by continuously tuning shared adapters and estimating mean shifts; and MOS (Sun et al., 2025) optimizes inference efficiency by merging adapter parameters via a self-optimization retrieval mechanism. LoRA-based methods, such as InfLoRA (Liang & Li, 2024) minimizes inter-task interference by imposing orthogonal constraints to isolate low-rank subspaces; CL-LoRA (He et al., 2025) splits the adaptation process by dedicating initial blocks to shared feature learning via early-exit distillation, while employing gradient reassignment on deeper blocks to preserve task-specific knowledge; LoRA-DRS (Liu & Chang, 2025) maintains representation stability by explicitly subtracting task-specific variations to construct a drift-resistant space; and MACIL (Wu et al., 2025) mitigates semantic drift by incorporating mean shift compensation and covariance calibration to align class representations across tasks. Despite these advances, existing PTM-based CIL methods still struggle with stability-plasticity dilemma. We propose GR-LoRA, a novel orthogonality-based low-rank adaptation method that simultaneously ensuring stability and enhancing plasticity through adaptive recycle gradient.

### 2.2. Out-of-Distribution Detection.

Out-of-Distribution (OOD) detection (Sun et al., 2022; Yang & Xu, 2025) is critical for ensuring model reliability by identifying samples deviating from the training distribution. Existing approaches generally fall into two paradigms: post-hoc scoring and training-time regularization. Post-hoc methods derive uncertainty scores from pre-trained classifiers. Baseline metrics, such as entropy (Liu et al., 2020), operate on the assumption that models exhibit higher confidence for In-Distribution (ID) samples than OOD samples. In modular CIL, where distinct task-specific models are maintained, identifying the correct module is intrinsically an OOD problem (Aljundi et al., 2017). Specifically, TUNA (Wang et al., 2025) has demonstrated that entropy minimization can effectively select the appropriate task-specific expert during inference. Conversely, training-time methods modify the optimization objective to enforce compact decision boundaries. A seminal approach is Outlier Exposure (Hendrycks et al., 2019; Du et al., 2022), which leverages large-scale auxiliary datasets to explicitly penalize the model’s confidence on anomalies. Adapting this to CIL, methods such as TPL (Lin et al., 2024) leverage stored data from past tasks as surrogate OOD classes, thereby reinforcing the boundary discrimination between current and historical task.

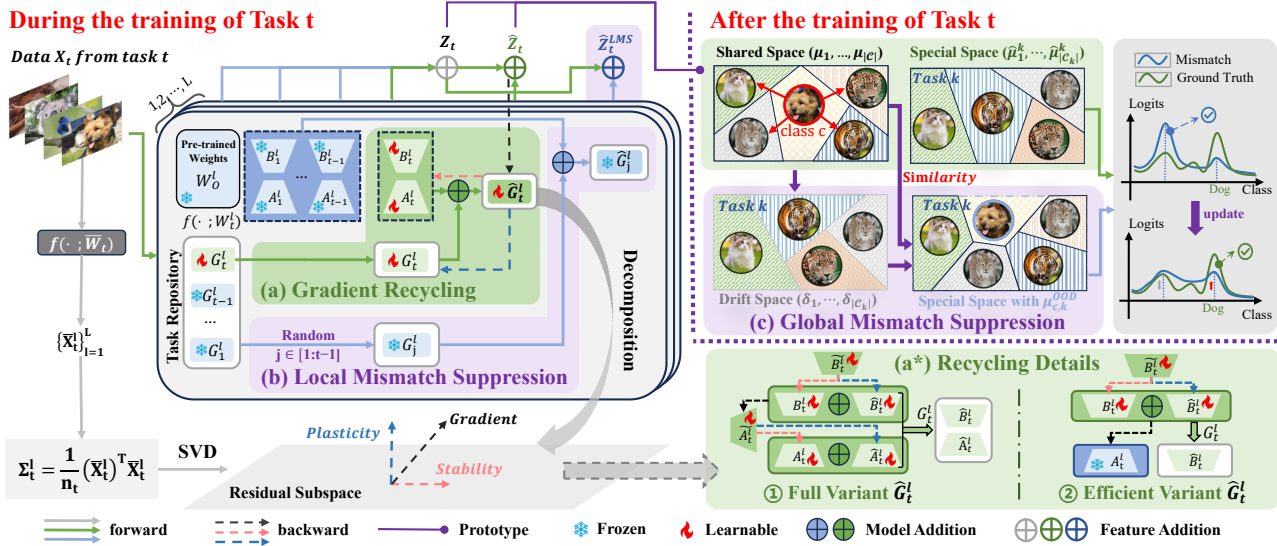


Figure 2. Illustration of the proposed framework. Before training on task  $t$ , a residual subspace is constructed from the current task data. During training of task  $t$ , the pre-trained backbone and all historical LoRA modules are frozen. We apply (a) **Gradient Recycling**: the learnable modules  $G_t$ ,  $A_t$ , and  $B_t$  are aggregated to form  $\hat{G}_t$  in the forward pass, while the backward gradients are projected onto the residual subspace and decomposed, as detailed in (a\*) **Recycling Details**. Simultaneously, (b) **Local Mismatch Suppression** is employed by feeding the current data through historical modules. After training on task  $t$ , (c) **Global Mismatch Suppression** leverages statistical class prototypes to synthesize out-of-distribution features together with in-distribution features for classifier retraining.

### 3. Method

#### 3.1. Preliminaries

The CIL problem is formulated as a stream of  $T$  tasks with disjoint label spaces. Let  $\mathcal{D}_t = \{(\mathbf{x}_{t,i}, y_{t,i})\}_{i=1}^{N_t}$ , denote the dataset for the  $t$ -th task, where  $\mathbf{x}_{t,i} \in \mathcal{X}_t$  and  $y_{t,i} \in \mathcal{Y}_t$ , and the label spaces satisfy  $\mathcal{Y}_t \cap \mathcal{Y}_{t'} = \emptyset$  for all  $t \neq t'$ . Access to historical data  $\mathcal{D}_{1:t-1}$  is strictly prohibited.

**During the training of task  $t$ .** For each task  $t$ , the pre-trained backbone  $W_0$  and all accumulated task parameters  $\{B_j A_j\}_{j=1}^{t-1}$  are frozen. A new trainable low-rank branch  $\Delta W_t = B_t A_t$  is introduced, with  $A_t \in \mathbb{R}^{r \times d}$  initialized to zero and  $B_t \in \mathbb{R}^{d \times r}$  initialized from a Gaussian distribution, with rank  $r \ll d$ . Consequently, the weight matrix at layer  $l$  for the task  $t$  is given by

$$W_t^l = W_{t-1}^l + B_t^l A_t^l = W_0^l + \sum_{j=1}^t B_j^l A_j^l. \quad (1)$$

Given the input  $X_t = \{x_{t,i}\}_{i=1}^B \in \mathcal{X}_t$ , the network output  $Z_t$  is obtained via layer-wise propagation as  $Z_t^l = f(X_t^l; W_t^l)$ ,  $X_t^{l+1} = \sigma(Z_t^l)$ , where  $\sigma(\cdot)$  is nonlinear activation, and  $l \in \{1, \dots, L\}$ ,  $X_t^1 = X_t$ . The model is optimized by minimizing the Cross Entropy (CE) loss:

$$\mathcal{L}_{CE}^l = CE(Z_t, Y_t), \quad (2)$$

where  $Y_t \in \mathcal{Y}_t$  denotes the corresponding labels.

**After the training of task  $t$ .** Following (Tan et al., 2024; Wu et al., 2025), the unified classifier is retrained after

each task using pseudo-features sampled from Gaussian distributions of all learned classes. Specifically, at each task, we compute class prototypes  $\mu_c$  and covariance matrices  $\Sigma_c$  for all classes  $c \in \mathcal{Y}_t$  based on the current feature  $Z_t$ . We then draw  $s_c$  pseudo-features  $h^c \sim \mathcal{N}(\mu_c, \Sigma_c)$  for each class. The classification head is subsequently optimized using the CE loss on these synthetic samples:

$$\mathcal{L}_{CA} = -\frac{1}{s_c |\mathcal{C}|} \sum_{c=1}^{|\mathcal{C}|} \sum_{i=1}^{s_c} CE(h_i^c, c) \quad (3)$$

where  $\mathcal{C}$  denotes the set of all accumulated classes.

#### 3.2. Overview

As illustrated in Figure 2, our framework follows a two-stage optimization scheme (Zhang et al., 2023a). During the training of task  $t$ , we employ Gradient Recycling alongside Local Mismatch Suppression (LMS). However, the accumulation of non-orthogonal gradients within task-specific modules inevitably biases the classifier. To rectify this, after the training of task  $t$ , we employ Classifier Alignment (CA) using generated samples. Within this phase, we incorporate Global Mismatch Suppression (GMS), which synthesizes OOD prototypes from ID statistics to assist CA in recalibrating global decision boundaries.

#### 3.3. Gradient Recycling

In this section, we introduce Gradient Recycling, which recycles the gradients discarded by orthogonal constraints

(Liang & Li, 2024; Liu & Chang, 2025) into a task-specific lightweight modules, enhancing plasticity of current task.

**Residual Subspace Construction.** Before learning task  $t$ , the feature space of the new task is isolated by temporarily removing the accumulated LoRA parameters from previous tasks (Liu & Chang, 2025). Specifically, the effective backbone weights at layer  $l$  are adjusted as  $\bar{\mathbf{W}}_t^l = \mathbf{W}_0^l - \sum_{j=1}^{t-1} \mathbf{B}_j^l \mathbf{A}_j^l$ . Forwarding the current task data  $\mathcal{D}_t$  through this subtraction-adjusted model yields layer-wise features  $\bar{\mathbf{X}}_t^l$ . The empirical covariance matrix is then computed as  $\Sigma_t^l = \frac{1}{N_t} (\bar{\mathbf{X}}_t^l)^\top \bar{\mathbf{X}}_t^l$ , where  $N_t$  denotes the number of samples in task  $t$ . Principal directions of the residual feature space are obtained via singular value decomposition (SVD) (Baker, 2005),  $\mathbf{U}_t^l \Lambda_t^l (\mathbf{U}_t^l)^\top = \text{SVD}(\Sigma_t^l)$ , where  $\mathbf{U}_t^l$  contains orthogonal eigenvectors and  $\Lambda_t^l$  represents the associated singular values. Based on these principal directions, residual subspace basis at layer  $l$  is defined by the top- $k$  components,  $\mathbf{P}_t^l = (\mathbf{U}_t^l)_k$ . Accordingly, the matrix  $\mathbf{P}_t^l (\mathbf{P}_t^l)^\top$  defines an orthogonal projector onto the residual subspace, which restricts parameter updates to directions that minimally interfere with previously learned tasks.

**Gradient Recycling Optimization.** Based on projector  $\mathbf{P}_t^l (\mathbf{P}_t^l)^\top$ , optimization for task  $t$  is performed through Gradient Recycling Module  $\hat{G}_t$ , as shown in Figure 2(a). Let  $\mathbf{g}_{t,s}^l$  denote gradient at layer  $l$  and training step  $s$ , and gradient decomposition induced by projection is given by

$$\mathbf{g}_{t,s}^l = \mathbf{P}_t^l (\mathbf{P}_t^l)^\top \mathbf{g}_{t,s}^l + (\mathbf{I} - \mathbf{P}_t^l (\mathbf{P}_t^l)^\top) \mathbf{g}_{t,s}^l, \quad (4)$$

where  $\mathbf{I}$  denotes the identity matrix, and the two terms correspond to orthogonal components responsible for stability and residual components accounting for plasticity.

$\hat{G}_t$  is implemented via a dual-component LoRA architecture, as shown in Figure 2(a\*) full variant. The effective LoRA parameters are decomposed into a shared orthogonal component  $\{\mathbf{A}_t^l, \mathbf{B}_t^l\}$  and a task-specific component  $\{\hat{\mathbf{A}}_t^l, \hat{\mathbf{B}}_t^l\}$ . The forward mapping is given by

$$\begin{aligned} \hat{W}_t^l &= W_0^l + \sum_{j=1}^{t-1} \mathbf{B}_j^l \mathbf{A}_j^l + (\mathbf{B}_t^l + \hat{\mathbf{B}}_t^l) (\mathbf{A}_t^l + \hat{\mathbf{A}}_t^l) \\ &= W_{t-1}^l + \hat{\mathbf{B}}_t^l \hat{\mathbf{A}}_t^l = W_{t-1}^l + \hat{G}_t^l, \end{aligned} \quad (5)$$

where  $\tilde{\mathbf{B}}_t^l$  and  $\tilde{\mathbf{A}}_t^l$  denote the effective unconstrained LoRA parameters. The task-specific components  $\hat{\mathbf{A}}_t^l$  and  $\hat{\mathbf{B}}_t^l$  are initialized to zero, ensuring that the initial training dynamics coincide with standard LoRA (Hu et al., 2022). By gradient additivity (Rahel & Hubert, 1991; Müller & Yao, 2010), the orthogonal component  $\mathbf{P}_t^l (\mathbf{P}_t^l)^\top \mathbf{g}_{t,s}^l$  is assigned to the shared parameter module  $\{\mathbf{A}_t^l, \mathbf{B}_t^l\}$  to preserve stability across tasks, while the residual component  $(\mathbf{I} - \mathbf{P}_t^l (\mathbf{P}_t^l)^\top) \mathbf{g}_{t,s}^l$  is assigned to the task-specific module  $\{\hat{\mathbf{A}}_t^l, \hat{\mathbf{B}}_t^l\}$  to facilitate plasticity for the current task. Sum-

ming both updates recovers the original unconstrained gradient, with interference confined to the task-specific module.

To improve training efficiency, a simplified variant of  $\hat{G}_t$  is adopted, as illustrated in Figure 2(a\*). Exploiting the commutativity of low-rank updates under mild conditions (Zhang et al., 2023b), the input projection  $A_t$  is frozen and only the task-specific output projection  $\hat{B}_t$  is learned, which is mathematically equivalent to the full formulation while substantially reducing computational cost. The LoRA weight formulation simplifies to:

$$\begin{aligned} \hat{W}_t^l &= W_0^l + \sum_{j=1}^{t-1} B_j^l A_j^l + (B_t^l + \hat{B}_t^l) \cdot A_t^l \\ &= W_{t-1}^l + \hat{B}_t^l A_t^l = W_{t-1}^l + \hat{G}_t^l. \end{aligned} \quad (6)$$

Regardless of the architectural variant, model optimization is performed using the CE loss:

$$\mathcal{L}_{CE} = CE(\hat{Z}_t, Y^t), \quad (7)$$

where  $\hat{Z}_t = f(X^t; \hat{W}_t)$  denotes the model output, and  $\hat{W}_t$  denotes the collection of parameters across all layers. After training, task-specific components are consolidated into the task repository, represented as  $G_t = \{\hat{B}_t, A_t\}$  for the full variant and  $G_t = \{\hat{B}_t\}$  for the efficient variant.

### 3.4. Mismatch Suppression

Upon the completion of task, the system must identify the optimal task-specific LoRA module for each input  $x$  without access to task identity (Wan & Yang, 2025). Let  $\tilde{G}_j$  denote the composite weight matrix for the  $j$ -th task, defined as  $\tilde{G}_j = W_0^l + \sum_{j'=1, j' \neq j}^t B_{j'}^l A_{j'}^l + \hat{G}_j^l$ . A critical challenge in CIL is that historical modules often exhibit high confidence on current task data due to semantic overlap, acting as an OOD disturbance (Xu & Yang, 2025). To mitigate this, we employ entropy as a selection metric (Wang et al., 2025). The optimal module is determined by minimizing the prediction entropy:

$$\tilde{G}^* = \arg \min_{\tilde{G}_j \in \{\tilde{G}_1, \dots, \tilde{G}_t\}} - \sum_{c=1}^{\mathcal{C}} f_c(x; \tilde{G}_j) \log f_c(x; \tilde{G}_j), \quad (8)$$

where  $f_c(x; \tilde{G}_j)$  denotes the predicted probability of class  $c$  given input  $x$  and  $\mathcal{C}$  denotes all accumulated classes.

However, sole reliance on intrinsic entropy is insufficient, as mismatched historical modules may still exhibit high confidence on current task data. To address this, we propose two Mismatch Suppression strategies, comprising LMS and GMS. As illustrated in Figure 2, LMS is integrated into the training phase of task  $t$  while GMS is applied during the classifier alignment phase after task  $t$ .

**Local Mismatch Suppression.** To enhance the discriminability of the selection metric during the training

phase, we propose LMS, which integrates an Outlier Exposure (Hendrycks et al., 2019; Miao et al., 2024) mechanism in training phase. As illustrated in Figure 2(b), LMS treats current data  $X^t$  as OOD samples relative to all historical modules  $\{\tilde{G}_j\}_{j=1}^{t-1}$ . Crucially, instead of enforcing a generic uniform distribution that risks degrading ID feature discriminability, we utilize the ground truth label as a directional OOD target. This allows us to employ the standard CE loss to concentrate probability mass onto this foreign class, effectively suppressing activations on the module’s native classes while preserving historical performance. To reduce the computational overhead of traversing all historical modules, we uniformly sample a single mismatched module index  $j$  per mini-batch to impose this constraint. Formally, given the prediction  $\hat{Z}_t^{\text{LMS}} = f(X^t; \tilde{G}_j)$ , the LMS objective for task  $t$  is given by:  $\mathcal{L}_{\text{LMS}} = CE(\hat{Z}_t^{\text{LMS}}, Y_t)$ . Consequently, the total loss for during the training of task  $t$  is formulated as:

$$L_{\text{total}} = L_{\text{CE}} + L_{\text{LMS}}. \quad (9)$$

**Global Mismatch Suppression.** While LMS effectively prevents historical modules from overfitting to the current data, it provides only local regularization. It fails to capture the converse failure mode, where task-specific modules of the current task over-generalize to classes from previous tasks, leading to global task confusion. To address this limitation, we propose GMS, which synthesizes virtual OOD features across all task-specific modules to enforce a globally consistent decision boundary, as shown in Figure 2(c). Specifically, let  $\mu_c$  denote the prototype of class  $c$  in the shared space, while  $\hat{\mu}_c$  and  $\hat{\Sigma}_c$  denote its prototype and covariance in the task-specific subspace, respectively. The prototype drift vector as  $\delta_c = \hat{\mu}_c - \mu_c$ , which captures the geometric shift between the two spaces (Zhou et al., 2024; Fukuda et al., 2025; Li et al., 2025). Based on these statistics and inspiration, we first compute the semantic affinity  $\alpha_{c',c}$  between the target class  $c$  and the task’s native classes  $c' \in \mathcal{C}_k$  using cosine similarity in the shared space:

$$\alpha_{c',c} = \frac{\exp(\text{sim}(\mu_{c'}, \mu_c)/\tau)}{\sum_{z \in \mathcal{C}_k} \exp(\text{sim}(\mu_z, \mu_c)/\tau)}. \quad (10)$$

Next, we estimate the virtual OOD prototype  $\mu_{c,k}^{\text{OOD}}$  by transferring the weighted geometric drift of these native classes to the target class  $c$ :

$$\mu_{c,k}^{\text{OOD}} = \mu_c + \sum_{c' \in \mathcal{C}_k} \alpha_{c',c} \cdot \delta_{c'}. \quad (11)$$

This process is extended globally to construct a set of OOD prototypes  $P_c^{\text{OOD}} = \{\mu_{c,k}^{\text{OOD}} \mid k \in [1, t], k \neq \text{Task}(c)\}$  for every class  $c$  across all disjoint task subspaces.

Finally, to align the classifier, we sample  $s_c/t$  pseudo-OOD samples  $h^{c,k} \sim \mathcal{N}(\mu_{c,k}^{\text{OOD}}, \hat{\Sigma}_c)$  for each task  $k$  subspace.

The classification head is retrained using the CE loss to enforce consistent prediction across different task subspaces:

$$\mathcal{L}_{\text{GMS}} = -\frac{1}{s_c |\mathcal{C}|} \sum_{c=1}^{|\mathcal{C}|} \sum_{i=1}^{s_c/t} \sum_{k=1}^t CE(h_i^{c,k}, c). \quad (12)$$

By further incorporating the standard classifier alignment loss on available data, the overall objective for classifier retraining is given by:

$$\mathcal{L}_{\text{head}} = \mathcal{L}_{\text{CA}} + \mathcal{L}_{\text{GMS}}. \quad (13)$$

### 3.5. Theoretical Analysis

To analyze the plasticity of the proposed method, we study the attainable population risk under different parameter constraints from an optimization perspective. We assume that the population risk  $R_t(w)$  is  $\mu$ -strongly convex and  $L$ -smooth with respect to the model parameters.

**Theorem 3.1.** *For each task  $t = 1, \dots, T$ , let  $\widehat{\mathbf{W}}_{\perp,t}$  denote the solution obtained under the orthogonal constraint induced by the layer-wise projectors  $\{\mathbf{\Pi}_t^l\}_{l=1}^L$ , and let  $\mathbf{W}_t^*$  be the unconstrained population minimizer of  $R_t$ . Then, with probability at least  $1 - \delta$ , it holds that*

$$\begin{aligned} \sum_{t=1}^T \left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_t^*) \right) &\leq \sum_{t=1}^T \sum_{l=1}^L \left( 4\mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l) \right. \\ &\quad \left. + \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2 + 2\sqrt{\frac{\log(2TL/\delta)}{2n_t}} \right), \end{aligned} \quad (14)$$

where  $\mathbf{W}_{\perp,t}^*$  is the population minimizer under orthogonal constraint and  $\mathfrak{R}_{n_t}(\cdot)$  denotes the Rademacher complexity.

As established in Theorem 3.1, when model updates are constrained to orthogonal subspace, the cumulative performance gap relative to the unconstrained optimum is dominated by  $\|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2$ . Consequently, even under perfect samples and optimization, orthogonal constraints induce an unavoidable plasticity gap.

**Theorem 3.2.** *Suppose the model parameters are augmented in the LoRA form  $\mathbf{W}^l = \mathbf{W}_0^l + \sum_{j=1}^{t-1} \mathbf{B}_j^l \mathbf{A}_j^l + (\mathbf{B}_t^l + \hat{\mathbf{B}}_t^l)(\mathbf{A}_t^l + \hat{\mathbf{A}}_t^l)$ , where  $\mathbf{B}_t^l \mathbf{A}_t^l$  denotes the shared module learned under orthogonal constraints, and  $\hat{\mathbf{B}}_t^l \hat{\mathbf{A}}_t^l$  is a task-specific module allocated to task  $t$ . Let  $\mathbf{W}_{\text{GR},t}^*$  denote the population minimizer under this parameterization. Then the cumulative optimality gap across tasks satisfies*

$$\begin{aligned} \sum_{t=1}^T \left( R_t(\mathbf{W}_{\text{GR},t}^*) - R_t(\mathbf{W}_t^*) \right) &\leq \sum_{t=1}^T \sum_{l=1}^L \\ &\quad \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) - \Pi_{\mathcal{S}_t^l}(\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*))\|_F^2, \end{aligned} \quad (15)$$

where  $\Pi_{\mathcal{S}_t^l}$  denotes the projection onto  $\mathcal{S}_t^l = \text{range}(\mathbf{B}_t^l \hat{\mathbf{A}}_t^l)$ .

Table 1. Last and average performance results on CIFAR-100 and ImageNet-R under the long-term setting (20 and 50 tasks). The mean and standard deviation of three trials are provided. We compare all methods using the same ViT-B/16-IN21K backbone, seeds, and class orders. Red denotes the best result and blue denotes the second-best result in each column.

Method	CIFAR-100				ImageNet-R			
	20		50		20		50	
	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
L2P (Wang et al., 2022b)	79.51 $\pm$ 0.67	85.50 $\pm$ 1.23	73.91 $\pm$ 1.67	81.90 $\pm$ 0.98	69.64 $\pm$ 0.42	75.28 $\pm$ 0.57	55.89 $\pm$ 1.59	62.98 $\pm$ 2.89
DualPrompt (Wang et al., 2022a)	80.44 $\pm$ 1.38	86.96 $\pm$ 1.98	76.66 $\pm$ 0.74	85.18 $\pm$ 0.92	66.61 $\pm$ 0.58	72.45 $\pm$ 0.37	61.50 $\pm$ 0.86	68.63 $\pm$ 1.31
CODA-Prompt (Smith et al., 2023)	81.36 $\pm$ 0.88	88.17 $\pm$ 0.61	55.45 $\pm$ 0.48	68.39 $\pm$ 0.53	69.96 $\pm$ 0.50	75.34 $\pm$ 0.85	48.89 $\pm$ 0.90	55.59 $\pm$ 2.67
SLCA (Zhang et al., 2023a)	89.62 $\pm$ 0.18	93.03 $\pm$ 1.09	87.90 $\pm$ 0.17	91.96 $\pm$ 1.30	75.53 $\pm$ 0.42	80.65 $\pm$ 1.16	68.95 $\pm$ 4.46	71.15 $\pm$ 10.57
InfLoRA (Liang & Li, 2024)	81.59 $\pm$ 0.94	87.33 $\pm$ 2.32	55.19 $\pm$ 2.83	69.96 $\pm$ 3.55	73.01 $\pm$ 0.82	79.76 $\pm$ 0.83	61.91 $\pm$ 1.31	71.20 $\pm$ 0.54
EASE (Zhou et al., 2024)	86.32 $\pm$ 0.48	90.83 $\pm$ 1.14	77.43 $\pm$ 3.11	84.59 $\pm$ 1.22	73.78 $\pm$ 0.47	80.28 $\pm$ 0.67	68.53 $\pm$ 0.04	75.49 $\pm$ 1.06
SSIAT (Tan et al., 2024)	90.07 $\pm$ 0.56	93.54 $\pm$ 0.83	87.33 $\pm$ 1.46	91.63 $\pm$ 0.89	78.31 $\pm$ 0.53	82.35 $\pm$ 0.52	74.52 $\pm$ 0.35	79.01 $\pm$ 0.54
CL-LoRA (He et al., 2025)	83.75 $\pm$ 1.39	88.98 $\pm$ 1.87	71.14 $\pm$ 2.44	78.80 $\pm$ 2.72	77.20 $\pm$ 0.66	83.45 $\pm$ 0.56	68.36 $\pm$ 0.82	76.80 $\pm$ 1.27
LoRA-DRS (Liu & Chang, 2025)	88.76 $\pm$ 0.22	92.35 $\pm$ 0.90	86.51 $\pm$ 1.05	91.17 $\pm$ 0.97	74.96 $\pm$ 0.17	80.66 $\pm$ 0.67	72.17 $\pm$ 0.27	78.05 $\pm$ 0.70
MOS (Sun et al., 2025)	89.48 $\pm$ 0.39	92.97 $\pm$ 1.01	87.05 $\pm$ 0.90	91.44 $\pm$ 1.21	75.04 $\pm$ 0.75	80.39 $\pm$ 0.87	66.92 $\pm$ 0.34	74.81 $\pm$ 0.41
MACIL (Wu et al., 2025)	90.31 $\pm$ 0.17	93.47 $\pm$ 0.78	85.09 $\pm$ 0.87	90.89 $\pm$ 1.24	79.46 $\pm$ 0.15	84.25 $\pm$ 0.51	70.10 $\pm$ 1.02	77.47 $\pm$ 0.55
<b>E-GR-LoRA</b>	<b>91.14<math>\pm</math>0.02</b>	<b>94.11<math>\pm</math>0.86</b>	<b>89.76<math>\pm</math>0.23</b>	<b>93.16<math>\pm</math>0.77</b>	79.21 $\pm$ 0.22	<b>84.41<math>\pm</math>0.32</b>	<b>76.64<math>\pm</math>0.39</b>	<b>81.94<math>\pm</math>0.82</b>
<b>GR-LoRA</b>	<b>91.46<math>\pm</math>0.11</b>	<b>94.41<math>\pm</math>0.91</b>	<b>90.03<math>\pm</math>0.28</b>	<b>93.38<math>\pm</math>0.96</b>	<b>80.23<math>\pm</math>0.27</b>	<b>85.05<math>\pm</math>0.47</b>	<b>76.74<math>\pm</math>0.63</b>	<b>82.64<math>\pm</math>0.84</b>

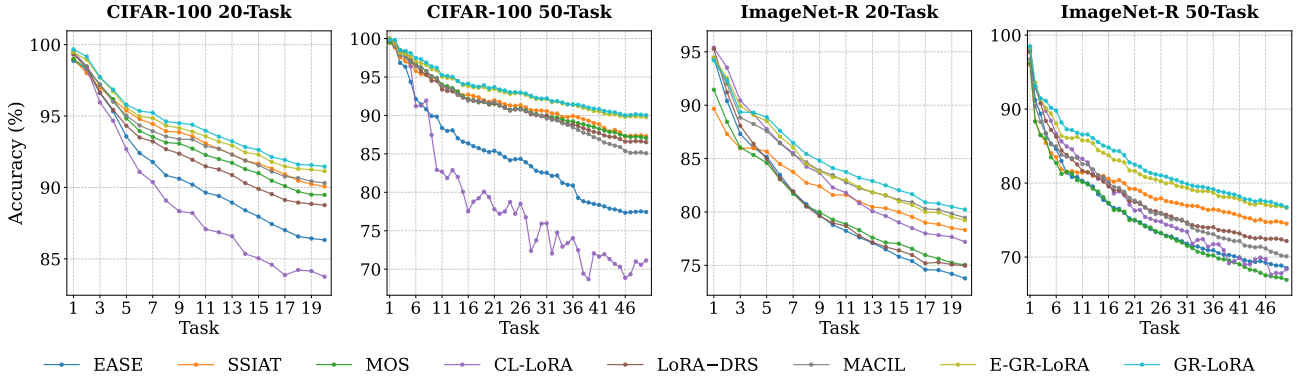


Figure 3. The performance of each learning session under different settings of ImageNet-R and CIFAR100. All methods are initialized with ViT-B/16-IN21k. These curves are plotted by calculating the average performance across three different seeds.

As established in Theorem 3.2, by introducing task-specific module to recycle the residual gradient components, the plasticity gap is no longer governed by the entire discarded gradient, but only by the portion that cannot be represented within the subspace spanned by the task-specific module. This result demonstrates that explicitly recycling non-orthogonal gradient components into task-specific module effectively restores the model’s usable optimization space, thereby mitigating the plasticity gap induced by orthogonal constraints. The detailed proofs are provided in Appendix C.

## 4. Experiments

### 4.1. Setup

**Datasets.** We train and evaluate our method on four widely used CIL benchmarks. The benchmarks include CIFAR-100 (Krizhevsky et al., 2009) with 100 categories, and ImageNet-R (Hendrycks et al., 2021a), CUB-200 (Wah

et al., 2011), and ImageNet-A (Hendrycks et al., 2021b), each comprising 200 categories. Following standard CIL protocols (Liu & Chang, 2025), all datasets are evaluated under 10-task, 20-task, and 50-task settings.

**Evaluation metrics.** Following CIL protocols (Liu & Chang, 2025), we evaluate performance using two metrics: the final accuracy  $\mathcal{A}_{Last}$  and the average accuracy  $\mathcal{A}_{Avg}$ . Let  $a_{i,j}$  denote the test accuracy on the  $j$ -th task ( $j \leq i$ ) after learning the  $i$ -th task. The average accuracy at stage  $i$  is defined as  $\mathcal{A}_i = \frac{1}{i} \sum_{j=1}^i a_{i,j}$ . Accordingly, the final accuracy is given by  $\mathcal{A}_{Last} = \frac{1}{T} \sum_{j=1}^T a_{T,j}$ , where  $T$  denotes the total number of tasks, and the overall average accuracy is computed as  $\mathcal{A}_{Avg} = \frac{1}{T} \sum_{i=1}^T \mathcal{A}_i$ .

**Baselines.** We consider four categories of PTM-based CIL methods for comparison: (1) prompt-based methods, including L2P (Wang et al., 2022b), DualPrompt (Wang et al., 2022a), and CODA-Prompt (Smith et al., 2023); (2)

Table 2. Last and average performance results on four benchmark datasets (10 tasks) are reported. The mean and standard deviation of three trials are provided. We compare all methods using the same ViT-B/16-IN21K backbone, seeds, and class orders. **Red** denotes the best result and **blue** denotes the second-best result in each column.

Method	CIFAR-100		ImageNet-A		ImageNet-R		CUB-200	
	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
L2P (Wang et al., 2022b)	84.06 $\pm$ 0.88	88.26 $\pm$ 1.34	44.04 $\pm$ 0.93	51.24 $\pm$ 2.26	72.34 $\pm$ 0.17	77.36 $\pm$ 0.64	67.02 $\pm$ 1.90	79.62 $\pm$ 1.60
DualPrompt (Wang et al., 2022a)	86.93 $\pm$ 0.24	91.13 $\pm$ 0.32	53.19 $\pm$ 0.74	64.59 $\pm$ 0.08	69.10 $\pm$ 0.62	74.28 $\pm$ 0.66	68.48 $\pm$ 0.47	80.59 $\pm$ 1.50
CODA-Prompt (Smith et al., 2023)	83.21 $\pm$ 3.39	87.71 $\pm$ 3.17	52.08 $\pm$ 0.12	63.92 $\pm$ 0.12	73.31 $\pm$ 0.50	78.47 $\pm$ 0.53	77.23 $\pm$ 1.12	81.90 $\pm$ 0.85
SLCA (Zhang et al., 2023a)	91.26 $\pm$ 0.37	94.29 $\pm$ 0.92	61.05 $\pm$ 0.63	68.88 $\pm$ 2.31	79.35 $\pm$ 0.28	83.29 $\pm$ 0.46	84.68 $\pm$ 0.09	90.77 $\pm$ 0.79
InfLoRA (Liang & Li, 2024)	86.20 $\pm$ 0.70	90.58 $\pm$ 1.52	47.75 $\pm$ 0.51	58.13 $\pm$ 0.56	75.88 $\pm$ 0.32	81.90 $\pm$ 0.65	69.04 $\pm$ 1.25	81.83 $\pm$ 0.45
EASE (Zhou et al., 2024)	88.22 $\pm$ 0.44	92.02 $\pm$ 0.76	54.93 $\pm$ 1.14	63.92 $\pm$ 0.76	75.91 $\pm$ 0.17	81.38 $\pm$ 0.29	85.04 $\pm$ 1.42	90.93 $\pm$ 1.03
SSIAT (Tan et al., 2024)	91.48 $\pm$ 0.24	94.28 $\pm$ 0.90	62.58 $\pm$ 1.58	<b>70.73<math>\pm</math>1.44</b>	79.54 $\pm$ 0.24	83.67 $\pm$ 0.57	89.83 $\pm$ 0.53	93.76 $\pm$ 0.52
CL-LoRA (He et al., 2025)	87.47 $\pm$ 0.60	91.37 $\pm$ 1.30	57.12 $\pm$ 1.31	68.17 $\pm$ 1.91	79.78 $\pm$ 0.17	85.10 $\pm$ 0.67	76.28 $\pm$ 2.70	86.81 $\pm$ 1.24
LoRA-DRS (Liu & Chang, 2025)	90.03 $\pm$ 0.07	93.24 $\pm$ 0.97	57.58 $\pm$ 0.79	66.72 $\pm$ 0.79	75.96 $\pm$ 0.36	81.82 $\pm$ 0.85	87.58 $\pm$ 0.28	92.30 $\pm$ 0.61
MOS (Sun et al., 2025)	91.54 $\pm$ 0.45	94.18 $\pm$ 1.15	57.54 $\pm$ 0.37	64.50 $\pm$ 1.44	77.65 $\pm$ 0.50	81.94 $\pm$ 0.66	89.88 $\pm$ 0.29	93.52 $\pm$ 0.61
MACIL (Wu et al., 2025)	91.86 $\pm$ 0.22	94.44 $\pm$ 0.96	63.15 $\pm$ 0.17	<b>70.54<math>\pm</math>1.79</b>	<b>81.82<math>\pm</math>0.22</b>	<b>85.76<math>\pm</math>0.32</b>	<b>90.23<math>\pm</math>0.13</b>	<b>93.78<math>\pm</math>0.40</b>
<b>E-GR-LoRA</b>	<b>92.00<math>\pm</math>0.03</b>	<b>94.56<math>\pm</math>0.96</b>	<b>63.22<math>\pm</math>0.86</b>	69.96 $\pm$ 2.20	80.59 $\pm$ 0.21	85.13 $\pm$ 0.29	89.79 $\pm$ 0.15	93.72 $\pm$ 0.51
<b>GR-LoRA</b>	<b>91.97<math>\pm</math>0.17</b>	<b>94.65<math>\pm</math>0.97</b>	<b>63.60<math>\pm</math>0.41</b>	70.24 $\pm$ 1.85	<b>82.09<math>\pm</math>0.18</b>	<b>86.20<math>\pm</math>0.28</b>	<b>89.91<math>\pm</math>0.44</b>	<b>93.85<math>\pm</math>0.73</b>

adapter-based methods, such as SSIAT (Tan et al., 2024), EASE (Zhou et al., 2024), and MOS (Sun et al., 2025); (3) LoRA-based methods, including InfLoRA (Liang & Li, 2024), LoRA-DRS (Liu & Chang, 2025), CL-LoRA (He et al., 2025), and MACIL (Wu et al., 2025); and (4) the fine-tuning method SLCA (Zhang et al., 2023a).

**Implementation Details.** We adopt ViT-B/16 (Dosovitskiy et al., 2021) pretrained on ImageNet-21K (Russakovsky et al., 2015) as the backbone and integrate LoRA with rank  $r = 10$  into the key and value projections of all attention layers. Following prior work (Liang & Li, 2024; Liu & Chang, 2025), the model is optimized with Adam and a cosine annealing schedule. All results are reported as the mean and standard deviation over three runs with different random seeds. More implementation details and hyperparameter settings are provided in the Appendix A. The source code is available at <https://github.com/njustkmg/ICML26-GR-LoRA>.

## 4.2. Main Results

**Performance on Long Task Sequences.** We evaluate the effectiveness of GR-LoRA under long task sequences, with results summarized in Table 1. The experimental results lead to the following observations: (1) GR-LoRA consistently outperforms all competing methods across CIFAR-100 and the more challenging ImageNet-R benchmark under both the 20-task and 50-task settings, achieving superior performance. This demonstrates its strong robustness to error accumulation and distributional drift in long task CIL. (2) Compared with orthogonality-based methods, GR-LoRA attains a more favorable stability-plasticity trade-off by explicitly recycling discarded gradients into task-specific modules, rather than suppressing them through hard constraints. This design allows the model to preserve previously ac-

quired knowledge while maintaining sufficient plasticity for learning new tasks, and the resulting advantage becomes increasingly pronounced as the task sequence length grows. (3) The efficient variant E-GR-LoRA consistently ranks second across all settings, indicating that the proposed gradient recycling remains effective even under reduced computational budgets and confirming its robustness to practical resource constraints. Additional results on long task sequences are provided in Appendix B.1.

**Versatility on Standard Benchmarks.** To assess the generalizability of GR-LoRA, we further evaluate it under standard short-task settings with 10 tasks on four benchmark datasets, as reported in Table 2. The experimental results can be summarized as follows: (1) GR-LoRA demonstrates consistently strong performance across all evaluated settings, indicating robust generalization under short task sequences; (2) on more challenging benchmarks, including ImageNet-A and the fine-grained CUB-200, GR-LoRA remains competitive and maintains stable performance. Detailed per-task results are reported in Appendix B.2.

## 4.3. Ablation Study

**Impact of each component.** As shown in Table 3, we analyze the impact of each proposed component on ImageNet-R. For a fair comparison, following SLCA (Zhang et al., 2023a), we adopt Classifier Alignment applied after training as the baseline, which corresponds to the component in GMS that focuses solely on preserving ID classification performance. The experimental results lead to the following observations: (1) under the presence of CA, recycling the components discarded during orthogonal projection into task-specific modules yields a clear performance improvement, demonstrating the effectiveness of gradient reuse; (2) to mitigate potential inference-time degradation caused by

mis-selection of task-specific module, the two proposed suppression strategies consistently improve model performance, validating their necessity in our method; (3) compared with the ‘‘Single LoRA’’ strategy that allocates an independent adapter for each task, our method achieves significantly better performance, highlighting the importance of sharing orthogonal knowledge across tasks rather than isolating parameters for individual tasks. Additional ablation results on other Datasets are provided in Appendix B.3, further confirming the generalizability of these components across different datasets.

Table 3. Ablation study of individual component contributions on the 20-task ImageNet-R benchmark. **CA** denotes Classifier Alignment, **LMS** denotes Local Mismatch Suppression, and **GMS** denotes Global Mismatch Suppression.  $\checkmark$  and  $\times$  indicate whether the corresponding component is enabled or disabled, respectively.

Components Ablations	Components			$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
	CA	LMS	GMS		
<b>LoRA-DRS</b>	$\checkmark$	$\times$	$\times$	$77.90 \pm 0.26$	$82.50 \pm 0.43$
<b>Single LoRA</b>	$\checkmark$	$\checkmark$	$\checkmark$	$70.50 \pm 0.93$	$76.79 \pm 0.53$
<b>GR-LoRA</b>	$\checkmark$	$\times$	$\times$	$78.22 \pm 0.44$	$83.77 \pm 0.21$
	$\checkmark$	$\checkmark$	$\times$	$79.12 \pm 0.24$	$84.33 \pm 0.22$
	$\checkmark$	$\times$	$\checkmark$	$79.65 \pm 0.30$	$84.73 \pm 0.05$
	$\checkmark$	$\checkmark$	$\checkmark$	$80.23 \pm 0.27$	$85.05 \pm 0.47$

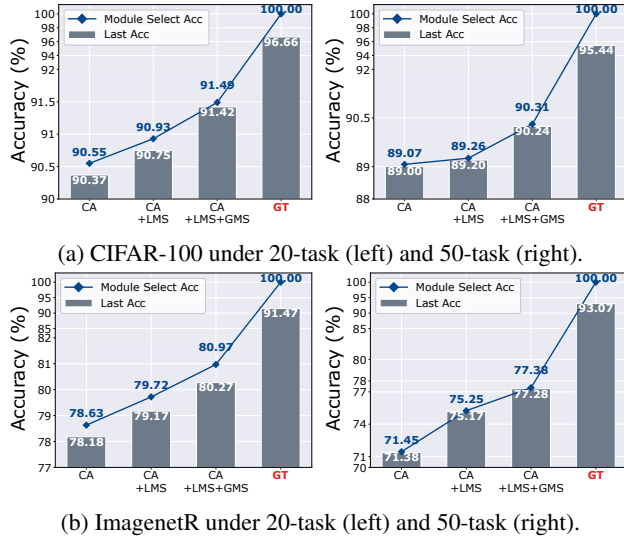


Figure 4. Module selection accuracy and Last accuracy across different datasets under the 20-task and 50-task setting.

**Module Select Accuracy Analysis.** As shown in Figure 4, we analyze the evolution of module selection accuracy for each task on CIFAR-100 and ImageNet-R under the 20-task and 50-task setting after integrating LMS and GMS. The results show that the proposed suppression strategies significantly improve the final model performance by increasing the accuracy of task-specific module selection at inference

time. Specifically, the two suppression strategies encourage each module to exhibit correct behavior on features corresponding to its associated task, while reducing confidence on mismatched features originating from other tasks. This design is analogous to a  $K+1$  classification in OOD detection (Miao et al., 2024), and effectively suppresses erroneous module selection at inference time.

Table 4. Comparison on the 20-task ImageNet-R benchmark in terms of Multiply-Accumulate operations (MACs), learnable parameters (LP), and model performance.

Method	MACs (G)	LP (M)	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
L2P	35.85	0.20	69.64	75.28
CODA-Prompt	33.73	3.99	69.96	75.34
SLCA	16.86	85.81	75.53	80.65
LoRA-DRS	18.31	0.38	74.96	80.66
SSIAT	17.10	1.20	78.31	82.35
MOS	16.92	0.76	75.04	80.39
MACIL	21.51	1.19	79.46	84.25
E-GR-LoRA	18.27	0.38	79.21	84.41
GR-LoRA	18.39	0.74	80.23	85.05
GR-LoRA( $r=5$ )	17.62	0.38	80.22	84.98
GR-LoRA( $r=20$ )	19.91	1.48	80.53	85.34

**Computational Overhead.** As show in Table 4 we compare the Multiply-Accumulate operations computation (Molchanov et al., 2017), trainable parameters, and accuracy metrics ( $\mathcal{A}_{Last}$  and  $\mathcal{A}_{Avg}$ ). The results show that GR-LoRA consistently achieves superior accuracy while incurring only modest training and computation overhead, demonstrating a favorable performance–efficiency trade-off. Moreover, the efficient variant E-GR-LoRA substantially reduces the number of trainable parameters with only negligible performance degradation, making it particularly suitable for resource-constrained scenarios. In addition, sensitivity analysis with respect to the rank  $r$  shows that the proposed method is largely insensitive to rank variations. Notably, setting  $r = 5$  attains training efficiency comparable to E-GR-LoRA while preserving higher accuracy, further highlighting the robustness of the proposed design.

Table 5. Inference latency comparison on the 50-task ImageNet-R benchmark.

Method	Inference Time (s)		$\mathcal{A}_{Last}$
	Per Image	Total Test	
LoRA-DRS	1.78	15.93	72.37
MOS	2.56	400.85	66.95
MACIL	2.01	19.04	71.13
GR-LoRA	2.58	427.32	77.43

We further analyze the inference latency of GR-LoRA and

recent SOTA methods on the 50-task ImageNet-R benchmark, as reported in Table 5. Since our task-inference strategy relies on an entropy-based selection mechanism, each input needs to be evaluated by multiple task-specific branches during inference. This design introduces non-negligible time overhead, which increases with the number of tasks. However, such linear computational scaling is a systemic limitation of entropy-based routing architectures rather than a limitation specific to GR-LoRA. Importantly, the accuracy gains brought by GR-LoRA are sufficient to compensate for this additional inference cost.

## 5. Conclusion

In this paper, we propose GR-LoRA to resolve the stability-plasticity dilemma in CIL via a Gradient Recycling mechanism. By utilizing task-specific LoRA modules to recycle non-orthogonal gradient updates and employing entropy-based selection, our method theoretically overcomes the limitations of strict orthogonality while balancing stability and plasticity. To ensure reliable selection, we further propose two suppression strategies consisting of LMS and GMS that minimizes mismatched module activation. Extensive experiments validate the superiority of GR-LoRA, particularly in long-sequence tasks. Additionally, we demonstrate that both our efficient E-GR-LoRA variant and simple rank adjustments offer flexible trade-offs between high parameter efficiency and superior performance. Despite these advantages, the entropy-based selection mechanism introduces additional inference latency as the number of tasks increases. Reducing this latency while preserving task-specific routing benefits remains important future work.

## Acknowledgement

The authors are grateful to the Area Chairs and the anonymous reviewers for their constructive comments. This work is partially supported by the NSFC (62276131), Natural Science Foundation of Jiangsu Province of China under Grant (BK20240081).

## Impact Statement

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

Aljundi, R., Chakravarty, P., and Tuytelaars, T. Expert gate: Lifelong learning with a network of experts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3366–3375, 2017.

Baker, K. Singular value decomposition tutorial. *The Ohio State University*, pp. 22, 2005.

CHEN, Q., Shui, C., Han, L., and Marchand, M. On the stability-plasticity dilemma in continual meta-learning: Theory and algorithm. In *Neural Information Processing Systems*, 2023.

De Lange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., and Tuytelaars, T. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houshy, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.

Du, X., Wang, Z., Cai, M., and Li, Y. VOS: learning what you don’t know by virtual outlier synthesis. 2022.

French, R. M. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.

French, R. M. and Ferrara, A. Modeling time perception in rats: Evidence for catastrophic interference in animal learning. In *Proceedings of the Twenty-first Annual Conference of the Cognitive Science Society*, pp. 173–178. Psychology Press, 2020.

Fukuda, T., Kera, H., and Kawamoto, K. Adapter merging with centroid prototype mapping for scalable class-incremental learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 4884–4893, 2025.

Han, Z., Gao, C., Liu, J., Zhang, J., and Zhang, S. Q. Parameter-efficient fine-tuning for large models: A comprehensive survey. *Trans. Mach. Learn. Res.*, 2024.

He, J., Duan, Z., and Zhu, F. Cl-lora: Continual low-rank adaptation for rehearsal-free class-incremental learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 30534–30544, 2025.

Hendrycks, D., Mazeika, M., and Dietterich, T. G. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.

Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8340–8349, 2021a.

- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15262–15271, 2021b.
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pp. 2790–2799. PMLR, 2019.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al. Lora: Low-rank adaptation of large language models. *International Conference on Learning Representations*, 1(2):3, 2022.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Li, T., Wan, F., Lin, Y., and Yang, Y. Feature drift oriented distribution reconstruction for imbalanced class incremental learning. In *European Conference on Artificial Intelligence*, pp. 193–200, 2025.
- Liang, Y.-S. and Li, W.-J. Inflora: Interference-free low-rank adaptation for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23638–23647, 2024.
- Lin, H., Shao, Y., Qian, W., Pan, N., Guo, Y., and Liu, B. Class incremental learning via likelihood ratio based task prediction. In *International Conference on Learning Representations*, 2024.
- Liu, W., Wang, X., Owens, J., and Li, Y. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, pp. 21464–21475, 2020.
- Liu, X. and Chang, X. Lora subtraction for drift-resistant space in exemplar-free continual learning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 15308–15318, 2025.
- Mermillod, M., Bugaiska, A., and Bonin, P. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects, 2013.
- Miao, W., Pang, G., Bai, X., Li, T., and Zheng, J. Out-of-distribution detection in long-tailed recognition with calibrated outlier class learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- Molchanov, P., Tyree, S., Karras, T., Aila, T., and Kautz, J. Pruning convolutional neural networks for resource efficient inference. In *International Conference on Learning Representations*, 2017.
- Müller, H.-G. and Yao, F. Additive modelling of functional gradients. *Biometrika*, pp. 791–805, 2010.
- Pan, S. The instability of matching with overconfident agents. *Games Econ. Behav.*, 113:396–415, 2019.
- Parisi, G. I., Kemker, R., Part, J. L., Kanan, C., and Wermter, S. Continual lifelong learning with neural networks: A review. *Neural networks*, 113:54–71, 2019.
- Polovinkin, E. S. Strongly convex analysis. *Sbornik: Mathematics*, pp. 259, 1996.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. Pmlr, 2021.
- Rahel, F. J. and Hubert, W. A. Fish assemblages and habitat gradients in a rocky mountain–great plains stream: biotic zonation and additive patterns of community change. *Transactions of the American Fisheries Society*, pp. 319–332, 1991.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, pp. 211–252, 2015.
- Smith, J. S., Karlinsky, L., Gutta, V., Cascante-Bonilla, P., Kim, D., Arbelles, A., Panda, R., Feris, R., and Kira, Z. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11909–11919, 2023.
- Steiner, A., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., and Beyer, L. How to train your vit? data, augmentation, and regularization in vision transformers. *Trans. Mach. Learn. Res.*, 2022.
- Sun, H.-L., Zhou, D.-W., Zhao, H., Gan, L., Zhan, D.-C., and Ye, H.-J. Mos: Model surgery for pre-trained model-based class-incremental learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 20699–20707, 2025.
- Sun, Y., Ming, Y., Zhu, X., and Li, Y. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning*, pp. 20827–20840. PMLR, 2022.
- Tan, Y., Zhou, Q., Xiang, X., Wang, K., Wu, Y., and Li, Y. Semantically-shifted incremental adapter-tuning is a continual vitransformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 23252–23262, 2024.

- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.
- Wan, F. and Yang, Y. Probabilistic group mask guided discrete optimization for incremental learning. In *International Conference on Machine Learning*, 2025.
- Wang, Y., Zhou, D.-W., and Ye, H.-J. Integrating task-specific and universal adapters for pre-trained model-based class-incremental learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 806–816, 2025.
- Wang, Z., Zhang, Z., Ebrahimi, S., Sun, R., Zhang, H., Lee, C.-Y., Ren, X., Su, G., Perot, V., Dy, J., et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European conference on computer vision*, pp. 631–648. Springer, 2022a.
- 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. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 139–149, 2022b.
- Wu, F., Cheng, L., Tang, S., Zhu, X., Fang, C., Zhang, D., and Wang, M. Navigating semantic drift in task-agnostic class-incremental learning. In *International Conference on Machine Learning*, 2025.
- Xu, H. and Yang, Y. ITP: instance-aware test pruning for out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 21743–21751, 2025.
- Yang, Y. and Xu, H. Strengthen out-of-distribution detection capability with progressive self-knowledge distillation. In *International Conference on Machine Learning*, 2025.
- Zhang, G., Wang, L., Kang, G., Chen, L., and Wei, Y. Slca: Slow learner with classifier alignment for continual learning on a pre-trained model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19148–19158, 2023a.
- Zhang, L., Zhang, L., Shi, S., Chu, X., and Li, B. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. *CoRR*, 2023b.
- Zhou, D., Yang, Y., and Zhan, D. Learning to classify with incremental new class. *IEEE Trans. Neural Networks Learn. Syst.*, pp. 2429–2443, 2022.
- Zhou, D.-W., Sun, H.-L., Ye, H.-J., and Zhan, D.-C. Expandable subspace ensemble for pre-trained model-based class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23554–23564, 2024.

## A. Implementation Details.

We utilized a single NVIDIA RTX 4090 GPU for all experimental evaluations. For the network architecture, we adopt ViT-B/16 (Dosovitskiy et al., 2021) with  $N = 12$  transformer blocks pretrained on ImageNet-21K (Russakovsky et al., 2015) as our backbone architecture across all experiments. Consistent with prior works (Liang & Li, 2024; Liu & Chang, 2025), We set LoRA rank to  $r = 10$ , and integrate the LoRA architecture into the key and value components of all the attention layers in the transformer. The model was optimized using the Adam optimizer with a learning rate of 0.0005 and a cosine annealing scheduler, with a batch size of 64. Across different datasets: each task is trained for 50 epochs on ImageNet-R and ImageNet-A, and 20 epochs on CIFAR100 and CUB-200. Following standard protocol, we report the mean and standard deviation over three runs with different random seeds. The use of random seeds introduces variability in class order across runs, making the evaluation of model performance more challenging.

## B. Additional Experiments.

### B.1. Performance on Long Task Sequences.

To validate the generalizability of our method across all datasets, particularly in long-sequence scenarios, we conducted additional experiments on the challenging CUB-200 and ImageNet-A datasets under long-task settings (20 and 50 sessions). As shown in Table 6, our method consistently achieves the best performance across all settings, significantly outperforming other SOTA methods. Furthermore, to better demonstrate the superiority of our approach, we report the evolution curves of  $\mathcal{A}_{Last}$  throughout the entire training process in Figure 5.

Table 6. Last and average performance results on CUB-200 and ImageNet-A under the long-term setting (20 and 50 tasks). The mean and standard deviation of three trials are provided. We compare all methods using the same ViT-B/16-IN21K backbone, seeds, and class orders. Red denotes the best result and blue denotes the second-best result in each column.

Method	CUB-200				ImageNet-A			
	20		50		20		50	
	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
CODA-Prompt (Smith et al., 2023)	66.41 $\pm$ 0.81	78.10 $\pm$ 1.87	46.25 $\pm$ 0.68	63.25 $\pm$ 2.69	45.40 $\pm$ 1.39	54.55 $\pm$ 0.86	30.35 $\pm$ 0.96	41.53 $\pm$ 0.64
SLCA (Zhang et al., 2023a)	82.48 $\pm$ 0.53	90.14 $\pm$ 1.02	78.47 $\pm$ 1.85	88.38 $\pm$ 0.24	55.01 $\pm$ 2.66	63.59 $\pm$ 2.20	49.31 $\pm$ 0.93	56.72 $\pm$ 1.82
EASE (Zhou et al., 2024)	84.51 $\pm$ 1.67	91.02 $\pm$ 1.02	86.46 $\pm$ 0.14	92.38 $\pm$ 0.25	50.32 $\pm$ 2.11	61.77 $\pm$ 1.60	37.06 $\pm$ 0.50	49.80 $\pm$ 0.24
SSIAT (Tan et al., 2024)	89.06 $\pm$ 0.66	93.52 $\pm$ 0.55	83.64 $\pm$ 1.67	91.84 $\pm$ 0.47	60.52 $\pm$ 2.07	68.87 $\pm$ 2.29	51.44 $\pm$ 0.57	61.63 $\pm$ 2.61
LoRA-DRS (Liu & Chang, 2025)	87.50 $\pm$ 0.26	92.66 $\pm$ 0.57	86.97 $\pm$ 0.13	92.70 $\pm$ 0.30	55.02 $\pm$ 1.74	64.52 $\pm$ 1.16	53.28 $\pm$ 1.75	63.44 $\pm$ 2.78
MOS (Sun et al., 2025)	89.42 $\pm$ 0.22	93.66 $\pm$ 0.42	89.28 $\pm$ 0.25	93.47 $\pm$ 0.43	55.26 $\pm$ 0.92	64.53 $\pm$ 1.02	52.29 $\pm$ 1.58	62.91 $\pm$ 1.10
MACIL (Wu et al., 2025)	88.63 $\pm$ 0.61	93.52 $\pm$ 0.43	82.06 $\pm$ 0.48	91.04 $\pm$ 0.41	59.40 $\pm$ 0.76	67.79 $\pm$ 1.51	47.86 $\pm$ 1.78	59.96 $\pm$ 1.44
<b>GR-LoRA</b>	<b>89.76<math>\pm</math>0.25</b>	<b>94.08<math>\pm</math>0.55</b>	<b>89.68<math>\pm</math>0.26</b>	<b>93.94<math>\pm</math>0.53</b>	<b>62.37<math>\pm</math>0.34</b>	<b>69.30<math>\pm</math>1.08</b>	<b>59.71<math>\pm</math>0.13</b>	<b>67.23<math>\pm</math>2.16</b>

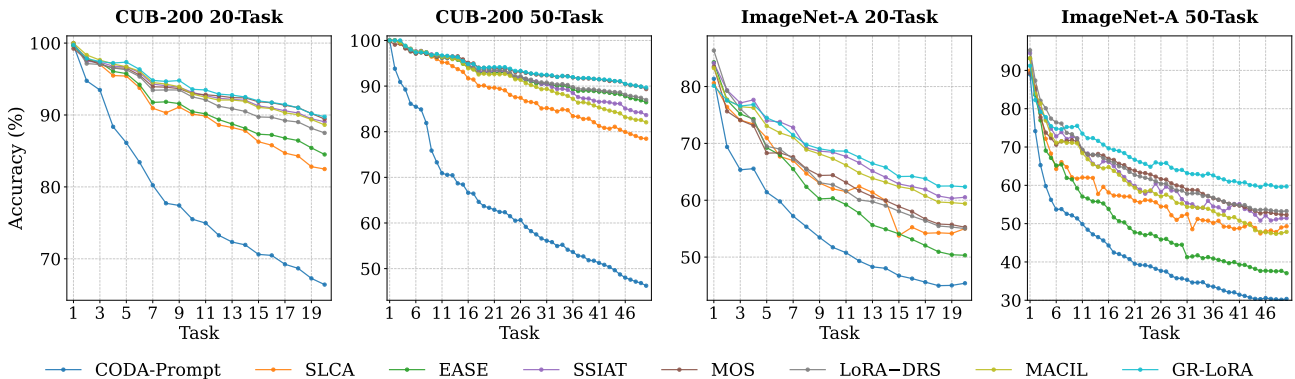


Figure 5. The performance of each learning session under different settings of CUB-200 and ImageNet-A. All methods are initialized with ViT-B/16-IN21k. These curves are plotted by calculating the average performance across three different seeds.

## B.2. Versatility on Standard Benchmarks.

To illustrate the performance evolution during the incremental learning process on standard benchmarks, we plot the accuracy curves ( $\mathcal{A}_{Last}$ ) after each task, as shown in Figure 6. It is evident that our method consistently maintains superior performance throughout the entire training sequence, occupying the leading position among all state-of-the-art baselines.

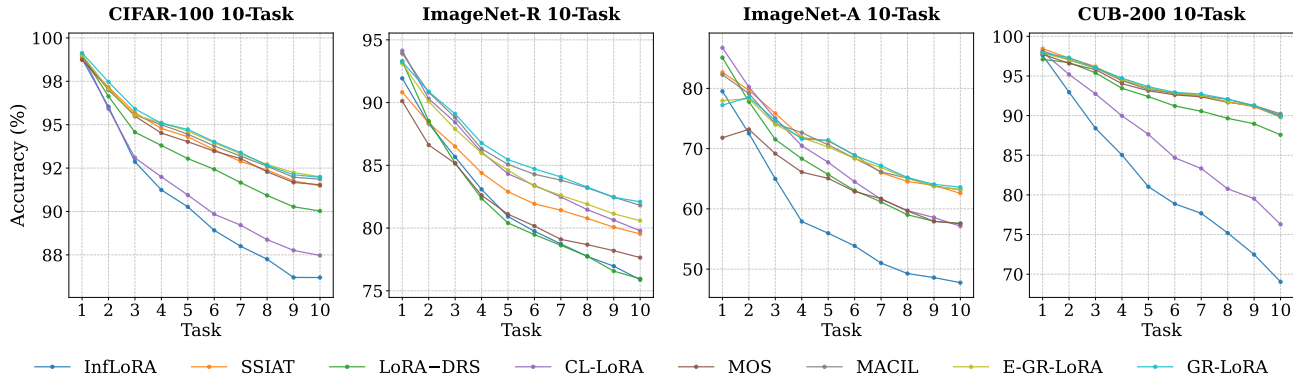


Figure 6. The performance of each learning session under four benchmark datasets (10 tasks). All methods are initialized with ViT-B/16-IN21k. These curves are plotted by calculating the average performance across three different seeds.

## B.3. Ablation Study on other Datasets.

Additional ablation studies are conducted on CIFAR-100 and ImageNet-A under the 20-task CIL setting. As shown in Table 7, progressively incorporating the proposed components consistently improves performance across both datasets. These results validate the effectiveness of each module and demonstrate the generalizability of GR-LoRA under diverse and more challenging data distributions.

Table 7. Ablation study of individual component contributions on CIFAR-100 and ImageNet-A under the 20-task CIL setting.

Method	CA	LMS	GMS	CIFAR-100		ImageNet-A	
				$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$	$\mathcal{A}_{Last}$	$\mathcal{A}_{Avg}$
LoRA-DRS	✓	×	×	90.34	93.93	55.69	64.39
Single LoRA	✓	✓	✓	88.06	92.41	58.53	64.38
GR-LoRA	✓	×	×	90.41	94.26	58.20	66.71
	✓	✓	×	90.74	94.45	59.12	67.91
	✓	×	✓	91.25	94.77	61.29	70.09
	✓	✓	✓	<b>91.42</b>	<b>94.84</b>	<b>62.08</b>	<b>70.55</b>

## B.4. Storage Overhead.

Storing class prototypes and covariance statistics is a common and effective strategy in CIL, as it helps mitigate the bias of the classifier toward newly learned tasks and improves the retention of previously acquired knowledge. Compared with standard prototype-based classifier or classifier alignment methods, GR-LoRA introduces only negligible additional storage overhead. Specifically, our method requires storing one extra shared-space prototype  $\mu_c \in \mathbb{R}^{1 \times 768}$  per class. This additional prototype is used only for constructing OOD prototypes. Importantly, GR-LoRA does not require storing any extra covariance matrix, the generated OOD samples reuse the same private-space covariance matrix  $\hat{\Sigma}_c$  as the original ID samples. As shown in Table 8, the storage cost of GR-LoRA is 2.365 MB per class, compared with 2.362 MB for standard Classify Alignment. Therefore, the additional overhead is only about 3 KB per class, which is negligible in practice.

Table 8. Storage overhead comparison of representative prototype-based CIL methods per class with feature dimension  $d = 768$ .

Method	Stored Variables	Parameters	Storage (MB)
Prototype Classifier	$\mu$	$d$	0.003
Classifier Alignment	$\mu + \Sigma$	$d + d^2$	2.362
GR-LoRA	$\mu + \hat{\mu} + \hat{\Sigma}$	$2d + d^2$	2.365

### C. Theoretical proof process.

To analyze the plasticity of the proposed method, we study the attainable population risk under different parameter constraints from an optimization perspective. We assume that the population risk  $R_t(\mathbf{W})$  is  $\mu$ -strongly convex and  $L$ -smooth with respect to the model parameters  $\mathbf{W}$ . Let  $\mathbf{P}_t^l \in \mathbb{R}^{d \times k}$  denote an orthonormal basis of the residual subspace at layer  $l$  for task  $t$ , i.e.,  $(\mathbf{P}_t^l)^\top \mathbf{P}_t^l = \mathbf{I}_k$ . The corresponding orthogonal projection operator is defined as  $\mathbf{\Pi}_t^l = \mathbf{P}_t^l (\mathbf{P}_t^l)^\top$ , and the complementary projection is  $\mathbf{Q}_t^l = \mathbf{I} - \mathbf{\Pi}_t^l$ .

**Theorem C.1.** For each task  $t = 1, \dots, T$ , let  $\widehat{\mathbf{W}}_{\perp,t}$  denote the solution obtained under the orthogonal constraint induced by the layer-wise projectors  $\{\mathbf{\Pi}_t^l\}_{l=1}^L$ , and let  $\mathbf{W}_t^*$  be the unconstrained population minimizer of  $R_t$ . Then, with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \sum_{t=1}^T \left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_t^*) \right) &\leq \sum_{t=1}^T \sum_{l=1}^L \left( \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2 \right. \\ &\quad \left. + 4\mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l) + 2\sqrt{\frac{\log(2TL/\delta)}{2n_t}} \right), \end{aligned} \quad (16)$$

where  $\mathbf{W}_{\perp,t}^*$  is the population minimizer under the orthogonal constraint and  $\mathfrak{R}_{n_t}(\cdot)$  denotes the Rademacher complexity.

*Proof.* For each task  $t$ , let the data distribution be  $\mathcal{P}_t$  over  $(\mathbf{x}, y)$ . Define the population risk and empirical risk (w.r.t. samples  $\mathcal{D}_t = \{(\mathbf{x}_{i,t}, y_{i,t})\}_{i=1}^{n_t}$ ) as

$$R_t(\mathbf{W}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}_t} [\ell(f(\mathbf{x}; \mathbf{W}), y)], \quad \widehat{R}_t(\mathbf{W}) := \frac{1}{n_t} \sum_{i=1}^{n_t} \ell(f(\mathbf{x}_{i,t}; \mathbf{W}), y_{i,t}). \quad (17)$$

Let the orthogonal constraint set be

$$\mathcal{S}_t := \left\{ \mathbf{W} : \mathbf{W}^l = \mathbf{W}_{t-1}^l + \mathbf{\Pi}_t^l \mathbf{U}^l \text{ for some } \mathbf{U}^l, \forall l \in \{1, \dots, L\} \right\}, \quad (18)$$

i.e., the layer-wise update  $\Delta \mathbf{W}^l := \mathbf{W}^l - \mathbf{W}_{t-1}^l$  is restricted to  $\text{range}(\mathbf{\Pi}_t^l)$ . Define

$$\mathbf{W}_{\perp,t}^* := \arg \min_{\mathbf{W} \in \mathcal{S}_t} R_t(\mathbf{W}), \quad \widehat{\mathbf{W}}_{\perp,t} := \arg \min_{\mathbf{W} \in \mathcal{S}_t} \widehat{R}_t(\mathbf{W}), \quad \mathbf{W}_t^* := \arg \min_{\mathbf{W}} R_t(\mathbf{W}). \quad (19)$$

Assume  $R_t(\mathbf{W})$  is  $\mu$ -strongly convex and  $L$ -smooth in  $\mathbf{W}$  (Polovinkin, 1996). For each task  $t$ ,

$$R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_t^*) = \underbrace{\left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) \right)}_{\text{estimation / generalization}} + \underbrace{\left( R_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_t^*) \right)}_{\text{approximation / constraint-induced}}. \quad (20)$$

Summing (20) over  $t = 1, \dots, T$  reduces the proof to bounding the two terms on the right-hand side. We first state a standard consequence of  $\mu$ -strong convexity.

**Lemma C.2.** If a differentiable function  $F$  is  $\mu$ -strongly convex, then for its unique minimizer  $\mathbf{W}^* = \arg \min_{\mathbf{W}} F(\mathbf{W})$ ,

$$F(\mathbf{W}) - F(\mathbf{W}^*) \leq \frac{1}{2\mu} \|\nabla F(\mathbf{W})\|_F^2, \quad \forall \mathbf{W}. \quad (21)$$

*Proof of Lemma C.2.* By  $\mu$ -strong convexity, for all  $\mathbf{U}, \mathbf{V}$ ,

$$F(\mathbf{U}) \geq F(\mathbf{V}) + \langle \nabla F(\mathbf{V}), \mathbf{U} - \mathbf{V} \rangle + \frac{\mu}{2} \|\mathbf{U} - \mathbf{V}\|_F^2. \quad (22)$$

Set  $\mathbf{U} = \mathbf{W}^*$  and  $\mathbf{V} = \mathbf{W}$ , rearrange:

$$F(\mathbf{W}) - F(\mathbf{W}^*) \leq \langle \nabla F(\mathbf{W}), \mathbf{W} - \mathbf{W}^* \rangle - \frac{\mu}{2} \|\mathbf{W} - \mathbf{W}^*\|_F^2. \quad (23)$$

Apply Cauchy–Schwarz and the inequality  $ab - \frac{\mu}{2}b^2 \leq \frac{1}{2\mu}a^2$  with  $a = \|\nabla F(\mathbf{W})\|_F$  and  $b = \|\mathbf{W} - \mathbf{W}^*\|_F$ :

$$\langle \nabla F(\mathbf{W}), \mathbf{W} - \mathbf{W}^* \rangle - \frac{\mu}{2} \|\mathbf{W} - \mathbf{W}^*\|_F^2 \leq \|\nabla F(\mathbf{W})\|_F \|\mathbf{W} - \mathbf{W}^*\|_F - \frac{\mu}{2} \|\mathbf{W} - \mathbf{W}^*\|_F^2 \leq \frac{1}{2\mu} \|\nabla F(\mathbf{W})\|_F^2. \quad (24)$$

This yields (21).  $\square$

We now bound  $R_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_t^*)$ . Because  $\mathcal{S}_t$  is an affine subspace whose feasible directions at layer  $l$  equal  $\text{range}(\mathbf{\Pi}_t^l)$ , the first-order optimality condition for the constrained minimizer  $\mathbf{W}_{\perp,t}^*$  can be written as: for every layer  $l$  and for every feasible perturbation  $\Delta \mathbf{W}^l \in \text{range}(\mathbf{\Pi}_t^l)$ ,

$$\langle \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*), \Delta \mathbf{W}^l \rangle = 0. \quad (25)$$

Since  $\text{range}(\mathbf{\Pi}_t^l)$  is exactly the set  $\{\mathbf{\Pi}_t^l \mathbf{U}^l : \mathbf{U}^l\}$ , (25) is equivalent to

$$\langle \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*), \mathbf{\Pi}_t^l \mathbf{U}^l \rangle = 0, \quad \forall \mathbf{U}^l. \quad (26)$$

Using  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$  and cyclicity of trace,

$$0 = \text{tr}((\nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*))^\top \mathbf{\Pi}_t^l \mathbf{U}^l) = \text{tr}((\mathbf{\Pi}_t^l)^\top \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*))^\top \mathbf{U}^l), \quad \forall \mathbf{U}^l. \quad (27)$$

Hence,

$$(\mathbf{\Pi}_t^l)^\top \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) = \mathbf{0}. \quad (28)$$

Because  $\mathbf{\Pi}_t^l$  is an orthogonal projector, it is symmetric:  $(\mathbf{\Pi}_t^l)^\top = \mathbf{\Pi}_t^l$ . Therefore,

$$\mathbf{\Pi}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) = \mathbf{0}, \quad \implies \quad \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) = (\mathbf{I} - \mathbf{\Pi}_t^l) \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) = \mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*). \quad (29)$$

Now apply Lemma C.2 to  $F(\mathbf{W}) = R_t(\mathbf{W})$  at  $\mathbf{W} = \mathbf{W}_{\perp,t}^*$ :

$$R_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_t^*) \leq \frac{1}{2\mu} \|\nabla R_t(\mathbf{W}_{\perp,t}^*)\|_F^2. \quad (30)$$

Finally, expand the full gradient norm layer-wise and use (29):

$$\|\nabla R_t(\mathbf{W}_{\perp,t}^*)\|_F^2 = \sum_{l=1}^L \|\nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2 = \sum_{l=1}^L \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2. \quad (31)$$

Thus we obtain the constraint-induced gap:

$$R_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_t^*) \leq \sum_{l=1}^L \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2. \quad (32)$$

We now bound  $R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*)$ . Let  $\mathcal{F}_{\perp,t}$  be the hypothesis class induced by parameters in  $\mathcal{S}_t$ :

$$\mathcal{F}_{\perp,t} := \{(\mathbf{x} \mapsto \ell(f(\mathbf{x}; \mathbf{W}), y)) : \mathbf{W} \in \mathcal{S}_t\}. \quad (33)$$

Assume we have an additive upper bound on the Rademacher complexity:

$$\mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}) \leq \sum_{l=1}^L \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l). \quad (34)$$

We use the standard uniform convergence bound: for any  $\delta_t \in (0, 1)$ , with probability at least  $1 - \delta_t$ ,

$$\sup_{f \in \mathcal{F}_{\perp,t}} \left( R_t(f) - \widehat{R}_t(f) \right) \leq 2 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}) + \sqrt{\frac{\log(1/\delta_t)}{2n_t}}. \quad (35)$$

Similarly, with probability at least  $1 - \delta_t$ ,

$$\sup_{f \in \mathcal{F}_{\perp,t}} \left( \widehat{R}_t(f) - R_t(f) \right) \leq 2 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}) + \sqrt{\frac{\log(1/\delta_t)}{2n_t}}. \quad (36)$$

By a union bound over (35) and (36), with probability at least  $1 - 2\delta_t$ ,

$$\sup_{f \in \mathcal{F}_{\perp,t}} |R_t(f) - \widehat{R}_t(f)| \leq 2 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}) + \sqrt{\frac{\log(1/\delta_t)}{2n_t}}. \quad (37)$$

Now use the ERM property:

$$\widehat{R}_t(\widehat{\mathbf{W}}_{\perp,t}) \leq \widehat{R}_t(\mathbf{W}_{\perp,t}^*). \quad (38)$$

Then, on the event (37),

$$\begin{aligned} R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) &= \left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - \widehat{R}_t(\widehat{\mathbf{W}}_{\perp,t}) \right) + \left( \widehat{R}_t(\widehat{\mathbf{W}}_{\perp,t}) - \widehat{R}_t(\mathbf{W}_{\perp,t}^*) \right) + \left( \widehat{R}_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_{\perp,t}^*) \right) \\ &\leq \left| R_t(\widehat{\mathbf{W}}_{\perp,t}) - \widehat{R}_t(\widehat{\mathbf{W}}_{\perp,t}) \right| + 0 + \left| \widehat{R}_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_{\perp,t}^*) \right| \\ &\leq 2 \sup_{\mathbf{W} \in \mathcal{S}_t} |R_t(\mathbf{W}) - \widehat{R}_t(\mathbf{W})| \\ &\leq 4 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}) + 2 \sqrt{\frac{\log(1/\delta_t)}{2n_t}}. \end{aligned} \quad (39)$$

Combining with (34), we further obtain

$$R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) \leq \sum_{l=1}^L \left( 4 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l) + 2 \sqrt{\frac{\log(1/\delta_t)}{2n_t}} \right). \quad (40)$$

Choose  $\delta_t$  so that the total failure probability across all tasks and layers is at most  $\delta$ . A convenient choice is  $\delta_t := \frac{\delta}{TL}$ , and we already paid a factor 2 in (37), hence the logarithmic term becomes  $\log(2TL/\delta)$ . Substituting  $\delta_t = \delta/(TL)$  into (40) yields, simultaneously for all  $t$  (and all layers accounted in the sum),

$$R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) \leq \sum_{l=1}^L \left( 4 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l) + 2 \sqrt{\frac{\log(2TL/\delta)}{2n_t}} \right), \quad (41)$$

with probability at least  $1 - \delta$  after union bounding over  $t = 1, \dots, T$ . Finally, plug (32) and (41) into (20), and sum over  $t = 1, \dots, T$ :

$$\begin{aligned} \sum_{t=1}^T \left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) \right) &\leq \sum_{t=1}^T \left( R_t(\widehat{\mathbf{W}}_{\perp,t}) - R_t(\mathbf{W}_{\perp,t}^*) + R_t(\mathbf{W}_{\perp,t}^*) - R_t(\mathbf{W}_{\perp,t}^*) \right) \\ &\leq \sum_{t=1}^T \sum_{l=1}^L \left( 4 \mathfrak{R}_{n_t}(\mathcal{F}_{\perp,t}^l) + 2 \sqrt{\frac{\log(2TL/\delta)}{2n_t}} + \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*)\|_F^2 \right). \end{aligned} \quad (42)$$

This is exactly the desired bound.  $\square$

As established in Theorem C.1, when model updates are restricted to the orthogonal subspace defined by  $\Pi_t^l$ , the cumulative performance gap relative to the unconstrained optimum is dominated by  $\|\mathbf{Q}_t^l \nabla R_t\|_F^2$ . Consequently, even under perfect samples and optimization, orthogonal constraints induce an unavoidable plasticity gap.

**Theorem C.3.** *Suppose the model parameters are augmented in the LoRA form  $\mathbf{W}^l = \mathbf{W}_0^l + \sum_{j=1}^{t-1} \mathbf{B}_j^l \mathbf{A}_j^l + (\mathbf{B}_t^l + \hat{\mathbf{B}}_t^l)(\mathbf{A}_t^l + \hat{\mathbf{A}}_t^l)$ , where  $\mathbf{B}_t^l \mathbf{A}_t^l$  denotes the shared module learned under orthogonal constraints, and  $\hat{\mathbf{B}}_t^l \hat{\mathbf{A}}_t^l$  is a task-specific module allocated to task  $t$ . Let  $\mathbf{W}_{\text{GR},t}^*$  denote the population minimizer under this augmented parameterization. Then the cumulative optimality gap across tasks satisfies*

$$\sum_{t=1}^T \left( R_t(\mathbf{W}_{\text{GR},t}^*) - R_t(\mathbf{W}_t^*) \right) \leq \sum_{t=1}^T \sum_{l=1}^L \frac{1}{2\mu} \left\| \mathbf{Q}_t^l \nabla_{\mathbf{w}^l} R_t(\mathbf{W}_{\perp,t}^*) - \Pi_{S_t^l}(\mathbf{Q}_t^l \nabla_{\mathbf{w}^l} R_t(\mathbf{W}_{\perp,t}^*)) \right\|_F^2, \quad (43)$$

where  $\Pi_{S_t^l}$  denotes the projection onto the subspace  $S_t^l = \text{range}(\hat{\mathbf{B}}_t^l \hat{\mathbf{A}}_t^l)$ .

*Proof.* For each layer  $l$ , define  $\mathbf{w}^l := \text{vec}(\mathbf{W}^l) \in \mathbb{R}^{d_l}$  and stack all layers as  $\mathbf{w} := ((\mathbf{w}^1)^\top, \dots, (\mathbf{w}^L)^\top)^\top \in \mathbb{R}^d$ , where  $d = \sum_{l=1}^L d_l$ . We use the standard inner product  $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^\top \mathbf{b}$  and its matrix form  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$ . For any operator  $\mathbf{M}$ , denote its orthogonal projector by the same symbol. Under the quadratic assumption, for each task  $t$  there exists a symmetric positive definite matrix  $\mathbf{H}_t \in \mathbb{R}^{d \times d}$  such that

$$R_t(\mathbf{w}) = R_t(\mathbf{w}_t^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_t^*)^\top \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t^*), \quad \nabla R_t(\mathbf{w}) = \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t^*), \quad (44)$$

with  $\mu \mathbf{I} \preceq \mathbf{H}_t \preceq L \mathbf{I}$ . In particular,  $\mathbf{w}_t^*$  is the unique unconstrained minimizer and  $\nabla R_t(\mathbf{w}_t^*) = \mathbf{0}$ . Let  $\mathcal{T}_{\perp,t}$  denote the tangent subspace induced by the layer-wise orthogonal projectors  $\{\Pi_t^l\}_{l=1}^L$ :

$$\mathcal{T}_{\perp,t} := \left\{ \Delta \mathbf{w} = ((\Delta \mathbf{w}^1)^\top, \dots, (\Delta \mathbf{w}^L)^\top)^\top : \Delta \mathbf{w}^l \in \text{range}(\Pi_t^l), \forall l \right\}. \quad (45)$$

Define the corresponding block-diagonal projector  $\Pi_{\perp,t}$  onto  $\mathcal{T}_{\perp,t}$  and  $\mathbf{Q}_{\perp,t} := \mathbf{I} - \Pi_{\perp,t}$ . By construction, on each layer block  $l$ ,  $(\Pi_{\perp,t})|_l = \Pi_t^l$  and  $(\mathbf{Q}_{\perp,t})|_l = \mathbf{Q}_t^l$ . Let  $\mathbf{w}_{\perp,t}^*$  be the constrained minimizer:

$$\mathbf{w}_{\perp,t}^* := \arg \min_{\mathbf{w} \in \mathbf{w}_{t-1} + \mathcal{T}_{\perp,t}} R_t(\mathbf{w}). \quad (46)$$

Since the feasible set is an affine subspace and  $R_t$  is differentiable and strongly convex,  $\mathbf{w}_{\perp,t}^*$  is unique and satisfies the first-order optimality condition:

$$\Pi_{\perp,t} \nabla R_t(\mathbf{w}_{\perp,t}^*) = \mathbf{0} \quad \iff \quad \nabla R_t(\mathbf{w}_{\perp,t}^*) = \mathbf{Q}_{\perp,t} \nabla R_t(\mathbf{w}_{\perp,t}^*). \quad (47)$$

Define the *discarded gradient* at  $\mathbf{w}_{\perp,t}^*$ :

$$\mathbf{g}_t := \mathbf{Q}_{\perp,t} \nabla R_t(\mathbf{w}_{\perp,t}^*) = \nabla R_t(\mathbf{w}_{\perp,t}^*), \quad \mathbf{g}_t^l := \mathbf{Q}_t^l \nabla_{\mathbf{w}^l} R_t(\mathbf{w}_{\perp,t}^*), \quad (48)$$

so that  $\|\mathbf{g}_t\|_2^2 = \sum_{l=1}^L \|\mathbf{g}_t^l\|_2^2$ . Under the GR-LoRA, the update directions include: (i) the orthogonal directions  $\mathcal{T}_{\perp,t}$  (shared orthogonal module), and (ii) additional task-specific directions in each layer  $l$  belonging to the subspace

$$S_t^l := \text{range}(\hat{\mathbf{B}}_t^l \hat{\mathbf{A}}_t^l) \subseteq \text{range}(\mathbf{Q}_t^l), \quad (49)$$

where the inclusion in (49) is the standard GR-LoRA design assumption (task-specific module spans residual directions). Let  $\mathcal{T}_{\text{GR},t}$  denote the expanded tangent space:  $\mathcal{T}_{\text{GR},t} := \left\{ \Delta \mathbf{w} : \Delta \mathbf{w}^l \in \text{range}(\Pi_t^l) \oplus S_t^l, \forall l \right\}$ . Let  $\Pi_{\text{GR},t}$  be the orthogonal projector onto  $\mathcal{T}_{\text{GR},t}$ , and  $\mathbf{Q}_{\text{GR},t} := \mathbf{I} - \Pi_{\text{GR},t}$ . Let  $\mathbf{w}_{\text{GR},t}^*$  be the population minimizer over the GR feasible affine space:

$$\mathbf{w}_{\text{GR},t}^* := \arg \min_{\mathbf{w} \in \mathbf{w}_{t-1} + \mathcal{T}_{\text{GR},t}} R_t(\mathbf{w}). \quad (50)$$

Again, uniqueness holds and the first-order optimality condition yields

$$\Pi_{\text{GR},t} \nabla R_t(\mathbf{w}_{\text{GR},t}^*) = \mathbf{0} \quad \iff \quad \nabla R_t(\mathbf{w}_{\text{GR},t}^*) = \mathbf{Q}_{\text{GR},t} \nabla R_t(\mathbf{w}_{\text{GR},t}^*). \quad (51)$$

For a quadratic  $R_t$  of the form (44), minimizing  $R_t$  over an affine subspace  $\mathbf{w}_0 + \mathcal{T}$  is equivalent to projecting the error  $\mathbf{w}_t^* - \mathbf{w}_0$  in the  $\mathbf{H}_t$ -geometry. A standard characterization is:

$$\mathbf{w}_{\mathcal{T}}^* = \arg \min_{\mathbf{w} \in \mathbf{w}_0 + \mathcal{T}} R_t(\mathbf{w}) \iff \mathbf{w}_{\mathcal{T}}^* = \mathbf{w}_0 + \mathbf{\Pi}_{\mathcal{T}}^{(\mathbf{H}_t)}(\mathbf{w}_t^* - \mathbf{w}_0), \quad (52)$$

where  $\mathbf{\Pi}_{\mathcal{T}}^{(\mathbf{H}_t)}$  denotes the  $\mathbf{H}_t$ -orthogonal projection onto  $\mathcal{T}$  (i.e., projection under the inner product  $\langle \mathbf{a}, \mathbf{b} \rangle_{\mathbf{H}_t} = \mathbf{a}^\top \mathbf{H}_t \mathbf{b}$ ). Moreover, the optimality gap admits the exact expression

$$R_t(\mathbf{w}_{\mathcal{T}}^*) - R_t(\mathbf{w}_t^*) = \frac{1}{2} \|\mathbf{Q}_{\mathcal{T}}^{(\mathbf{H}_t)}(\mathbf{w}_0 - \mathbf{w}_t^*)\|_{\mathbf{H}_t}^2, \quad (53)$$

where  $\mathbf{Q}_{\mathcal{T}}^{(\mathbf{H}_t)} = \mathbf{I} - \mathbf{\Pi}_{\mathcal{T}}^{(\mathbf{H}_t)}$  and  $\|\mathbf{v}\|_{\mathbf{H}_t}^2 := \mathbf{v}^\top \mathbf{H}_t \mathbf{v}$ . We will use a *Euclidean* upper bound of (53) via the spectral lower bound  $\mathbf{H}_t \succeq \mu \mathbf{I}$ :

$$\|\mathbf{v}\|_{\mathbf{H}_t}^2 \leq \frac{1}{\mu} \|\mathbf{H}_t \mathbf{v}\|_2^2 \quad \text{since} \quad \mathbf{v}^\top \mathbf{H}_t \mathbf{v} \leq \frac{1}{\mu} \mathbf{v}^\top \mathbf{H}_t^2 \mathbf{v} = \frac{1}{\mu} \|\mathbf{H}_t \mathbf{v}\|_2^2. \quad (54)$$

Apply (53) to  $\mathcal{T} = \mathcal{T}_{\text{GR},t}$  and  $\mathbf{w}_0 = \mathbf{w}_{t-1}$ :

$$R_t(\mathbf{w}_{\text{GR},t}^*) - R_t(\mathbf{w}_t^*) = \frac{1}{2} \|\mathbf{Q}_{\text{GR},t}^{(\mathbf{H}_t)}(\mathbf{w}_{t-1} - \mathbf{w}_t^*)\|_{\mathbf{H}_t}^2. \quad (55)$$

Now define the constrained-only solution  $\mathbf{w}_{\perp,t}^*$  which corresponds to tangent space  $\mathcal{T}_{\perp,t}$ . By the same formula,

$$R_t(\mathbf{w}_{\perp,t}^*) - R_t(\mathbf{w}_t^*) = \frac{1}{2} \|\mathbf{Q}_{\perp,t}^{(\mathbf{H}_t)}(\mathbf{w}_{t-1} - \mathbf{w}_t^*)\|_{\mathbf{H}_t}^2. \quad (56)$$

Crucially, the *residual gradient* at  $\mathbf{w}_{\perp,t}^*$  satisfies

$$\mathbf{g}_t = \nabla R_t(\mathbf{w}_{\perp,t}^*) = \mathbf{H}_t(\mathbf{w}_{\perp,t}^* - \mathbf{w}_t^*). \quad (57)$$

Because  $\mathbf{w}_{\perp,t}^* \in \mathbf{w}_{t-1} + \mathcal{T}_{\perp,t}$ , we can write

$$\mathbf{w}_{\perp,t}^* - \mathbf{w}_t^* = (\mathbf{w}_{t-1} - \mathbf{w}_t^*) + \Delta_{\perp,t} \quad \text{for some} \quad \Delta_{\perp,t} \in \mathcal{T}_{\perp,t}. \quad (58)$$

Projecting onto the orthogonal complement of  $\mathcal{T}_{\text{GR},t}$  eliminates the  $\mathcal{T}_{\perp,t}$  component since  $\mathcal{T}_{\perp,t} \subseteq \mathcal{T}_{\text{GR},t}$ :

$$\mathbf{Q}_{\text{GR},t}(\mathbf{w}_{\perp,t}^* - \mathbf{w}_t^*) = \mathbf{Q}_{\text{GR},t}(\mathbf{w}_{t-1} - \mathbf{w}_t^*). \quad (59)$$

Now use (55) and the bound (54) with  $\mathbf{v} = \mathbf{Q}_{\text{GR},t}(\mathbf{w}_{t-1} - \mathbf{w}_t^*)$ :

$$\begin{aligned} R_t(\mathbf{w}_{\text{GR},t}^*) - R_t(\mathbf{w}_t^*) &= \frac{1}{2} \|\mathbf{Q}_{\text{GR},t}^{(\mathbf{H}_t)}(\mathbf{w}_{t-1} - \mathbf{w}_t^*)\|_{\mathbf{H}_t}^2 \leq \frac{1}{2} \|\mathbf{Q}_{\text{GR},t}(\mathbf{w}_{t-1} - \mathbf{w}_t^*)\|_{\mathbf{H}_t}^2 \\ &\leq \frac{1}{2\mu} \left\| \mathbf{H}_t \mathbf{Q}_{\text{GR},t}(\mathbf{w}_{t-1} - \mathbf{w}_t^*) \right\|_2^2. \end{aligned} \quad (60)$$

Using (59) and (57),

$$\mathbf{H}_t \mathbf{Q}_{\text{GR},t}(\mathbf{w}_{t-1} - \mathbf{w}_t^*) = \mathbf{H}_t \mathbf{Q}_{\text{GR},t}(\mathbf{w}_{\perp,t}^* - \mathbf{w}_t^*) = \mathbf{Q}_{\text{GR},t} \mathbf{H}_t(\mathbf{w}_{\perp,t}^* - \mathbf{w}_t^*) = \mathbf{Q}_{\text{GR},t} \mathbf{g}_t, \quad (61)$$

where the last equality uses that  $\mathbf{Q}_{\text{GR},t}$  is an *Euclidean* orthogonal projector and  $\mathbf{g}_t$  is a vector; thus  $\mathbf{Q}_{\text{GR},t} \mathbf{g}_t$  is well-defined. Plugging (61) into (60) yields

$$R_t(\mathbf{w}_{\text{GR},t}^*) - R_t(\mathbf{w}_t^*) \leq \frac{1}{2\mu} \|\mathbf{Q}_{\text{GR},t} \mathbf{g}_t\|_2^2. \quad (62)$$

By construction,  $\mathcal{T}_{\text{GR},t}$  expands  $\mathcal{T}_{\perp,t}$  by adding  $S_t^l$  in each layer. Under the design assumption  $S_t^l \subseteq \text{range}(\mathbf{Q}_t^l)$  in (49), we have an *orthogonal direct sum* per layer:

$$\text{range}(\mathbf{\Pi}_t^l) \perp S_t^l, \quad \text{range}(\mathbf{\Pi}_t^l) \oplus S_t^l \subseteq \mathbb{R}^{d_t}. \quad (63)$$

Hence, the Euclidean orthogonal projector onto  $\text{range}(\mathbf{\Pi}_t^l) \oplus \mathcal{S}_t^l$  is

$$\mathbf{\Pi}_{\text{GR},t}^l = \mathbf{\Pi}_t^l + \Pi_{\mathcal{S}_t^l}, \quad \mathbf{Q}_{\text{GR},t}^l = \mathbf{I} - \mathbf{\Pi}_{\text{GR},t}^l = \mathbf{I} - \mathbf{\Pi}_t^l - \Pi_{\mathcal{S}_t^l}. \quad (64)$$

Now recall from (48) that  $\mathbf{g}_t^l \in \text{range}(\mathbf{Q}_t^l)$ , i.e.,  $\mathbf{\Pi}_t^l \mathbf{g}_t^l = \mathbf{0}$ . Therefore,

$$\mathbf{Q}_{\text{GR},t}^l \mathbf{g}_t^l = (\mathbf{I} - \mathbf{\Pi}_t^l - \Pi_{\mathcal{S}_t^l}) \mathbf{g}_t^l = \mathbf{g}_t^l - \Pi_{\mathcal{S}_t^l}(\mathbf{g}_t^l). \quad (65)$$

Stacking all layers,

$$\|\mathbf{Q}_{\text{GR},t} \mathbf{g}_t\|_2^2 = \sum_{l=1}^L \|\mathbf{Q}_{\text{GR},t}^l \mathbf{g}_t^l\|_2^2 = \sum_{l=1}^L \|\mathbf{g}_t^l - \Pi_{\mathcal{S}_t^l}(\mathbf{g}_t^l)\|_2^2. \quad (66)$$

Returning to matrix notation,  $\|\cdot\|_2$  on  $\text{vec}(\cdot)$  equals the Frobenius norm  $\|\cdot\|_F$  on matrices, so

$$\|\mathbf{g}_t^l - \Pi_{\mathcal{S}_t^l}(\mathbf{g}_t^l)\|_2^2 = \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) - \Pi_{\mathcal{S}_t^l}(\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*))\|_F^2. \quad (67)$$

Combine (62) and (66):

$$R_t(\mathbf{W}_{\text{GR},t}^*) - R_t(\mathbf{W}_t^*) \leq \sum_{l=1}^L \frac{1}{2\mu} \|\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*) - \Pi_{\mathcal{S}_t^l}(\mathbf{Q}_t^l \nabla_{\mathbf{W}^l} R_t(\mathbf{W}_{\perp,t}^*))\|_F^2. \quad (68)$$

Summing the above inequality over  $t = 1, \dots, T$  yields the claimed result.  $\square$

As established in Theorem C.3, by introducing task-specific low-rank modules to recycle residual gradient components, the plasticity gap is no longer governed by the entire discarded gradient  $\mathbf{Q}_t^l \nabla R_t$ , but only by the portion that lies outside the representational subspace  $\mathcal{S}_t^l$ . This result demonstrates that Gradient Recycling effectively restores the usable optimization geometry under orthogonal constraints, thereby mitigating the intrinsic plasticity gap induced by subspace projection.