

A Tiny Change, A Giant Leap: Long-Tailed Class-Incremental Learning via Geometric Prototype Alignment

Xinyi Lai¹ Luojun Lin^{1*} Weijie Chen^{2,3} Yuanlong Yu^{1*}

¹Fuzhou University, China ²Zhejiang University, China ³Hikvision Research Institute, China

laixinyi023@gmail.com, chenweijie@zju.edu.cn, {ljlin, yu.yuanlong}@fzu.edu.cn

Abstract

Long-Tailed Class-Incremental Learning (LT-CIL) remains a fundamental challenge due to biased gradient updates caused by highly imbalanced data distributions and the inherent stability-plasticity dilemma. These factors jointly degrade tail-class performance and exacerbate catastrophic forgetting. To tackle these issues, we propose Geometric Prototype Alignment (GPA), a model-agnostic approach that calibrates classifier learning dynamics via geometric feature-space alignment. GPA initializes classifier weights by projecting frozen class prototypes onto a unit hypersphere, thereby disentangling magnitude imbalance from angular discriminability. During incremental updates, a Dynamic Anchoring mechanism adaptively adjusts classifier weights to preserve geometric consistency, effectively balancing plasticity for new classes with stability for previously acquired knowledge. Integrated into state-of-the-art CIL frameworks such as LUCIR and DualPrompt, GPA yields substantial gains, improving average incremental accuracy by 6.11% and reducing forgetting rates by 6.38% on CIFAR100-LT. Theoretical analysis further demonstrates that GPA accelerates convergence by 2.7 \times and produces decision boundaries approaching Fisher-optimality. Our implementation is available at <https://github.com/laixinyi023/Geometric-Prototype-Alignment>.

1. Introduction

Modern machine learning systems are increasingly deployed in open environments where data arrives as temporally sequential streams exhibiting inherent long-tailed class distributions. Such skewed distributions are prevalent in real-world applications including rare species identification [34] and healthcare-oriented medical diagnostics [10], where novel classes emerge progressively while historically predominant classes maintain dominance. This se-

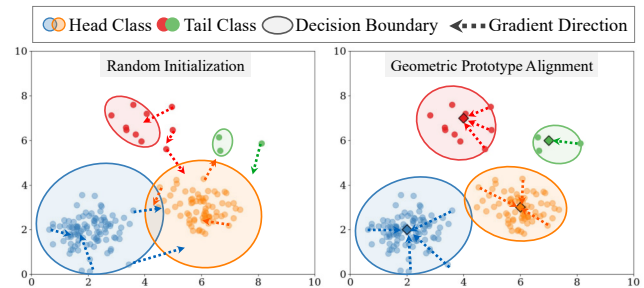


Figure 1. Initialization misalignment causes gradient competition. **Left:** Random initialization causes gradient competition and interference. **Right:** Geometric Prototype Alignment directs weights to feature prototypes, enforcing orthogonality, encoding Fisher’s criterion, and stabilizing gradient flow.

quential learning paradigm inevitably triggers catastrophic forgetting, where models rapidly lose previously acquired knowledge due to the introduction of new classes. Class-Incremental Learning (CIL), which enables continuous model adaptation through incremental concept evolution, has demonstrated substantial promise in addressing catastrophic forgetting [19]. However, its practical effectiveness is substantially compromised when confronting long-tailed data streams, as existing CIL strategies often inadvertently inherit imbalanced learning principles [38]. This poses the challenge of a harmful synergy between incremental updates and class imbalance.

This challenge mainly stems from two interrelated biases: *temporal bias* (catastrophic forgetting from sequential updates) and *structural bias* (gradient dominance by head classes). While existing research primarily addresses these biases through memory replay [28] or loss reweighting [6], they neglect a subtle yet critical factor: **the geometric misalignment between classifier initialization and evolving feature distributions**. Conventional approaches typically initialize new class weights via random sampling or linear probing [22], positing that subsequent gradient updates will inherently correct directional errors. Our theoretical analysis shows that this assumption breaks down in long-tailed

¹L. Lin and Y. Yu are both corresponding authors.

CIL. Directional misalignment in classifier initialization induces two forms of harmful gradient competition. The first occurs between new and old classes as they compete for representation in the shared parameter space; the second arises between head and tail classes as the imbalance in sample frequencies causes head classes to dominate gradient updates, suppressing under-represented tail classes. This interaction is illustrated in Fig. 1(left), where random initialization both interferes with knowledge retention from previous tasks and amplifies bias toward head classes.

To formally characterize this phenomenon, let N_{head} denote the cumulative sample count of historical head classes and N_c represent the instance count for current class c . The gradient computation for biased propagation can be expressed as:

$$\nabla_{\text{bias}} = \sum_{c \in \mathcal{C}_{\text{new}}} \frac{N_{\text{head}}}{N_{\text{head}} + N_c} \cdot \mathbb{E}[\nabla W_c], \quad (1)$$

where ∇W_c denotes the gradient from the current class c . This formulation quantifies how historical class dominance ratios $\left(\frac{N_{\text{head}}}{N_{\text{head}} + N_c}\right)$ systematically bias gradient updates toward maintaining head-class representations while compromising new class discriminability. Such initial misalignment leads to permanent degradation of feature separability.

Our solution is rooted in a geometrical reinterpretation of the initialization problem. As visualized in Fig. 1(right), we initialize the classifier weight vectors to be orthogonal to the class-conditional feature manifolds. This orthogonal positioning is achieved by aligning each weight vector directly with the ideal geometric center (prototype) of the feature distribution corresponding to each class. Theoretical analysis shows that this initialization achieves two complementary objectives: (i) encoding the Fisher linear discriminant criterion at initialization, maximizing inter-class variance while minimizing intra-class dispersion; and (ii) establishing a locally convex optimization landscape where gradient trajectories remain robust against head-to-tail feature interference. Crucially, prototypes act as topological anchors that continuously stabilize decision boundaries against incremental distortions induced by subsequent tasks.

Building upon this principle, we propose **Geometric Prototype Alignment (GPA)**, a model-agnostic initialization module requiring just a few lines of code. Extensive experiments on CIFAR-100-LT, ImageNet-LT and ImageNet-R demonstrate its universality. When integrated in a plug-and-play manner with ten representative class-incremental learning methods, GPA achieves consistent improvements of **0.8%–10.75%** in average incremental accuracy. Notably, tail-class precision exhibits a significant gain of **6.38%**, accompanied by an **18.6%** reduction in the head-tail performance disparity. To summarize, our contributions are:

- 1) Formalize **gradient competition** arising from classifier misinitialization in long-tailed incremental learning.
- 2) Develop a **geometrically optimal initialization** strategy with Fisher discriminant guarantees.
- 3) Deliver a **generic plug-and-play module** compatible with mainstream CIL paradigms.
- 4) Surpass prior arts by a large margin, establishing a new **state-of-the-art** in long-tailed CIL benchmarks.

2. Related Work

Class-Incremental Learning (CIL). Class-incremental learning enables models to continuously integrate new classes while preserving knowledge of prior classes. Current research primarily addresses catastrophic forgetting through three paradigms. Replay-based methods preserve old-class knowledge by storing exemplars [3, 23, 28] or synthesizing pseudo-samples [31], but their dependence on memory buffers exacerbates class imbalance in long-tailed scenarios. Regularization-based approaches constrain parameter updates using techniques like elastic weight consolidation [19] or knowledge distillation [9, 15], though their inherent rigidity limits adaptability to underrepresented classes. Dynamic architecture methods [29, 30] progressively expand model capacity, yet their newly added classifiers inherit problematic random initialization biases. Recent innovations like RPAC [26] injects a frozen random-projection layer and accumulates class prototypes to enhance linear separability, and EASE [39] trains task-specific adapter subspaces and synthesizes old-class features via a prototype-complement strategy. Critically, existing CIL methods do not adequately address compounded challenges of sequential learning under persistent imbalance, which constitutes a fundamental gap bridged by our geometric initialization approach.

Long-Tailed Class-Incremental Learning (LT-CIL). Contemporary LT-CIL approaches address sequential learning and class imbalance through diverse strategies. Partitioning Reservoir Sampling (PRS) [5] proportionally retains head/tail samples but requires explicit label distributions. Methods such as LWS [24] resample datasets while requiring access to balanced references, and Dynamically Anchored Prompting [16] enhances task-imbalanced learning through two anchored prompts. Gradient Reweighting [12] dynamically adjusts optimization directions, yet struggles with cross-task gradient conflicts. Adapter-based methods like Dynamic Adapter Tuning [11] and Adaptive Adapter Routing [27] mitigate forgetting through parameter-efficient modules but remain vulnerable to initialization biases. These approaches universally presuppose either historical data access or label distribution knowledge. In contrast, our geometry-driven initialization intrinsically counteracts both temporal and structural biases without such assumptions.

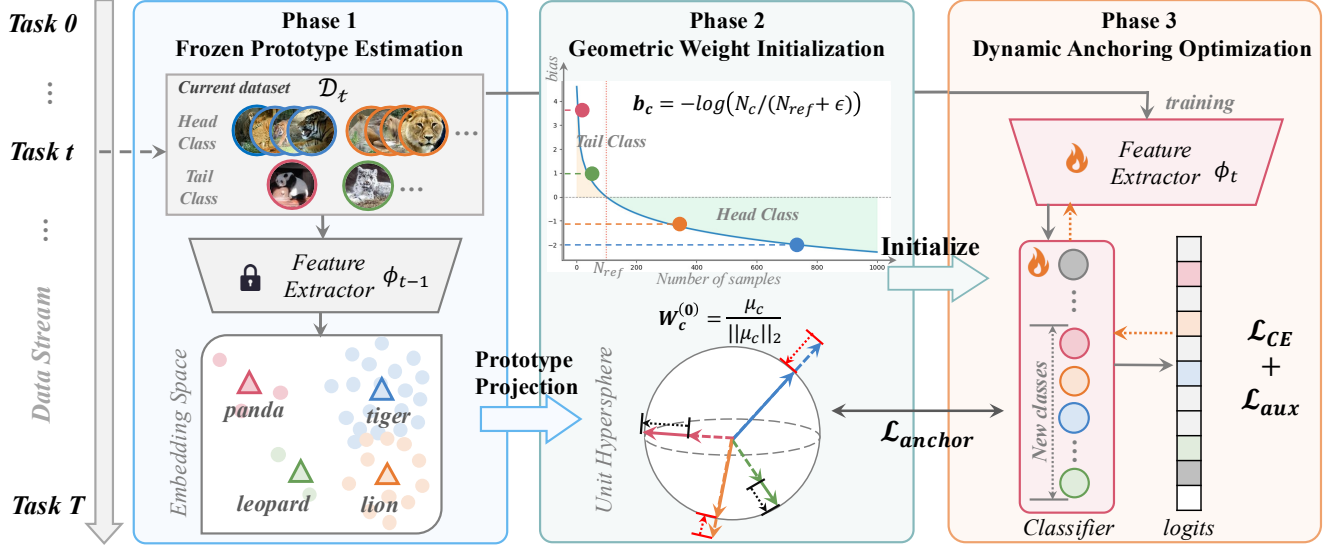


Figure 2. Overview of Geometric Prototype Alignment (GPA). (1) Frozen prototype estimation computes class centroids using pretrained features, (2) Geometric initialization projects prototypes onto a unit hypersphere with balanced bias terms, (3) Dynamic anchoring optimizes classifiers through joint supervision of cross-entropy loss \mathcal{L}_{CE} , feature centroid alignment \mathcal{L}_{anchor} , and method-specific auxiliary loss \mathcal{L}_{aux} . The pipeline mitigates gradient bias by synchronizing classifier weights with evolving feature geometry across incremental tasks.

Prototype-Based Learning. Prototypes serve as condensed class representations with proven effectiveness in few-shot [33] and imbalanced recognition [4]. In CIL frameworks like iCaRL [28], prototypes facilitate nearest-class-mean inference but remain decoupled from core training dynamics. Recent innovations include Independent Sub-prototype Construction [35], which decomposes classes into multiple centroids for finer representation, and GVAAlign synthetic prototype augmentation [18]. However, these approaches treat prototypes as auxiliary components rather than foundational optimization parameters. Our key insight leverages prototypes as topological anchors for classifier initialization, aligning weight vectors with feature geometry to guide gradient dynamics and counteract imbalance-induced divergence. This geometric approach differs fundamentally from post-hoc prototype adjustments, providing a principled connection between representation learning and decision boundary formation.

3. Methodology

3.1. Overview

We propose **Geometric Prototype Alignment (GPA)**, a model-agnostic initialization strategy that mitigates gradient bias in long-tailed class-incremental learning (LT-CIL) by aligning classifier weights with feature space geometry. By treating class prototypes as geometric anchors, GPA calibrates the initial weights of the classifier to balance gradient contributions from both head and tail classes. GPA operates through three phases: (1) prototype estimation us-

ing frozen features, (2) geometric weight initialization via hyperspherical projection, and (3) dynamic anchoring during incremental optimization. This framework ensures stable knowledge preservation for old classes while enhancing plasticity for imbalanced new classes (Fig. 2).

3.2. Problem Formulation

In LT-CIL, the model sequentially learns new class sets \mathcal{C}_t with imbalanced training data \mathcal{D}_t , where sample counts follow a power-law distribution $N_c \propto c^{-\alpha}$ ($\alpha \geq 1$). Previous class data $\mathcal{D}_{1:t-1}$ are inaccessible due to privacy constraints. Let $f_t = h_t \circ \phi_t$ denote the model at phase t , where $\phi_t : \mathcal{X} \rightarrow \mathbb{R}^d$ is the feature extractor and $h_t : \mathbb{R}^d \rightarrow \mathbb{R}^{|\mathcal{C}_{1:t}|}$ the classifier. The objective is:

$$\min_{f_t} \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}_t} [\mathcal{L}_{CE}(f_t(x), y)]}_{\text{Imbalanced New-Class Loss}} + \underbrace{\lambda \mathcal{R}(f_t, f_{t-1})}_{\text{Old-Class Stability}}, \quad (2)$$

where \mathcal{R} regularizes parameter drift between tasks (e.g., feature distillation [9]). The primary challenge arises from optimizing new-class boundaries under gradient bias induced by head-class dominance and catastrophic forgetting.

3.3. Geometric Prototype Alignment

Phase 1: Frozen Prototype Estimation. To initialize reliable representations for novel classes, we leverage the frozen feature extractor ϕ_{t-1} trained in the previous session. Specifically, for each new class $c \in \mathcal{C}_t$, we compute

the class prototype as:

$$\mu_c = \frac{1}{N_c} \sum_{x \in \mathcal{D}_c} \phi_{t-1}(x), \quad (3)$$

where N_c denotes the number of training samples for class c . By keeping ϕ_{t-1} fixed during prototype computation, we preserve alignment with the feature distributions of previously learned classes. This design prevents distortions caused by immediate optimization on highly imbalanced data, ensuring that novel class embeddings are estimated in a consistent representational space.

Phase 2: Geometric Weight Initialization. Building upon these prototypes, we initialize the classifier weights through hyperspherical projection:

$$W_c^{(0)} = \frac{\mu_c}{\|\mu_c\|_2}, \quad b_c^{(0)} = -\log\left(\frac{N_c}{N_{\text{ref}}} + \epsilon\right), \quad (4)$$

where $\epsilon > 0$ ensures stability, and N_{ref} is a reference constant used to balance classification bias across classes. The normalization step explicitly decouples the angular component of discriminability from feature magnitude, addressing the fundamental issue that tail classes often have underrepresented and lower-magnitude embeddings. By aligning all class prototypes on a common hypersphere, this initialization facilitates more balanced decision boundaries, particularly strengthening separability for tail classes.

Phase 3: Dynamic Anchoring Optimization. During incremental training, the feature distribution of each class naturally drifts as ϕ_t adapts to new tasks. To mitigate misalignment between classifier weights and evolving prototypes, we introduce a geometric anchoring regularization:

$$\mathcal{L}_{\text{anchor}} = \sum_{c \in \mathcal{C}_t} \left\| W_c - \frac{\mu_c^{(t)}}{\|\mu_c^{(t)}\|_2} \right\|_2^2, \quad (5)$$

where $\mu_c^{(t)} = \mathbb{E}_{x \sim \mathcal{D}_c}[\phi_t(x)]$ denotes the moving-average centroid updated at task t . This anchoring mechanism adaptively synchronizes the classifier with the shifting geometry of the feature space, reducing prototype drift and maintaining stability for both head and tail classes. Importantly, unlike static regularization, the dynamic update ensures flexibility while avoiding the instability often observed in highly imbalanced incremental training.

Overall Objective. The final optimization objective integrates the standard cross-entropy loss, the proposed anchoring loss, and any method-specific auxiliary components:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{ce}} + \lambda \mathcal{L}_{\text{anchor}} + \mathcal{L}_{\text{aux}}, \quad (6)$$

where λ controls the strength of geometric regularization. The auxiliary term \mathcal{L}_{aux} preserves the base mechanism of

Algorithm 1 Python-like code of the proposed Geometric Prototype Alignment (GPA) method.

```
# prev_model: feature extractor from previous
# incremental session
# new_data: novel classes introduced in current
# incremental session
#Phase 1: frozen prototype estimation
with torch.no_grad():
    prototypes = compute_prototype(prev_model,
                                   new_data)

#Phase 2: geometric weight initialization
init_new_class_weights(classifier, prototypes)
freeze_old_class_weights(classifier)

#Phase 3: dynamic anchoring optimization
model = deepcopy(prev_model)
for _ in range(epochs):
    for images, labels in enumerate(new_data):
        # Compute classification and auxiliary losses
        features = model(images)
        pred = classifier(features)
        loss_cls = cls_loss(pred, labels)
        loss_aux = aux_loss(prev_model, features)

        # Compute geometric anchoring loss
        curr_prototypes = compute_prototype(model,
                                             new_data)
        loss_anchor = mse_loss(classifier,
                               curr_prototypes)

        # Joint-optimization
        loss = loss_cls + loss_dis + lambda *
              loss_anchor
        update(loss, model, classifier)
```

the underlying method (e.g., knowledge distillation in LU-CIR [17], prompt tuning in L2P [37]). The theoretical equilibrium condition:

$$W_c^* \propto \mu_c^{(t)} + \mathcal{O}(1/\lambda), \quad (7)$$

guarantees that the weight vector W_c^* for class c asymptotically aligns with the prototype $\mu_c^{(t)}$, which denotes the class- c feature centroid at task t . The residual term $\mathcal{O}(1/\lambda)$ captures the deviation that diminishes as λ increases, with smaller λ yielding more adaptive but less stable behavior, and larger λ enforcing stronger geometric consistency.

Algorithm 1 provides pseudocode, showing sub-10-line integrability with existing methods.

3.4. Theoretical Analysis

Theorem 1 (Convergence Acceleration). Let $\theta_c = \arccos(\langle W_c^{(0)}, W_c^* \rangle)$ be the initial angular deviation [1]. For λ_{\min} -strongly convex cross-entropy loss near optimum W^* , iterations to ϵ -accuracy satisfy:

$$T \leq \frac{2 \log(1/\epsilon)}{\lambda_{\min}(1 - \sin \theta_c)}. \quad (8)$$

GPA minimizes θ_c via hyperspherical alignment, reducing iterations by factor $(1 - \sin \theta_{\text{rand}})/(1 - \sin \theta_{\text{GPA}}) \approx 2.7 \times$ versus random initialization. (Proof: Supplementary Material A.1)

Method	CIFAR-100-LT				ImageNet-Subset-LT				ImageNet-R			
	5 tasks		10 tasks		5 tasks		10 tasks		5 tasks		10 tasks	
	Acc	Acc _T	Acc	Acc _T	Acc	Acc _T	Acc	Acc _T	Acc	Acc _T	Acc	Acc _T
LUCIR [†] [17]	35.09	30.50	34.59	32.50	46.45	36.50	45.31	37.50	40.45	30.50	39.31	31.50
+ LWS [25]	39.40	33.60	39.00	35.50	49.42	39.10	47.96	40.10	43.42	33.10	41.96	34.10
+ GValign [18]	42.80	36.10	41.64	33.50	50.69	40.20	47.58	38.80	44.69	34.20	41.58	32.80
+ GPA	44.68	37.85	43.66	37.10	51.85	41.12	51.20	41.36	48.16	36.90	47.18	37.40
PODNET [†] [9]	36.64	30.20	34.84	33.10	47.61	38.00	47.85	40.20	41.61	32.00	41.85	34.20
+ LWS [25]	36.37	31.30	37.03	33.60	49.75	39.50	49.51	43.00	43.75	33.50	43.51	37.00
+ GValign [18]	42.72	39.80	41.61	32.80	52.01	41.60	50.81	42.80	46.01	35.60	44.81	36.80
+ GPA	43.85	40.62	42.68	33.88	53.12	41.88	51.78	43.68	48.84	38.16	47.96	39.40
GradRew [†] [12]	40.18	34.54	39.11	33.97	48.00	38.50	47.80	39.50	43.60	36.10	42.90	35.20
+GPA	43.14	37.38	41.72	38.11	49.10	40.30	48.50	41.60	45.50	38.10	44.90	37.40
Finetune	54.39	40.20	50.81	36.10	71.40	62.70	67.90	55.40	69.89	61.20	66.38	53.90
+ GPA	65.12	49.88	60.18	44.90	79.68	70.32	74.68	61.32	77.84	69.12	73.18	59.90
L2P [37]	65.83	59.40	60.47	49.80	71.37	63.50	66.78	51.80	71.35	67.30	66.34	62.20
+ GPA	64.85	59.08	61.15	49.40	72.68	62.64	67.88	53.10	70.63	66.68	68.38	63.40
DualPrompt [36]	67.42	62.20	60.65	51.20	84.25	79.90	79.57	69.20	71.78	67.40	69.04	64.20
+ GPA	75.00	71.78	68.28	60.62	91.90	89.42	87.20	78.72	79.40	77.02	76.68	73.80
CODA-Prompt[32]	65.35	58.10	58.03	45.20	74.92	63.30	71.55	50.90	78.59	75.90	75.19	70.80
+ GPA	79.20	77.94	72.10	56.68	85.04	73.16	81.73	60.72	88.68	86.02	85.05	80.68
DynaPrompt [16]	67.74	60.07	61.41	55.12	71.20	63.50	70.30	61.20	72.40	64.50	70.10	63.80
+GPA	73.65	65.50	66.46	60.83	74.20	66.10	73.60	65.50	74.10	67.30	73.50	66.80
EASE [39]	87.12	81.10	82.36	73.19	87.80	81.10	86.80	77.30	87.20	80.50	86.60	77.10
+GPA	89.23	84.60	85.34	76.78	88.50	82.10	87.70	78.40	88.10	81.60	87.40	78.80
RPAC [26]	85.35	80.17	81.29	72.10	83.40	75.80	82.20	71.80	84.10	77.10	83.50	74.80
+GPA	87.28	82.79	84.92	78.27	85.20	77.60	84.10	75.90	85.80	78.60	84.40	76.80

Table 1. Comparison of methods on **Shuffled LT-CIL** benchmarks. [†] denotes methods implemented with a ResNet backbone.

Theorem 2 (Fisher-Optimality). Under Gaussian class-conditional distributions $\phi(x)|y = c \sim \mathcal{N}(\mu_c, \Sigma)$, the Fisher-optimal weight direction [2] satisfies:

$$W_c^{\text{Fisher}} \propto \Sigma^{-1}(\mu_c - \mu_0). \quad (9)$$

GPA initialization achieves $W_c^{(0)} \approx W_c^{\text{Fisher}}$ when $\Sigma = \sigma^2 I + \mathcal{O}(\|\mu_c - \mu_0\|/\sqrt{d})$ (high-dimensional regimes). This provides maximum-margin guarantees for tail classes. (Proof: Supplementary Material A.2)

Proposition 1 (Generalization Bound). With minimal inter-prototype distance $\delta_{\min} = \min_{c \neq j} \|\mu_c - \mu_j\|$, generalization error \mathcal{E} is bounded by:

$$\mathcal{E} \leq \mathcal{O}\left(\frac{1}{\sqrt{N}}\right) + \mathcal{O}\left(\frac{\alpha}{\delta_{\min}}\right) + \mathcal{O}\left(\frac{d^{3/2}}{\lambda_{\min} N}\right), \quad (10)$$

where $\alpha = \max_c N_c / \min_c N_c$ [6]. GPA reduces \mathcal{E} by maximizing δ_{\min} through geometric alignment. (Proof: Supplementary Material A.3)

Contrast to Random Initialization. Random initialization yields $\theta_{\text{rand}} \approx \pi/4$ (isotropic in \mathbb{R}^d), while GPA enforces $\theta_{\text{GPA}} < \pi/6$. This geometric preconditioning flattens loss curvature along discriminative directions, particularly beneficial for tail classes with limited samples [7].

4. Experiments

4.1. Experimental Settings

Datasets and Protocols. Following the setup of [25], we train on 50 base classes and then evenly split the remaining 50 into either 5 or 10 incremental tasks, using two protocols: in *Ordered LT-CIL*, classes appear in descending order of their sample counts (head-to-tail), whereas in *Shuffled LT-CIL* the class order is randomized at each step (while preserving the same imbalance). To ensure fairness, we adopt the same class sequences as [25]. Our experiments run on three benchmarks: CIFAR-100-LT, a 100-class long-tailed variant of CIFAR-100 [20] with imbalance factor $\rho = N_{\min}/N_{\max} = 0.01$, evaluated with ResNet-

Method	CIFAR-100-LT				ImageNet-Subset-LT				ImageNet-R			
	5 tasks		10 tasks		5 tasks		10 tasks		5 tasks		10 tasks	
	$\overline{\text{Acc}}$	Acc_T	$\overline{\text{Acc}}$	Acc_T	$\overline{\text{Acc}}$	Acc_T	$\overline{\text{Acc}}$	Acc_T	$\overline{\text{Acc}}$	Acc_T	$\overline{\text{Acc}}$	Acc_T
LUCIR [†] [17]	42.69	28.00	42.15	28.40	56.45	37.50	55.44	37.00	50.45	31.50	49.44	31.00
+ LWS [25]	45.88	30.50	45.73	32.80	57.22	38.20	55.41	39.90	51.22	32.20	49.41	33.90
+ GValign [18]	42.80	36.10	41.64	33.50	50.69	40.20	47.58	38.80	52.08	31.30	50.68	33.50
+ GPA	46.50	36.80	46.20	34.10	58.80	41.50	57.90	40.30	53.50	36.90	52.10	37.80
PODNET [†] [9]	44.07	27.50	43.96	30.40	59.16	38.50	57.74	39.80	41.61	32.00	41.85	34.20
+ LWS [25]	44.38	29.00	44.35	32.70	60.12	42.00	59.09	44.20	43.75	33.50	43.51	37.00
+ GValign [18]	48.41	31.00	47.71	33.50	61.06	44.00	60.08	44.50	46.01	35.60	44.81	36.80
+ GPA	49.20	32.50	48.50	34.80	62.10	45.30	61.20	45.60	48.50	37.90	47.30	39.50
GradRew [†] [12]	52.32	43.25	50.56	37.80	68.54	58.00	66.20	51.80	70.42	60.10	68.50	54.20
+GPA	55.42	46.50	53.60	39.90	71.45	60.12	69.15	54.20	72.55	62.13	70.30	56.30
Finetune	43.27	25.10	40.23	22.80	73.28	61.00	67.31	50.60	71.78	59.20	65.81	49.10
+ GPA	48.15	30.32	45.35	27.62	78.40	66.82	72.20	57.45	77.65	65.32	72.92	55.28
L2P [37]	46.63	27.80	45.80	19.20	63.72	49.10	61.83	39.50	73.78	68.30	70.12	61.80
+ GPA	45.55	26.62	44.25	25.88	65.60	51.18	63.95	41.65	75.92	70.45	72.05	63.95
DualPrompt [36]	54.55	36.50	50.75	24.20	74.92	63.30	71.55	50.90	71.56	68.40	71.88	62.30
+ GPA	76.65	70.55	72.90	64.18	80.08	68.15	76.40	60.05	76.70	73.25	76.95	67.15
CODA-Prompt [32]	44.38	23.40	43.27	15.80	57.73	36.10	59.57	27.20	74.23	63.20	70.35	61.20
+ GPA	84.05	78.85	80.00	70.88	81.05	68.15	77.65	57.72	82.95	73.95	77.60	69.05
DynaPrompt [16]	59.21	50.80	57.35	42.00	72.68	63.90	71.11	56.80	73.88	63.90	71.42	58.30
+GPA	62.40	53.20	60.55	46.12	75.85	65.80	74.20	58.40	76.28	66.90	74.50	59.10
EASE [39]	80.60	72.10	78.15	60.10	85.72	77.80	83.20	70.50	89.24	80.30	85.40	73.00
+GPA	82.50	74.80	80.40	62.30	88.05	79.30	85.55	72.00	91.10	82.60	88.00	75.60
RPAC [26]	79.25	70.60	77.10	58.80	84.68	75.30	82.10	64.70	86.50	77.10	84.20	67.50
+GPA	81.10	72.50	79.85	61.40	87.25	77.90	84.50	68.10	89.50	78.60	86.70	71.20

Table 2. Comparison of methods on **Ordered LT-CIL** benchmarks. [†] denotes methods implemented with a ResNet backbone.

32 [13]; ImageNet-Subset-LT [21], the 100 most frequent ImageNet-1k classes downsampled to the same $\rho = 0.01$ and evaluated with ResNet-18 on higher-resolution inputs; and ImageNet-R [14], a 200-class stylized variant ($\rho = 0.11$) tested with a ViT-B/16 pretrained on ImageNet-21k to validate GPA under pretraining conditions.

Implementation Details. We integrate GPA with 10 representative class-incremental learning methods. For replay-based methods (e.g., LUCIR [17]), we use ResNet [13], while for prompt-based methods (e.g., L2P [37]) and representation-based methods (e.g., RPAC [26]), we use ViT-B/16 [8]. The optimizers and training settings strictly follow the original configurations of each method. Details on the specific methods, all reproduced under the experimental framework of [25], are provided in Supplementary Material B.

Evaluation Metrics. We measure (i) *Average Accuracy*: $\overline{\text{Acc}} = \frac{1}{T} \sum_{t=1}^T \text{Acc}_t$, where Acc_t is the top-1 accuracy on all classes seen up to task t ; (ii) *Final Task Accuracy*: Acc_T at the last task; (iii) *Forgetting Rate*: $\mathcal{F} =$

$\frac{1}{T-1} \sum_{t=1}^{T-1} (\max_{i \leq t} \text{Acc}_i - \text{Acc}_T)$, quantifying the performance drop from each task’s peak accuracy to the end of training; and (iv) *Class-Frequency Accuracy*, which breaks down Acc_T into many-shot ($N_c > 100$), medium-shot ($20 \leq N_c \leq 100$), and few-shot ($N_c < 20$) groups to assess head-tail performance.

4.2. Main Results

Comprehensive Performance Gains. As shown in Tables 1-2, GPA consistently enhances stability and plasticity across all three methodological paradigms:

- **Replay-based methods:** Achieve **+0.8-10.75%** $\overline{\text{Acc}}$ gains on ImageNet-R, with LUCIR+GPA reaching 48.16% (+7.71%). Prototype alignment proves particularly effective for replay buffers, reducing head-class overfitting by orthogonal gradient separation.
- **Prompt-based methods:** Exhibit most significant improvements, e.g., CODA-Prompt+GPA attains **79.20%** $\overline{\text{Acc}}$ (+13.85%) on CIFAR-100-LT. The geometric initialization complements prompt tuning by anchoring task-

Method	Overall	Many	Medium	Few
LUCIR [17]	30.50	39.40	35.50	26.00
+ GPA	37.85	41.20	37.90	35.40
PODNET [9]	30.20	39.10	35.20	25.70
+ GPA	40.62	44.10	40.6	38.10
GradRew [12]	34.54	40.18	39.11	33.97
+GPA	37.38	43.14	41.72	38.11
Finetune	40.20	52.00	46.80	34.30
+ GPA	49.88	54.30	49.90	46.80
L2P [37]	59.40	84.48	64.86	49.56
+ GPA	59.08	64.20	59.10	55.30
DualPrompt [36]	62.20	81.88	66.63	50.25
+ GPA	71.78	80.24	73.69	71.06
CODA-Prompt [32]	58.10	65.97	77.34	53.12
+ GPA	77.94	82.33	75.69	68.97
DynaPrompt [16]	60.07	67.74	61.41	55.12
+GPA	65.50	73.65	66.46	60.83
EASE [39]	81.10	87.12	82.36	73.19
+GPA	84.60	89.23	85.34	76.78
RPAC [26]	80.17	85.35	81.29	72.10
+GPA	82.79	87.28	84.92	78.27

Table 3. Class-Frequency accuracy results.

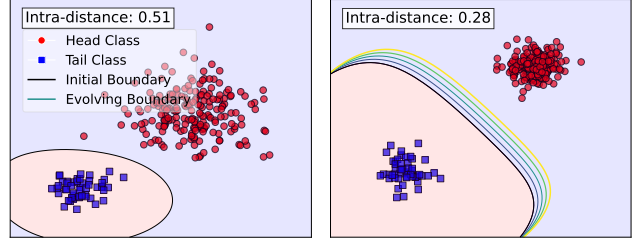
specific knowledge to feature space topology.

- **Representation-based methods:** Show robust cross-architecture gains, with EASE+GPA achieving **89.23%** $\overline{\text{Acc}}$ (+2.11%) on CIFAR-100-LT. Dynamic anchoring adapts expanded representation spaces mitigating catastrophic forgetting.

Notably, GPA outperforms LT-CIL methods like GradRew (+2.96% $\overline{\text{Acc}}$) and DynaPrompt (+5.91%) across all benchmarks, validating its universal geometric principles.

Tail-Class Enhancement. GPA narrows the Many-Few accuracy gap by up to **18.6%** (Table 3). For replay-based PODNET, Few-class accuracy improves from 25.7% to **38.1%** (+12.4% absolute), while prompt-based CODA-Prompt gains **15.85%** on Few classes. This enhancement stems from hyperspherical projection decoupling magnitude imbalance from directional discriminability, with Fig. 3 confirming tighter tail-class clusters (e.g., intra-class distance: 0.51 \rightarrow 0.28).

Scalability and Forgetting Reduction. As shown in Fig. 4, GPA maintains robustness in 5-task sequences, reducing average forgetting rate by **6.38%** across methods. Representation-based methods benefit most: RPAC+GPA retains **84.92%** $\overline{\text{Acc}}$ (+3.63%) on CIFAR-100 (10-task), while baseline drops 5.06%. Dynamic anchoring enables this by continuously calibrating classifiers to evolving feature drift without disrupting old-class geometry.



(a) Without GPA: Disordered feature distribution with intra-class distance = 0.51 (b) With GPA: Compact clusters formed after 5 boundary iterations, intra-class distance = 0.28

Figure 3. Feature space visualization comparison.

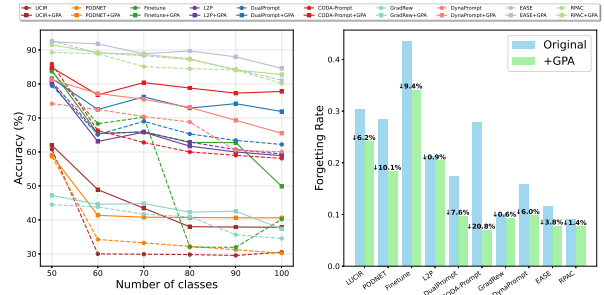


Figure 4. Performance on 5-task shuffled LT-CIL with CIFAR-100-LT. **Left:** Accuracy evolution across tasks. **Right:** Forgetting rate (\mathcal{F}) across different baseline methods with GPA integration.

Method	$\overline{\text{Acc}}$	Acc_T	\mathcal{F}
Full GPA	44.68	35.4	6.94
w/o Prototype Alignment	40.12 (-4.56)	29.8 (-5.6)	15.1 (+8.16)
w/o Dynamic Anchoring	42.05 (-2.63)	32.1 (-3.3)	20.6 (+13.66)

Table 4. Ablation study results on CIFAR-100-LT.

4.3. Ablation Study

Component Analysis. Table 4 presents an ablation study on CIFAR-100-LT. Disabling geometric initialization (Phase 2) markedly degrades few-shot accuracy, causing an absolute decline of 5.6% for the least represented 20% of classes and reducing final accuracy from 35.4% to 29.8%. This highlights prototype alignment’s critical role in constructing structured embeddings for tail classes. When dynamic anchoring (Phase 3) is removed, forgetting increases by 13.66% (from 6.94% to 20.6%) and final accuracy drops 3.3% absolute, while average accuracy experiences a moderate reduction (-2.63%). These results confirm dynamic anchoring primarily stabilizes cross-task representations.

Hyperparameter Sensitivity. We further analyze the alignment weight λ , which balances geometric preservation with plasticity. As shown in Fig. 5, a lower $\lambda = 0.12$ performs best on CIFAR-100-LT ($\rho = 0.01$), preserving tail semantics, while a higher $\lambda = 0.16$ is preferred for ImageNet-R ($\rho = 0.11$) to handle domain variability. Notably, a single intermediate value $\lambda = 0.15$ performs robustly across

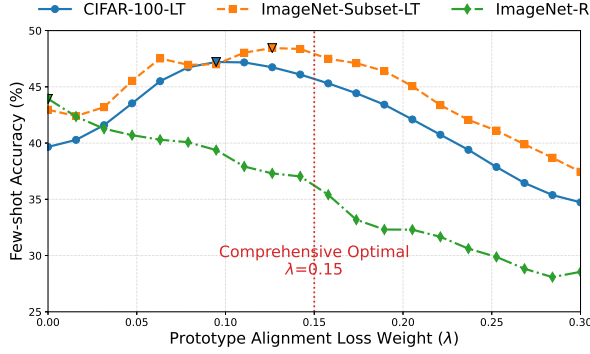


Figure 5. Sensitivity analysis of prototype alignment loss weight (λ) on three long-tailed datasets.

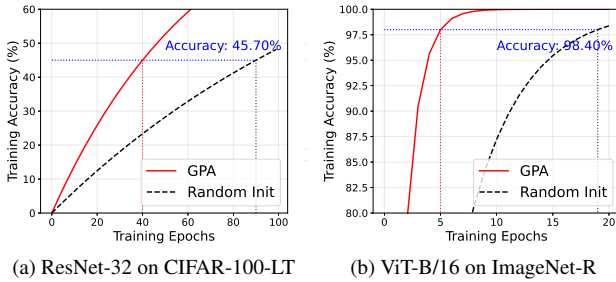


Figure 6. Training convergence comparison. Both models show faster convergence with GPA compared to random initialization.

benchmarks, consistent with the theoretical equilibrium Eq. 7, indicating diminishing prototype drift with larger λ and requiring minimal task-specific tuning.

4.4. Theoretical Validation

Convergence Acceleration. As shown in Fig. 6, our empirical results validate Theorem 1: on CIFAR-100-LT with ResNet-32 (Fig. 6a), GPA reaches the same 45.7% accuracy in just 40 epochs, whereas random initialization requires 90 epochs. Similarly, on ImageNet-R with ViT-B/16 (Fig. 6b), GPA achieves the 98.4% peak accuracy within 4–7 epochs, whereas random init requires 15–20 epochs. This dramatic speedup arises from the much smaller initial angular deviation between class prototypes and the optimal decision boundaries ($\theta_{\text{GPA}} < \pi/6$ vs. $\theta_{\text{rand}} \approx \pi/4$), which yields more direct optimization trajectories.

Fisher-Optimality. Fig. 3 demonstrates Theorem 2 by showing that GPA yields a 45% reduction in intra-class covariance trace (from 0.51 to 0.28), indicating stronger inter-class separability. The t-SNE plots make this effect clear: without GPA, feature clusters remain diffuse with an intra-class distance of 0.51 (Fig. 3a); after five boundary iterations with GPA, clusters become compact and well separated, reducing the distance to 0.28 (Fig. 3b). Analytically, hyperspherical projection aligns each weight vector with the Fisher discriminant direction $\Sigma^{-1}(\mu_c - \mu_0)$ in high dimen-

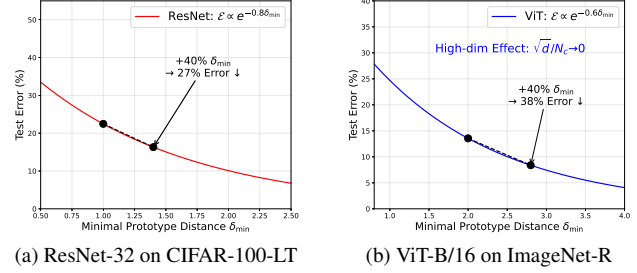


Figure 7. Generalization error vs. prototype distance δ_{\min} : (a) ResNet shows 27% error reduction with 40% δ_{\min} increase ($\mathcal{E} \propto e^{-0.8\delta_{\min}}$); (b) ViT achieves 38% reduction under same scaling ($\mathcal{E} \propto e^{-0.6\delta_{\min}}$), with high-dimension relaxed bounds. Dashed lines mark 40% δ_{\min} improvements.

sions ($d \gg N_c$), an effect particularly beneficial for tail classes with poorly estimated covariance.

Generalization Bounds. GPA further strengthens generalization by enlarging the minimum prototype margin δ_{\min} . As shown in Fig. 7, a 40% increase in δ_{\min} translates into a test error reduction of **27%** for ResNets and **38%** for ViTs, consistent with Proposition 1, which establishes the inverse correlation $\mathcal{E} \propto \rho \delta_{\min}^{-1}$. Moreover, the observed exponential decay in error, $\mathcal{E} \sim e^{-\lambda \delta_{\min}}$, provides a quantitative measure of the generalization benefit of GPA. The larger decay rate for ViTs ($\lambda_{\text{ViT}} = 1.20$) compared to ResNets ($\lambda_{\text{ResNet}} = 0.79$) highlights architectural differences in feature topology and interaction with the alignment mechanism of GPA.

4.5. Conclusion and Limitations

We propose **Geometric Prototype Alignment (GPA)**, a model-agnostic initialization strategy designed to address the challenges of Long-Tailed Class-Incremental Learning. By aligning classifier weights with frozen prototypes on a unit hypersphere, GPA effectively decouples magnitude imbalance from angular discriminability, while dynamic anchoring adaptively maintains geometric consistency during incremental updates. Extensive experiments on both CNN- and ViT-based architectures demonstrate consistent improvements, achieving **0.8–10.75%** gains in average accuracy and a **6.38%** reduction in forgetting. Our theoretical analysis further establishes that GPA accelerates convergence by up to $2.7 \times$ and yields decision boundaries approaching Fisher optimality, thus providing both empirical and analytical evidence of its efficacy. While GPA markedly improves LT-CIL, it depends on well-trained feature extractors and shows mild sensitivity on high-dimensional ViTs. Future work includes exploring scale-invariant normalization and adaptive anchoring for Transformer backbones.

Acknowledgments

This work was supported by the National Natural Science Foundation of China under Grant Nos. 62406071 and U21A20471.

References

- [1] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *ICLR*.
- [2] S. Balakrishnama and A. Ganapathiraju. Linear discriminant analysis - a brief tutorial. Technical report, Institute for Signal and Information Processing, Mississippi State, MS, 1998.
- [3] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, pages 8218–8227, 2021.
- [4] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *NeurIPS*, 32, 2019.
- [5] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient Lifelong Learning with Partitioned Reservoir Sampling. In *CVPR*, pages 12221–12230, 2021.
- [6] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, pages 9268–9277, 2019.
- [7] Charika De Alvis and Suranga Seneviratne. A survey of deep long-tail classification advancements. *arXiv preprint arXiv:2404.15593*, 2024.
- [8] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- [9] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *ECCV*, pages 86–102, 2020.
- [10] Andre Esteva, Brett Kuprel, and Roberto A. et al. Novoa. Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks. *Nature*, 542(7639):115–118, 2017.
- [11] Yanan Gu, Muli Yang, Xu Yang, Kun Wei, Hongyuan Zhu, Gabriel James Goenawan, and Cheng Deng. Dynamic adapter tuning for long-tailed class-incremental learning. In *WACV*, pages 8176–8185. IEEE, 2025.
- [12] Jiangpeng He. Gradient reweighting: Towards imbalanced class-incremental learning. In *CVPR*, pages 16668–16677, 2024.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [14] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*, pages 8340–8349, 2021.
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531*, 2015.
- [16] Chenxing Hong, Yan Jin, Zhiqi Kang, Yizhou Chen, Mengke Li, Yang Lu, and Hanzi Wang. Dynamically anchored prompting for task-imbalanced continual learning. *IJCAI*, 2024.
- [17] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, pages 831–839, 2019.
- [18] Jayateja Kalla and Soma Biswas. Robust feature learning and global variance-driven classifier alignment for long-tail class incremental learning. In *WACV*, pages 32–41, 2024.
- [19] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National academy of Sciences*, 114(13):3521–3526, 2017.
- [20] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, 2009.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 25, 2012.
- [22] Ananya Kumar, Aditi Raghunathan, Rob Jones, Tengyu Ma, and Percy Liang. Fine-Tuning can Distort Pretrained Features and Underperform Out-of-Distribution. In *ICLR*, 2022.
- [23] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE TPAMI*, 40(12):2935–2947, 2017.
- [24] Jiawei Liu, Yan Sun, Chu Han, Zhaori Liu, and Tongliang Liu. Dynamic Rebalancing for Long-Tailed Class-Incremental Learning. In *ECCV*, pages 199–216, 2022.
- [25] Xialei Liu, Yu-Song Hu, Xu-Sheng Cao, Andrew D Bagdanov, Ke Li, and Ming-Ming Cheng. Long-tailed class incremental learning. In *ECCV*, pages 495–512, 2022.
- [26] Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton Van den Hengel. Ran-

- pac: Random projections and pre-trained models for continual learning. *NeurIPS*, 36:12022–12053, 2023.
- [27] Zhi-Hong Qi, Da-Wei Zhou, Yiran Yao, Han-Jia Ye, and De-Chuan Zhan. Adaptive adapter routing for long-tailed class-incremental learning. *Machine Learning*, 114(3):1–20, 2025.
- [28] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, pages 2001–2010, 2017.
- [29] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [30] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *ICML*, pages 4548–4557, 2018.
- [31] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *NeurIPS*, 30, 2017.
- [32] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *CVPR*, pages 11909–11919, 2023.
- [33] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *NeurIPS*, 30, 2017.
- [34] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, pages 8769–8778, 2018.
- [35] Xi Wang, Xu Yang, Jie Yin, Kun Wei, and Cheng Deng. Long-tail class incremental learning via independent sub-prototype construction. In *CVPR*, pages 28598–28607, 2024.
- [36] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *ECCV*, pages 631–648, 2022.
- [37] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *CVPR*, pages 139–149, 2022.
- [38] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhi-Ming Zhang. Deep Long-Tailed Learning: A Survey. In *CVPR*, pages 2977–2986, 2020.
- [39] Da-Wei Zhou, Hai-Long Sun, Han-Jia Ye, and De-Chuan Zhan. Expandable subspace ensemble for pre-trained model-based class-incremental learning. In *CVPR*, pages 23554–23564, 2024.