

FLYPROMPT: BRAIN-INSPIRED RANDOM-EXPANDED ROUTING WITH TEMPORAL-ENSEMBLE EXPERTS FOR GENERAL CONTINUAL LEARNING

Hongwei Yan¹ Guanglong Sun¹ Kanglei Zhou² Qian Li³ Liyuan Wang^{2*} Yi Zhong^{1*}

¹School of Life Sciences, IDG/McGovern Institute for Brain Research, Tsinghua University

²Department of Psychological and Cognitive Sciences, Tsinghua University

³School of Medicine, Shenzhen Campus of Sun Yat-Sen University

{yanhw22, sgl23}@mails.tsinghua.edu.cn, liqian255@mail.sysu.edu.cn
{zhoukanglei, liyuanwang, zhongyithu}@tsinghua.edu.cn

ABSTRACT

General continual learning (GCL) challenges intelligent systems to learn from single-pass, non-stationary data streams without clear task boundaries. While recent advances in continual parameter-efficient tuning (PET) of pretrained models show promise, they typically rely on multiple training epochs and explicit task cues, limiting their effectiveness in GCL scenarios. Moreover, existing methods often lack targeted design and fail to address two fundamental challenges in continual PET: how to allocate expert parameters to evolving data distributions, and how to improve their representational capacity under limited supervision. Inspired by the fruit fly’s hierarchical memory system characterized by sparse expansion and modular ensembles, we propose FlyPrompt, a brain-inspired framework that decomposes GCL into two subproblems: *expert routing* and *expert competence improvement*. FlyPrompt introduces a randomly expanded analytic router for instance-level expert activation and a temporal ensemble of output heads to dynamically adapt decision boundaries over time. Extensive theoretical and empirical evaluations demonstrate FlyPrompt’s superior performance, achieving up to 11.23%, 12.43%, and 7.62% gains over state-of-the-art baselines on CIFAR-100, ImageNet-R, and CUB-200, respectively. Our source code is available at FlyGCL.

1 INTRODUCTION

General Continual Learning (GCL) (Buzzega et al., 2020; De Lange et al., 2021), aims to equip intelligent systems with the ability to learn continuously from non-stationary, single-pass data streams, where tasks may not have clear boundaries and can evolve over time. Unlike traditional Continual Learning (CL) (Wang et al., 2024b; Parisi et al., 2019), which assumes well-defined task boundaries and multiple training epochs, GCL presents a much more challenging problem, as it requires rapid adaptation, robust knowledge retention, and efficient resource usage under conditions of limited supervision and task ambiguity (Fig. 1). The ability to effectively tackle GCL has profound implications for real-world applications such as autonomous agents and personal assistants, where systems must learn from dynamic environments without clear task definitions.

Recent advances¹ in parameter-efficient tuning (PET) of pretrained models (PTMs) have shown promise in CL (Wang et al., 2022d;c; Smith et al., 2023), but they still face fundamental limitations under GCL conditions. Such methods introduce task-specific prompt experts to adapt PTMs incrementally, and typically rely on clear task cues and sufficient gradient updates to allocate and train expert modules (Wang et al., 2024a; 2022b). However, those assumptions no longer hold in GCL (Koh et al., 2021; Moon et al., 2023). We therefore identify two fundamental challenges that remain unresolved: (1) how to dynamically route inputs to appropriate experts without task labels or iterative training, and (2) how to ensure that each expert maintains strong and adaptive representations under sparse and imbalanced supervision. Both remain non-trivial and underexplored.

*Corresponding Authors: L. Wang and Y. Zhong.

¹Due to the page limit, we present a comprehensive summary of related work in Appendix B.

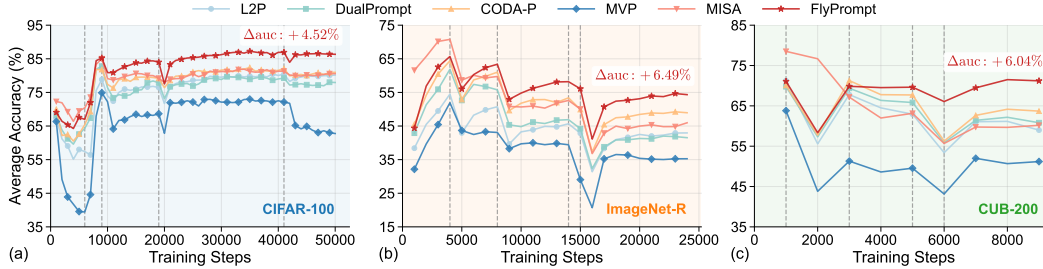


Figure 1: Any-time average accuracy of GCL methods over three datasets using Sup-21K. Dashed lines indicate task transition. Δauc , the improvement of area-under-curve score by FlyPrompt.

The complexity of GCL has also been extensively studied in biological systems, where organisms have evolved efficient strategies for lifelong learning in dynamic environments. The fruit fly *Drosophila* provides a compelling model: despite having fewer than 100,000 neurons, it exhibits robust memory consolidation, context-aware behavior, and stable learning under minimal supervision (Davis, 2023; Li et al., 2020; Modi et al., 2020; Oswald & Waddell, 2015). These capabilities are largely attributed to the mushroom body, a central brain structure that encodes sensory inputs via sparse random projections and organizes learning into modular, hierarchical compartments (Aso et al., 2014b; Dasgupta et al., 2017; Wang et al., 2021; 2022a; 2023b; Wang & Li, 2025) (see Fig. 4(Left)). Projection neurons (PNs) from the antennal lobe connect randomly to Kenyon cells (KCs) in mushroom body, yielding high-dimensional sparse codes that support input separation and routing even under noisy or overlapping conditions (Turner et al., 2008; Honegger et al., 2011). Furthermore, different KC subregions exhibit plasticity on distinct timescales (Aso et al., 2014a; Aso & Rubin, 2016), enabling both rapid adaptation and long-term consolidation (Cervantes-Sandoval et al., 2013; Bouzaiane et al., 2015; Wang & Li, 2025). These mechanisms closely mirror the goals of GCL, offering principled inspiration for tackling its core challenges.

Building upon these neurobiological principles and our preliminary analysis in Sec. 2.2, we propose to decompose the GCL challenges into two essential subproblems: (1) *expert routing*, which aims to assign each input to an appropriate subnetwork (expert) under unknown and shifting task boundaries; and (2) *expert competence improvement*, which seeks to enhance the robustness and adaptability of each expert given limited training and imbalanced class exposure. To address these challenges, we introduce **FlyPrompt**, a brain-inspired framework that integrates two key components: (i) a *Random Expanded Analytic Router* (REAR) that mimics the fruit fly’s sparse expansion circuit to rapidly assign inputs to experts in a forward-only, closed-form manner; and (ii) a *Task-wise Experts with Temporal Ensemble* (TE^2) that captures knowledge across multiple time scales using exponential moving averages, mirroring the compartmental consolidation observed in the mushroom body.

FlyPrompt is supported by both theoretical analysis and empirical validation. Across diverse GCL benchmarks, including CIFAR-100, ImageNet-R, and CUB-200, it consistently outperforms state-of-the-art CL and GCL methods, achieving accuracy improvements of up to 11.23%, 12.43%, and 7.62%, respectively. By integrating biologically grounded design with principled algorithmic structure, FlyPrompt offers an interdisciplinary perspective on addressing the core challenges of GCL and also exemplifies the potential of the emerging field of NeuroAI (Zador et al., 2023).

2 PRELIMINARIES

In this section, we formulate GCL, and then evaluate PET methods in an instantiated GCL scenario.

2.1 PROBLEM FORMULATION

In CL, a model learns sequential tasks $t \in \{1, \dots, T\}$, each associated with a dataset $\mathcal{D}_t = (\mathbf{x}_t, \mathbf{y}_t)$ where $\mathbf{x}_t \in \mathcal{X}_t$ and $\mathbf{y}_t \in \mathcal{Y}_t$. The model comprises a backbone $f_\theta(\cdot)$ and an output head $g_\psi(\cdot)$, which together produce predictions $\hat{\mathbf{y}} = g_\psi(f_\theta(\mathbf{x}))$. The objective is to learn a unified mapping from input domains $\mathcal{X} = \bigcup_t \mathcal{X}_t$ to label spaces $\mathcal{Y} = \bigcup_t \mathcal{Y}_t$. Classical CL settings impose structural assumptions on the input or label space. Domain-incremental learning (DIL) assumes disjoint input domains with a shared label space ($\mathcal{X}_i \cap \mathcal{X}_j = \emptyset$, $\mathcal{Y}_i = \mathcal{Y}_j$), while task-incremental and class-incremental learning (TIL, CIL) assume disjoint label spaces ($\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset$), with TIL additionally

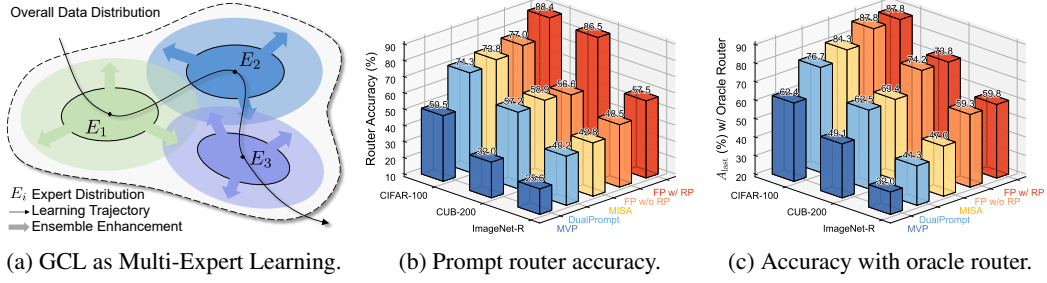


Figure 2: Empirical analysis of GCL. (a) A schematic illustration of GCL viewed as multi-expert collaboration. (b) Prompt selection accuracy for methods with explicit expert routing designs. (c) Final average accuracy ($A_{\text{last}}, \uparrow$) when using a test-time oracle to provide the correct prompt identity. Results evaluated across three benchmarks with Sup-21K. FP, FlyPrompt. RP, Random Projection.

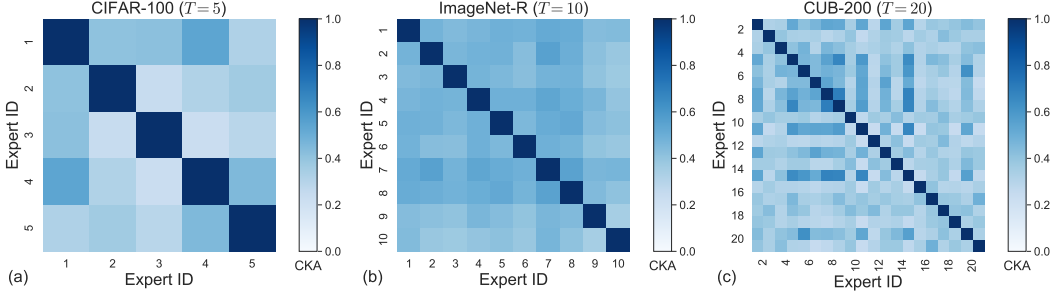


Figure 3: CKA similarity of feature representations between experts of MVP on three datasets.

providing task identity at test time (Van de Ven & Tolia, 2019). Under these assumptions, the learning objective can be decomposed into two orthogonal subproblems: *task identity prediction* (TIP), which selects an appropriate task-specific module, and *within-task prediction* (WTP), which performs classification under the selected module (Kim et al., 2022; Wang et al., 2023a).

GCL, however, lifts these assumptions and operates under substantially more challenging conditions. Tasks arrive as a one-pass data stream, and their label spaces may overlap with non-negligible probability: $\forall i \neq j, P(\mathcal{Y}_i \cap \mathcal{Y}_j \neq \emptyset) > 0$ (Koh et al., 2021). This entangles inter-task and intra-task interference, undermining the TIP-WTP orthogonality. Especially when using pretrained backbones, the strong priors encoded in PTMs already bias prediction before adaptation, causing task identity and class discrimination to co-evolve (Wang et al., 2023a; 2024a). Moreover, GCL introduces additional difficulties such as severe intra-task class imbalance (e.g., long-tailed distributions) and limited training iterations. These problems are compounded by memory constraints (Buzzega et al., 2020; De Lange et al., 2021), where storing past data is restricted or disallowed.

A representative instantiation of GCL is Si-Blurry (Moon et al., 2023), which explicitly partitions the global label space \mathcal{Y} into a disjoint subset \mathcal{Y}^D and a blurry subset \mathcal{Y}^B , with $\mathcal{Y} = \mathcal{Y}^D \cup \mathcal{Y}^B$ and $\mathcal{Y}^D \cap \mathcal{Y}^B = \emptyset$. The *disjoint class ratio* $r_D = |\mathcal{Y}^D|/|\mathcal{Y}|$ controls the proportion of task-specific classes, while the *blurry sample ratio* r_B determines how frequently classes in \mathcal{Y}^B reappear across tasks. This flexible design captures the stochasticity and heterogeneity of GCL, which has been validated in recent theoretical and empirical work (Mi et al., 2020; Zhuang et al., 2024; Kang et al., 2025). We therefore adopt Si-Blurry as the default GCL benchmark (see Appendix D for discussion about the task/session boundary information and our empirical rationality of using Si-Blurry).

2.2 ANALYSIS OF GCL METHODS WITH EXPERTS

Recent CL and GCL methods increasingly adopt PET techniques on top of PTMs. These methods can be seen as lightweight extensions of architecture-based CL (Zhu et al., 2021; Wang et al., 2023b), where instead of expanding full networks, they introduce trainable modules p (e.g., adapters, prompts, and LoRA) that act as semantic-aware adaptation experts to give instructed outputs $f_\theta(x; p)$. A common strategy is to maintain a pool of such experts and design a router to assign inputs to the appropriate ones. However, most existing methods, such as L2P (Wang et al., 2022d), DualPrompt (Wang et al., 2022c), MVP (Moon et al., 2023), CODA-P (Smith et al., 2023), and MISA (Kang et al., 2025), train these routing functions synchronously with the stream of incoming

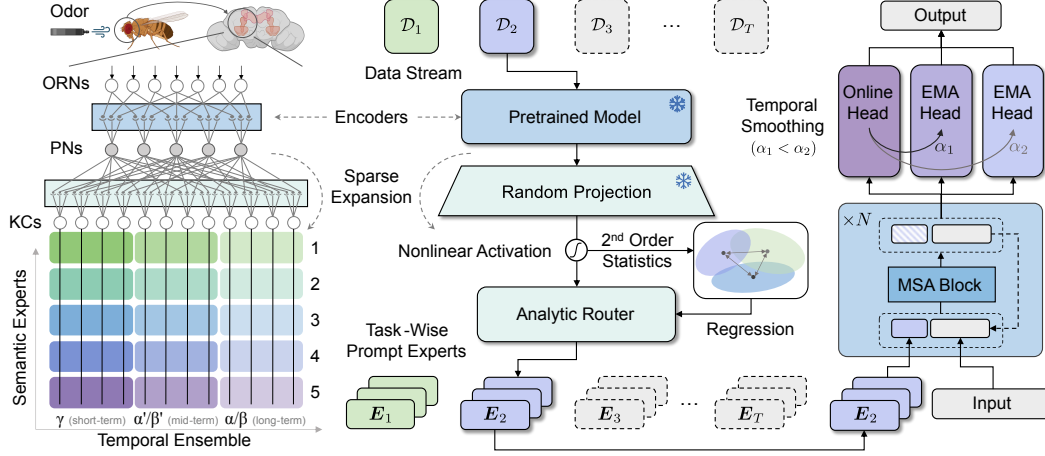


Figure 4: Method overview. Inspired by the fruit fly’s olfactory memory system (*Left*), FlyPrompt incorporates a random-expanded analytic router (*Middle*) and temporal ensemble-based experts (*Right*). ORNs, olfactory receptor neurons. PNs, projection neurons. KCs, Kenyon cells.

data, making them vulnerable to distributional shifts and limited iterations. These issues are especially pronounced in GCL, where blurry task boundaries and online constraints prohibit iterative tuning, and class imbalance further weakens the quality of per-task representations.

Based on this observation, we propose to decompose the GCL problem into two critical subproblems: **expert routing**, which determines which expert to assign to each input; and **expert competence**, which enhances the quality and robustness of each expert’s representation. Compared to the classical TIP-WTP formulation for CL with clearly segmented tasks, GCL’s blurry structure and overlapping class distributions make semantic-level experts a more suitable abstraction for adaptation (see Fig. 2a). Moreover, since the same class may appear in multiple tasks, the correspondence between classes and experts is inherently one-to-many, which further complicates routing.

To better understand these challenges, we conduct a preliminary empirical study from a multi-expert collaboration (details in Sec. 4.1). We first evaluate the routing accuracy of methods that explicitly predict expert identity (e.g., DualPrompt, MVP, MISA, and our FlyPrompt) after all GCL training tasks. A prediction is considered correct if the selected expert belongs to the set of experts previously trained on the true label y_t , acknowledging the overlap across tasks. As shown in Fig. 2b, existing routers based on similarity or contrastive losses still exhibit considerable limitations in expert selection. Next, we evaluate the final average accuracy under an oracle router that always selects a correct expert for each input (Fig. 2c). The results reveal that, even with perfect routing, previous methods still exhibit inferior performance, highlighting a second bottleneck: the limited competence of individual experts. In a PTM-based context, such competence depends not only on the representation space shaped by each expert module, but also on how well the output head can maintain consistent decision boundaries over time; even when an early expert’s encoder is frozen, a single head that keeps adapting to later data can gradually become misaligned with its fixed representation. To verify this point, our analysis of expert-specific representations using centered kernel alignment (CKA; Appendix F.4) in Fig. 3 confirms that experts indeed specialize in distinct feature subspaces, underscoring the need for accurate expert assignment. Together with the observed degradation under an oracle router, these results support decomposing GCL into the interacting subproblems above, and clarify that expert competence must account for both representation quality and decoding robustness.

3 FLYPROMPT: A BRAIN-INSPIRED GCL APPROACH

In this section, we propose **FlyPrompt**, an innovative brain-inspired approach designed to tackle the key challenges of GCL by explicitly improving **expert routing** and **expert competence**. As shown in Fig. 4, FlyPrompt consists of two core components: (i) a *Random Expanded Analytic Router* (REAR) that employs fixed random projections and closed-form updates to assign inputs to experts, inspired by the sparse expansion circuits in fruit flies, and (ii) *Task-wise Experts with Temporal Ensemble* (TE^2) that adaptively refine class boundaries over time to improve expert-level performance, reflecting modularized ensembles architecture in fruit flies’ neural systems.

3.1 RANDOM EXPANDED ANALYTIC ROUTER

Recent studies in CL have explored the use of Random Projection (RP) to construct forward-only learners with desirable properties such as rapid adaptation and immunity to catastrophic forgetting (Zhuang et al., 2022; 2024; McDonnell et al., 2024). When combined with nonlinear activation, RP can significantly improve the linear separability of input features by capturing high-order interactions without backpropagation. Notably, this mechanism aligns closely with the olfactory system of the fruit fly, where projection neurons (PNs) connect sparsely and randomly to a high-dimensional population of Kenyon cells (KCs) in the mushroom body, achieving a nearly 40-fold expansion. Global inhibition then enforces sparsity, allowing robust and efficient pattern separation (Dasgupta et al., 2017; Honegger et al., 2011). This neurobiological design inspires REAR, which mimics the biologically grounded random expansion to enable efficient, non-iterative expert selection in GCL.

Concretely, given a pretrained backbone $f_\theta(\cdot)$ that maps input \mathbf{x} to an embedding $\mathbf{h} = f_\theta(\mathbf{x}) \in \mathbb{R}^d$ of dimension d , we apply a fixed RP followed by a nonlinear activation:

$$\varphi(\mathbf{x}) = \sigma(f_\theta(\mathbf{x})\mathbf{R}) = \sigma(\mathbf{h}\mathbf{R}) \in \mathbb{R}^M, \quad (1)$$

where $\mathbf{R} \in \mathbb{R}^{d \times M}$ is a random matrix with $R_{i,j} \sim \mathcal{N}(0, 1)$, $M > d$, and $\sigma(\cdot)$ is an element-wise activation function (e.g., ReLU). The resulting feature $\varphi(\mathbf{x})$ is sparse and high-dimensional.

During online training, we associate each task t with a corresponding expert E_t . For each incoming batch $\mathcal{B}_i \subset \mathcal{D}_t$ of size B , we compute the projected features $\Phi_i \in \mathbb{R}^{B \times M}$, whose row vectors are $\{\varphi(\mathbf{x})^\top | (\mathbf{x}, y) \in \mathcal{B}_i\}$, and update two statistics: the *Gram matrix* $\mathbf{G} \in \mathbb{R}^{M \times M}$ capturing second-order feature correlations; and the *prototype matrix* $\mathbf{Q} \in \mathbb{R}^{M \times T}$ storing expert-wise feature sums:

$$\mathbf{G} \leftarrow \mathbf{G} + \Phi_i^\top \Phi_i, \quad \mathbf{Q} \leftarrow \mathbf{Q} + \Phi_i^\top \mathbf{C}_t, \quad (2)$$

where $\mathbf{C}_t \in \mathbb{R}^{B \times T}$ whose row vectors are the same one-hot embedding $\mathbf{c}_t \in \{0, 1\}^T$ for expert E_t . We then construct a router matrix $\mathbf{U} \in \mathbb{R}^{T \times M}$ by minimizing the following objective:

$$\mathcal{L}(\mathbf{U}) = \sum_{t=1}^T \sum_{(\mathbf{x}_t, y) \in \mathcal{D}_t} \|\varphi(\mathbf{x}_t)\mathbf{U}^\top - \mathbf{c}_t\|_2^2 + \lambda \|\mathbf{U}\|_F^2, \quad (3)$$

where $\lambda > 0$ is the regularization parameter. This objective encourages the router to map samples from task t to the corresponding expert E_t while maintaining numerical stability. Using the accumulated statistics \mathbf{G} and \mathbf{Q} , the closed-form solution to this optimization problem is given by:

$$\hat{\mathbf{U}}^\top = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{Q}. \quad (4)$$

The calculation of $\hat{\mathbf{U}}$ is only needed once upon evaluation, therefore this optimization process is efficient and lightweight compared to gradient-based routing mechanisms. At inference time, the routing score s and selected expert \hat{E} for an input \mathbf{x} given router $\hat{\mathbf{U}}$ is computed as:

$$\mathbf{s}(\mathbf{x}) = \varphi(\mathbf{x})\hat{\mathbf{U}}^\top \in \mathbb{R}^T, \quad \hat{E}(\mathbf{x}) = \arg \max_{t \leq T} s_t(\mathbf{x}). \quad (5)$$

This routing mechanism is efficient, biologically motivated, and requires no gradient updates, making it well-suited for GCL’s online, single-pass constraints. Unlike prior methods based on random expanded features, such as RanPAC (McDonnell et al., 2024) or ACIL (Zhuang et al., 2022), which apply closed-form ridge regression directly for final classification on fixed representations, REAR uses random projections solely for instance-level expert routing while keeping each expert’s prompts and heads fully trainable. Empirical comparison between the analytic router (REAR) and analytic classifier (RanPAC) under GCL benchmarks is shown in Appendix Tab. 13. And we further demonstrate the superiority of REAR upon alternative routing strategies in Appendix Tab. 17. We then summarize the core theoretical guarantee that explains why REAR yields reliable routing in the expanded sparse feature space. Full assumptions and proofs are included in Appendix E.1.

Theorem 1 (REAR, informal). *With high probability over the random expansion and the data stream, the population excess risk of the ridge router learned from online statistics admits the following decomposition:*

$$\mathcal{R}(\hat{\mathbf{U}}) - \mathcal{R}(\mathbf{U}^*) \lesssim \sqrt{\log(N)/M} + (\sqrt{N} \lambda)^{-1} + \lambda,$$

for suitable universal constants. Therefore, by increasing the expansion dimension M and the number of samples N , and choosing the regularization parameter λ to balance estimation error and bias, the population excess risk (and, under a fixed margin assumption on expert scores; see Appendix E.1, the misrouting probability) can be made arbitrarily small.

Interpretation. The first term reflects the approximation error due to finite random features; increasing M improves the expressive power of the expansion. The second term captures the variance arising from finite data, which diminishes as N grows or λ increases. The third term represents the bias introduced by ridge regularization, which stabilizes learning but limits expressiveness if it is too large. In practice, this decomposition implies that robust and forward-only routing can be achieved by employing sufficiently rich random expansions and moderate regularization, without requiring task-level or iterative refinement.

3.2 TASK EXPERTS AS TEMPORAL ENSEMBLES

To improve the competence of each expert under dynamic distributions, we draw inspiration from the fruit fly’s KCs in the mushroom body and their connections to output neurons. This brain structure integrates multi-timescale plasticity and hierarchical processing across subregions, where γ KCs mediate short-term memory, α'/β' KCs support intermediate memory, and α/β KCs are critical for long-term memory consolidation (Krashes et al., 2007; Cervantes-Sandoval et al., 2013; Bouzaiane et al., 2015). These KC subtypes are sequentially recruited during learning and exhibit compartment-specific modulation by dopamine neurons (Aso et al., 2014a; Oswald & Waddell, 2015; Aso et al., 2014b; Aso & Rubin, 2016), enabling temporally staged memory formation and retrieval. Inspired by this biological design, we equip each expert in FlyPrompt with a *temporal ensemble of output heads*, implemented using exponential moving averages (EMA) with varying decay rates.

Concretely, instead of using only one shadow head in naïve EMA, each expert E_t in FlyPrompt maintains a set of n EMA heads with decay rates $\{\alpha_j\}_{j=1}^n$, where $\alpha_j \neq \alpha_k$ for all $j \neq k$. Let the online head be parameterized as $\psi = (\mathbf{W}, \mathbf{b}) \in \mathbb{R}^{|\mathcal{Y}| \times d} \times \mathbb{R}^{|\mathcal{Y}|}$, and the j -th EMA head of expert E_t as $(\mathbf{W}_t^{(j)}, \mathbf{b}_t^{(j)})$. When a new task t begins, its prompt \mathbf{p}_t is initialized as the average of previously learned prompts as a warm start:

$$\mathbf{p}_t = \frac{1}{t-1} \sum_{i=1}^{t-1} \mathbf{p}_i \quad \text{for } t > 1,$$

with random initialization for $t = 1$. This average-prompt warm start provides a more informed initialization under single-pass GCL streams, where each expert only observes limited data, and empirically accelerates convergence and more compatible with blurry boundaries in which classes can reoccur across sessions. The EMA heads of E_t are initialized as clones of the current online head. During training, only online head ψ and prompt \mathbf{p}_t are updated using the cross-entropy:

$$\mathcal{L}_t(\mathbf{x}, \mathbf{y}) = \text{CE}(\mathbf{f}_\theta(\mathbf{x}; \mathbf{p}_t) \mathbf{W}^\top + \mathbf{b} + \mathbf{m}, \mathbf{y}), \quad (6)$$

where $\mathbf{m} \in \mathbb{R}^{|\mathcal{Y}|}$ is a non-parametric logit mask initialized for each data batch (\mathbf{X}, \mathbf{y}) . We set $m_c = 0$ and for any class $c \in \mathbf{y}$ encountered in the current batch, and set $m_{c'} = -\infty$ to suppress predictions on unseen labels $c' \notin \mathbf{y}$. This masking strategy mitigates interference from class imbalance both across and within tasks (Moon et al., 2023; Kang et al., 2025), evaluated in Tab. 15.

After each update step, the EMA heads are updated as:

$$\mathbf{W}_t^{(j)} \leftarrow \alpha_j \mathbf{W}_t^{(j)} + (1 - \alpha_j) \mathbf{W}, \quad \mathbf{b}_t^{(j)} \leftarrow \alpha_j \mathbf{b}_t^{(j)} + (1 - \alpha_j) \mathbf{b}. \quad (7)$$

At inference, the REAR module first selects an expert $e = \hat{E}(\mathbf{x})$. Using the associated prompt \mathbf{p}_e , we compute logits from the online and EMA heads:

$$\mathbf{z}^{(0)} = \mathbf{f}_\theta(\mathbf{x}; \mathbf{p}_e) \mathbf{W}^\top + \mathbf{b}, \quad (8)$$

$$\mathbf{z}^{(j)} = \mathbf{f}_\theta(\mathbf{x}; \mathbf{p}_e) \mathbf{W}_e^{(j)\top} + \mathbf{b}_e^{(j)}, \quad \forall j \in \{1, \dots, n\}. \quad (9)$$

We then ensemble all $n + 1$ heads by computing the SoftMax of each and taking their element-wise maximum, followed by logit masking:

$$\hat{\mathbf{z}}(\mathbf{x}) = \max_{j \in \{0, \dots, n\}} \text{softmax}(\mathbf{z}^{(j)} + \mathbf{m}), \quad \hat{\mathbf{y}}(\mathbf{x}) = \arg \max_c \hat{z}_c(\mathbf{x}). \quad (10)$$

This temporal ensemble mechanism enables FlyPrompt to integrate stable, long-term information via EMA heads while preserving rapid adaptation through the online head, mirroring biological memory consolidation and facilitating robust inference under non-stationary, imbalanced streams. Here, we also present a theoretical guarantee that supports the use of multiple EMA heads in GCL.

Table 1: Overall performance of representative methods over three GCL benchmarks across PTMs.

PTM	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Sup-21K	Seq FT	19.71 \pm 3.39	10.42 \pm 4.92	7.51 \pm 3.94	2.29 \pm 0.85	3.47 \pm 0.41	1.49 \pm 0.42
	Linear Probe	49.69 \pm 6.09	23.07 \pm 7.33	29.24 \pm 1.26	16.87 \pm 3.14	28.96 \pm 2.46	17.33 \pm 3.08
	Seq FT w/ SL	64.90 \pm 7.18	62.06 \pm 1.89	47.20 \pm 1.47	39.60 \pm 2.43	56.16 \pm 4.32	56.50 \pm 3.08
	L2P	76.23 \pm 2.73	79.11 \pm 1.43	44.40 \pm 1.03	42.03 \pm 1.72	64.30 \pm 2.18	61.42 \pm 2.13
	DualPrompt	76.04 \pm 3.32	76.62 \pm 0.74	46.13 \pm 1.94	40.80 \pm 1.04	65.03 \pm 2.24	62.43 \pm 1.78
	CODA-P	79.13 \pm 3.06	<u>80.91</u> \pm 0.70	<u>51.87</u> \pm 2.81	<u>48.09</u> \pm 2.75	<u>66.01</u> \pm 2.20	<u>62.90</u> \pm 2.46
	MVP	67.74 \pm 4.96	63.22 \pm 0.69	39.50 \pm 1.41	32.63 \pm 3.95	54.69 \pm 3.14	50.07 \pm 3.86
	MISA	<u>80.35</u> \pm 2.39	80.75 \pm 1.24	51.52 \pm 2.09	45.08 \pm 1.43	65.40 \pm 3.01	60.20 \pm 1.82
	FlyPrompt (Ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88
Sup-21K/1K	L2P	63.88 \pm 7.79	68.96 \pm 7.63	47.10 \pm 1.21	42.22 \pm 1.94	42.96 \pm 4.13	45.00 \pm 3.83
	DualPrompt	68.02 \pm 2.08	67.04 \pm 5.84	<u>52.80</u> \pm 1.21	47.39 \pm 1.60	<u>46.80</u> \pm 2.89	<u>46.39</u> \pm 2.76
	CODA-P	<u>69.29</u> \pm 2.52	<u>69.47</u> \pm 7.19	51.20 \pm 1.76	44.30 \pm 1.50	44.66 \pm 2.73	45.18 \pm 4.50
	MVP	64.69 \pm 3.77	51.29 \pm 7.56	48.99 \pm 2.01	38.12 \pm 5.20	44.10 \pm 2.81	33.97 \pm 9.62
	MISA	62.91 \pm 7.96	67.99 \pm 7.41	50.87 \pm 1.69	<u>47.75</u> \pm 2.87	42.76 \pm 2.33	44.05 \pm 1.94
	FlyPrompt (Ours)	78.48 \pm 1.31	80.39 \pm 3.54	62.01 \pm 2.32	56.55 \pm 3.94	54.45 \pm 4.67	55.50 \pm 3.55
iBOT-21K	L2P	56.82 \pm 8.42	<u>67.61</u> \pm 8.76	35.97 \pm 1.62	36.95 \pm 2.44	14.76 \pm 1.53	<u>24.51</u> \pm 4.82
	DualPrompt	<u>66.06</u> \pm 4.52	67.14 \pm 8.60	42.48 \pm 1.62	35.91 \pm 0.88	19.90 \pm 3.68	21.84 \pm 2.35
	CODA-P	62.13 \pm 7.17	63.38 \pm 9.98	<u>45.50</u> \pm 1.66	<u>39.44</u> \pm 1.35	17.72 \pm 5.33	20.82 \pm 7.66
	MVP	62.33 \pm 3.06	48.32 \pm 11.42	41.55 \pm 1.98	29.29 \pm 5.03	<u>28.73</u> \pm 3.18	23.62 \pm 9.51
	MISA	65.30 \pm 2.28	67.43 \pm 6.75	40.94 \pm 1.22	36.16 \pm 1.58	18.62 \pm 3.36	23.66 \pm 2.21
	FlyPrompt (Ours)	75.58 \pm 1.70	79.36 \pm 3.47	57.75 \pm 2.12	54.39 \pm 1.29	28.86 \pm 5.84	36.79 \pm 7.58
iBOT-1K	L2P	53.17 \pm 7.08	<u>62.28</u> \pm 8.19	38.29 \pm 2.65	39.86 \pm 0.95	19.20 \pm 2.21	31.21 \pm 5.24
	DualPrompt	52.39 \pm 3.21	53.56 \pm 6.10	45.76 \pm 1.63	39.19 \pm 0.65	29.32 \pm 3.15	30.53 \pm 5.33
	CODA-P	<u>59.29</u> \pm 4.03	61.30 \pm 6.73	<u>49.56</u> \pm 1.57	<u>42.64</u> \pm 2.78	27.57 \pm 2.83	33.61 \pm 4.52
	MVP	57.52 \pm 3.62	44.08 \pm 12.42	44.76 \pm 2.23	34.93 \pm 4.48	<u>33.81</u> \pm 3.50	26.32 \pm 9.97
	MISA	54.31 \pm 2.91	55.89 \pm 5.10	43.91 \pm 3.95	40.09 \pm 1.24	27.76 \pm 2.69	<u>33.74</u> \pm 2.11
	FlyPrompt (Ours)	70.14 \pm 1.76	74.84 \pm 4.26	61.50 \pm 1.66	57.18 \pm 1.36	38.75 \pm 5.72	45.00 \pm 4.19
DINO-1K	L2P	47.98 \pm 7.38	<u>59.13</u> \pm 6.32	35.81 \pm 1.37	36.58 \pm 1.31	21.18 \pm 2.01	32.47 \pm 6.10
	DualPrompt	52.12 \pm 4.01	55.71 \pm 6.11	43.03 \pm 1.12	35.40 \pm 1.40	27.80 \pm 4.21	29.49 \pm 4.24
	CODA-P	<u>54.69</u> \pm 4.49	58.91 \pm 5.43	<u>45.16</u> \pm 2.05	<u>38.23</u> \pm 2.02	29.22 \pm 2.97	31.85 \pm 7.47
	MVP	53.64 \pm 3.91	41.02 \pm 12.09	41.78 \pm 2.15	32.00 \pm 4.22	<u>33.44</u> \pm 3.43	26.02 \pm 10.29
	MISA	52.03 \pm 3.07	55.98 \pm 4.26	41.26 \pm 3.25	37.50 \pm 1.62	27.13 \pm 3.31	<u>33.08</u> \pm 4.10
	FlyPrompt (Ours)	65.92 \pm 2.74	72.66 \pm 4.52	57.29 \pm 2.40	54.72 \pm 1.89	37.38 \pm 5.86	44.66 \pm 2.35
MoCo v3-1K	L2P	28.17 \pm 7.08	39.07 \pm 11.31	17.43 \pm 1.71	16.27 \pm 5.43	12.42 \pm 2.31	20.00 \pm 7.36
	DualPrompt	53.33 \pm 4.65	58.20 \pm 7.73	36.69 \pm 1.74	30.24 \pm 1.94	19.88 \pm 3.35	21.93 \pm 4.30
	CODA-P	53.47 \pm 3.42	58.55 \pm 7.19	<u>39.89</u> \pm 2.71	31.72 \pm 4.86	20.09 \pm 2.52	24.10 \pm 6.48
	MVP	54.33 \pm 4.56	40.84 \pm 14.21	36.45 \pm 2.35	26.37 \pm 6.04	28.48 \pm 3.34	23.56 \pm 9.78
	MISA	<u>57.00</u> \pm 6.06	<u>62.18</u> \pm 3.94	38.85 \pm 4.27	<u>33.47</u> \pm 0.95	25.02 \pm 4.39	<u>27.68</u> \pm 4.35
	FlyPrompt (Ours)	64.12 \pm 5.18	71.51 \pm 8.48	52.32 \pm 1.50	49.06 \pm 1.35	<u>27.92</u> \pm 4.53	33.32 \pm 3.58

Theorem 2 (TE², informal). *For an EMA head with decay α and window $L = 1/(1 - \alpha)$, the parameter error at time t satisfies*

$$\mathbb{E} \|\widetilde{\mathbf{W}}_t^{(\alpha)} - \mathbf{W}_t^*\|^2 \lesssim \zeta^2/L + (LP_t)^2,$$

where ζ^2 bounds the online noise and P_t measures drift. A geometric EMA bank contains, at every time, a head that achieves a near-optimal bias-variance trade-off up to a constant factor.

Interpretation. The bound decomposes the parameter error into a variance term $O(\zeta^2/L)$, controlled by the effective window size, and a drift-induced bias term $O((LP_t)^2)$, which increases with nonstationarity. Larger L reduces variance but increases bias, creating a bias-variance trade-off. A geometric bank of EMA windows ensures that, at any time, one head is near the optimal trade-off for the current drift level. Intuitively, when the input stream contains segments with varying temporal dynamics, such as sudden shifts at session transitions or gradual changes within each task, different EMA heads can align better with different segments, leading to more adaptive predictions. In practice, two EMA heads with windows of 10 and 100 ($\alpha = 0.9, 0.99$) are sufficient (see Sec. 4.2).

4 EXPERIMENT

In this section, we first introduce the experiment setups and then present the experiment results.

Table 2: Comparison with prominent offline methods on three GCL benchmarks under Sup-21K.

Method	CIFAR-100		ImageNet-R		CUB-200	
	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
S-Prompt++	80.21 \pm 2.55	83.48 \pm 1.20	52.14 \pm 1.65	49.13 \pm 1.60	66.61 \pm 2.21	64.73 \pm 2.25
HiDe-Prompt	77.10 \pm 3.81	81.77 \pm 2.00	53.77 \pm 1.09	49.87 \pm 3.01	67.05 \pm 2.37	67.12 \pm 0.50
HiDe-LoRA	80.07 \pm 2.41	82.00 \pm 1.25	55.09 \pm 1.45	51.29 \pm 6.29	67.26 \pm 1.76	67.28 \pm 1.45
HiDe-Adapter	79.52 \pm 2.81	81.41 \pm 0.95	53.92 \pm 1.32	50.86 \pm 5.08	66.09 \pm 1.41	64.53 \pm 1.78
NoRGa	78.89 \pm 3.33	83.03 \pm 1.20	54.12 \pm 1.37	50.09 \pm 3.66	67.16 \pm 2.44	67.06 \pm 0.58
SD-LoRA	79.26 \pm 2.21	78.91 \pm 2.48	55.51 \pm 1.30	51.97 \pm 3.09	64.12 \pm 2.02	60.57 \pm 0.77
FlyPrompt (Ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88

4.1 EXPERIMENT SETUP

Benchmarks. We evaluate FlyPrompt under the Si-Blurry GCL setting (Moon et al., 2023; Kang et al., 2025) using three representative benchmarks: CIFAR-100 (Krizhevsky et al., 2009) (60K images, 100 classes), ImageNet-R (Hendrycks et al., 2021) (30K images, 200 classes), and CUB-200 (Wah et al., 2011) (12K images, 200 fine-grained classes). Unless specified otherwise, we adopt the default Si-Blurry configuration with disjoint class ratio $r_D = 50\%$ and blurry sample ratio $r_B = 10\%$, trained over five sessions. We report two widely used metrics: average anytime accuracy A_{auc} (evaluated every 1000 batches) and final accuracy A_{last} (measured after all sessions) (Koh et al., 2021). Additional experiments of different (r_D, r_B) and online CL are provided in Appendix F.1. Unless specified, all results are averaged over five runs (\pm standard deviation) with different seeds.

Baselines. We compare FlyPrompt against a diverse set of CL and GCL methods: (1) lower-bound baselines such as sequential fine-tuning (Seq FT, including the version with a slow learning rate, SL) (Zhang et al., 2023) and linear probing; (2) prompt-based CL baselines including L2P (Wang et al., 2022d), DualPrompt (Wang et al., 2022c) and CODA-P (Smith et al., 2023); (3) state-of-the-art GCL methods such as MVP (Moon et al., 2023) and MISA (Kang et al., 2025); (4) prominent offline CL methods S-Prompt++ (Wang et al., 2022b), HiDe (Wang et al., 2023a), NoRGa (Le et al., 2024) and SD-LoRA (Wu et al., 2025) in Tabs. 2 and 12. We also implement the online version of the analytic baseline RanPAC (McDonnell et al., 2024) in Tab. 3 and other variants in Tab. 13.

Implementation. We adopt the ViT-B/16 backbone pretrained on ImageNet-21K and ImageNet-1K, including strong supervised paradigms Sup-21K, Sup-21K/1K (Sup-21K fine-tuned on ImageNet-1K) (Ridnik et al., 2021; Dosovitskiy et al., 2020), and self-supervised paradigms iBOT-21K, iBOT-1K (Zhou et al., 2021), DINO-1K (Caron et al., 2021), and MoCo v3-1K (Chen et al., 2021). We set the projection dimension $M = 10^4$, and λ based on checkpoints: 10^4 (Sup-21K), 10^6 (Sup-21K/1K, MoCo), and 10^7 (iBOT, DINO). We use $n = 2$ EMA heads with decay rates 0.9, 0.99. More implementation details of baseline methods and GCL benchmark setup can be found in Appendix C.

4.2 EXPERIMENT RESULTS

Overall Performance. Tab. 1 summarizes GCL performance across all benchmarks. Prompt-based CL methods perform well with supervised backbones (e.g., Sup-21K, Sup-21K/1K), but degrade significantly under self-supervised ones (e.g., DINO, MoCo v3-1K), particularly on fine-grained benchmarks CUB-200. This highlights the challenge of extracting discriminative features without strong pretraining priors. MVP, which incorporates contrastive learning for improved expert selection, outperforms others under the fine-grained benchmark and self-supervised PTMs, reinforcing the importance of prompt routing. However, the contrastive loss yields limited performance gains in other cases due to the absence of a replay buffer. Fig. 1 presents the anytime accuracy during GCL. MISA benefits from stronger prompt initialization and achieves relatively higher performance at the early stage, but steadily declines due to parameter overwriting, eventually matching weaker baselines like CODA-P. This suggests that while good initialization helps, it alone is insufficient for sustained GCL performance. In contrast, **FlyPrompt** consistently outperforms all baselines across datasets and PTMs. It achieves up to 11.23%, 12.43%, and 7.62% improvements in A_{auc} ; 13.53%, 16.49%, and 12.28% in A_{last} on CIFAR-100, ImageNet-R, and CUB-200, respectively. As shown in Fig. 1, FlyPrompt maintains stable, high accuracy throughout GCL, with minimal drops during session transitions. This results confirm FlyPrompt as a new state-of-the-art for GCL.

Table 3: Ablation study of different components in FlyPrompt. “—” indicates not applicable.

PTM	FlyPrompt Components			CIFAR-100		ImageNet-R	
	REAR	Prompt Expert	EMA head	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Sup-21K	RP-Based Analytic Classifier (RanPAC [†])			69.91 \pm 3.88	79.92 \pm 0.07	47.29 \pm 0.70	47.33 \pm 0.12
	—	×	×	71.33 \pm 2.17	73.22 \pm 1.63	41.73 \pm 0.97	37.33 \pm 1.71
	—	×	✓	71.69 \pm 2.27	73.30 \pm 1.55	42.50 \pm 1.01	38.35 \pm 1.01
	×	✓	×	80.75 \pm 1.98	83.65 \pm 1.94	54.91 \pm 1.32	52.58 \pm 1.36
	×	✓	✓	82.17 \pm 2.07	83.75 \pm 1.86	55.90 \pm 1.37	53.65 \pm 0.92
	✓	✓	×	81.90 \pm 2.20	84.23 \pm 1.32	55.76 \pm 1.32	52.76 \pm 1.30
	✓	✓	✓	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91
Sup-21K/1K	RP-Based Analytic Classifier (RanPAC [†])			69.76 \pm 3.33	79.49 \pm 0.16	52.91 \pm 1.07	54.79 \pm 0.22
	—	×	×	57.82 \pm 6.94	62.67 \pm 8.55	46.24 \pm 0.72	39.06 \pm 2.81
	—	×	✓	60.76 \pm 6.77	63.39 \pm 8.81	49.88 \pm 0.93	41.65 \pm 1.38
	×	✓	×	70.09 \pm 3.52	67.62 \pm 5.80	52.07 \pm 1.35	44.13 \pm 3.61
	×	✓	✓	72.51 \pm 3.15	68.66 \pm 5.96	55.55 \pm 1.51	45.90 \pm 3.53
	✓	✓	×	71.28 \pm 2.58	69.73 \pm 5.78	53.12 \pm 2.19	44.69 \pm 3.65
	✓	✓	✓	78.48 \pm 1.31	80.39 \pm 3.54	62.01 \pm 2.32	56.55 \pm 3.94

Table 4: Effect of REAR and TE² on $A_{\text{auc}}(\%)$ performance for PTM-based CL methods using CIFAR-100 under Sup-21K. The numbers in parentheses indicate the difference from the baseline, and the arrow direction indicates an increase (\uparrow) or decrease (\downarrow). See Tab. 14 for complete results.

Method	Baseline	w/ REAR	w/ TE ²	w/ Both
DualPrompt	76.04 \pm 3.32	80.63 \pm 2.25 (\uparrow 4.59)	76.83 \pm 3.44 (\uparrow 0.79)	82.33 \pm 2.17 (\uparrow 6.29)
MVP	67.74 \pm 4.96	67.44 \pm 4.89 (\downarrow 0.30)	68.91 \pm 4.86 (\uparrow 1.17)	68.93 \pm 4.60 (\uparrow 1.19)
MISA	80.35 \pm 2.39	82.03 \pm 1.97 (\uparrow 1.68)	81.65 \pm 2.24 (\uparrow 1.30)	83.60 \pm 2.08 (\uparrow 3.25)
S-Prompt++	80.21 \pm 2.55	81.43 \pm 2.45 (\uparrow 1.21)	81.93 \pm 2.21 (\uparrow 1.72)	83.11 \pm 2.30 (\uparrow 2.90)
HiDe-Prompt	77.10 \pm 3.81	78.41 \pm 2.64 (\uparrow 1.31)	77.46 \pm 3.56 (\uparrow 0.36)	78.60 \pm 2.53 (\uparrow 1.51)
NoRGa	78.89 \pm 3.33	79.37 \pm 2.71 (\uparrow 0.48)	79.16 \pm 3.28 (\uparrow 0.27)	79.37 \pm 2.74 (\uparrow 0.48)

Ablation Study. To assess the contribution of each FlyPrompt component, we conduct a comprehensive ablation study of REAR for prompt selection, multi-prompt across tasks, and TE² for EMA head ensemble. Results in Tab. 3 show that each module provides consistent gains, with the full FlyPrompt achieving the best performance. We additionally include RanPAC[†], an analytic learner using random projections but no expert modularity, to simulate REAR without multi-expert routing. While this performs competitively under limited training, it falls short without expert specialization, underscoring the importance of both routing and competence. Notably, the gain from EMA heads alone is modest unless combined with REAR and prompt modularization, highlighting the synergy among bio-inspired components rather than simple additive effects. We further integrate our REAR and TE² components into a range of strong baseline models by replacing their routing and output head modules correspondingly. Results in Tab. 4 further demonstrate the consistent improvements when either component is added (more results across datasets and metrics are presented in Tab. 14).

Hyperparameter Sensitivity. Fig. 5 explores key hyperparameters in REAR. Increasing the projection dimension M improves performance, consistent with the theory that higher-dimensional spaces enable better feature separability and router performance (in Theorem 1), mirroring sparse expansion in the fruit fly mushroom body. However, since the memory cost grows linearly with M , we set $M = 10,000$ as a practical trade-off. The regularization parameter λ has smaller impact, with performance stable across several orders of magnitude. Full results across other PTMs are provided in Appendix F.5. We further analyze EMA decay rates with temporal ensemble. Tab. 6 shows that two EMA heads of 0.9, 0.99, combined with the online head, achieve the best trade-off across datasets. This aligns with neurobiological findings that the mushroom body maintains short-, mid-, and long-term memory modules in parallel. Among various ensemble strategies (Tab. 5), the “SoftMax + element-wise maximum” method is most effective and used by default. Detailed evaluations across other PTMs and configurations are provided in Appendices F.8 and F.10.

Detailed Analysis. Returning to the core challenges identified in Sec. 2.2, we revisit the roles of expert routing and expert competence improvement. As shown in Fig. 2, methods (e.g., FlyPrompt) that improve in these two areas correlate strongly with better overall GCL performance. In particular, Fig. 2b demonstrates the impact of random projection in boosting routing accuracy, while Fig. 2c highlights remaining headroom for improving expert competence. Despite introducing RP layer and tracking feature statistics, FlyPrompt adds minimal parameter overhead, i.e., just 0.83% more

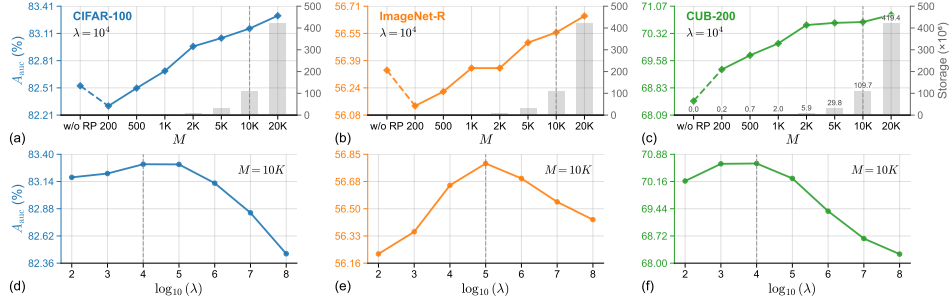


Figure 5: Analysis of hyperparameters in REAR. (a-c) Different projection dimension M with fixed $\lambda = 10^4$; we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$. Dashed lines indicate the optimal choice of each hyperparameter.

Table 5: Performance comparison of ensemble aggregation choices for TE² under Sup-21K.

Ensemble Method	CIFAR-100		ImageNet-R		CUB-200	
	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Mean	81.34 \pm 1.64	85.11 \pm 1.03	52.71 \pm 1.36	53.24 \pm 1.22	68.49 \pm 2.57	73.95 \pm 1.90
Max Prob.	82.29 \pm 2.25	84.95 \pm 1.20	55.56 \pm 1.38	53.53 \pm 1.40	68.00 \pm 2.50	66.56 \pm 1.60
Min Entropy	81.92 \pm 2.19	84.23 \pm 1.32	55.05 \pm 1.31	52.88 \pm 1.38	66.78 \pm 2.53	64.73 \pm 1.36
SoftMax+Mean	82.30 \pm 1.82	85.98 \pm 0.80	56.16 \pm 1.56	55.53 \pm 0.89	70.77 \pm 3.00	74.86 \pm 1.54
SoftMax+Max Prob.	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88
SoftMax+Min Entropy	83.11 \pm 2.34	86.50 \pm 0.64	55.94 \pm 1.41	54.24 \pm 1.34	69.86 \pm 2.80	71.51 \pm 1.79

Table 6: Performance comparison of different EMA decay rates for TE² under Sup-21K.

EMA Decay Rate	CIFAR-100		ImageNet-R	
	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Online head only	81.90 \pm 2.20	84.23 \pm 1.32	54.91 \pm 1.32	52.58 \pm 1.36
+0.9	82.81 \pm 2.28	86.36 \pm 0.54	56.36 \pm 1.52	55.09 \pm 0.89
+0.99	82.84 \pm 2.51	86.41 \pm 0.39	55.94 \pm 1.65	54.67 \pm 0.89
+0.999	81.80 \pm 2.37	84.39 \pm 0.83	55.19 \pm 1.39	53.52 \pm 0.80
+0.9,0.99	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91
+0.9,0.99,0.999	82.99 \pm 2.22	86.24 \pm 0.79	56.35 \pm 1.72	55.50 \pm 0.77

Table 7: Computational cost and overall performance using CIFAR-100 under Sup-21K.

Method	Total Param. (M, \downarrow)	Trainable Param. (M, \downarrow)	Time Delay (s/batch, \downarrow)	$A_{\text{auc}}(\%, \uparrow)$
L2P	86.01	0.22	5.17	76.23
DualPrompt	86.35	0.55	4.78	76.04
CODA-P	86.72	0.92	4.75	79.13
MVP	86.12	0.32	5.35	67.74
MISA	86.37	0.58	4.78	80.35
FlyPrompt (ours)	87.08	0.46	4.96	83.24

parameters than MISA on ViT-B/16, and incurs negligible increase in computational cost (see Tab. 7, more comprehensive comparison in Tab. 12 and detailed cost breakdown of components in Tab. 19). Together, these findings validate FlyPrompt’s effectiveness in resolving the GCL challenges.

5 CONCLUSION

We presented **FlyPrompt**, a biologically inspired framework for GCL, which addresses the core challenges of expert routing and expert competence improvement under blurred task boundaries and single-pass constraints. Grounded in the neurobiological principles of the fruit fly’s mushroom body, known for its sparse expansion, random connectivity, and multiscale modularity, FlyPrompt integrates a randomly expanded analytic router for non-iterative expert selection and a temporal ensemble of expert heads for robust adaptation over time. Theoretical analysis and empirical results across multiple GCL benchmarks demonstrate its strong performance and scalability.

While these results are encouraging, several limitations of the current work point to promising future directions. For instance, the temporal ensemble relies on a fixed composition of EMA decay rates, and adapting these dynamically to data drift could enhance robustness. Additionally, performance under extreme long-tailed distributions warrants further study. Looking forward, GCL is essential for deploying real-world learning systems, such as embodied agents, user-facing AI, and resource-constrained devices, where data is dynamic and supervision is limited. As continual adaptation is a natural strength of biological systems, the underlying principles they offer will continue to inspire future advances in GCL and beyond.

Acknowledgment. This work was supported by the STI2030-Major Projects 2022ZD0204900, the National Natural Science Foundation of China (NO. 62406160, 32530042, 32021002), the Beijing Natural Science Foundation L247011, and the Beijing Major Science and Technology Project No. Z251100008425003.

Ethics statement. I acknowledge that I and all co-authors of this work have read and commit to adhering to the ICLR Code of Ethics.

Reproducibility statement. We have included the source code with clear instructions, and will release them upon acceptance.

REFERENCES

- Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. In *CVPR*, pp. 11254–11263, 2019a.
- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *NeurIPS*, 32, 2019b.
- Yoshinori Aso and Gerald M Rubin. Dopaminergic neurons write and update memories with cell-type-specific rules. *Elife*, 5:e16135, 2016.
- Yoshinori Aso, Daisuke Hattori, Yang Yu, Rebecca M Johnston, Nirmala A Iyer, Teri-TB Ngo, Heather Dionne, LF Abbott, Richard Axel, Hiromu Tanimoto, et al. The neuronal architecture of the mushroom body provides a logic for associative learning. *elife*, 3:e04577, 2014a.
- Yoshinori Aso, Divya Sitaraman, Toshiharu Ichinose, Karla R Kaun, Katrin Vogt, Ghislain Belliard-Guérin, Pierre-Yves Plaçais, Alice A Robie, Nobuhiro Yamagata, Christopher Schnaitmann, et al. Mushroom body output neurons encode valence and guide memory-based action selection in *drosophila*. *elife*, 3:e04580, 2014b.
- Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, pp. 8218–8227, 2021.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Non-stationary stochastic optimization. *Operations research*, 63(5):1227–1244, 2015.
- Emna Bouzaiane, Severine Trannoy, Lisa Scheunemann, Pierre-Yves Plaçais, and Thomas Preat. Two independent mushroom body output circuits retrieve the six discrete components of *drosophila* aversive memory. *Cell Reports*, 11(8):1280–1292, 2015.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *NeurIPS*, 33:15920–15930, 2020.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pp. 9650–9660, 2021.
- Isaac Cervantes-Sandoval, Alfonso Martin-Peña, Jacob A Berry, and Ronald L Davis. System-like consolidation of olfactory memories in *drosophila*. *Journal of Neuroscience*, 33(23):9846–9854, 2013.
- Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, pp. 9640–9649, 2021.
- Sanjoy Dasgupta, Charles F Stevens, and Saket Navlakha. A neural algorithm for a fundamental computing problem. *Science*, 358(6364):793–796, 2017.
- Ronald L Davis. Learning and memory using *drosophila melanogaster*: a focus on advances made in the fifth decade of research. *Genetics*, 224(4):iyad085, 2023.
- Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *PAMI*, 44(7):3366–3385, 2021.

- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Dan Hendrycks, Steven Basart, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*, 2021.
- Kyle S Honegger, Robert AA Campbell, and Glenn C Turner. Cellular-resolution population imaging reveals robust sparse coding in the drosophila mushroom body. *Journal of neuroscience*, 31(33):11772–11785, 2011.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Zhiqi Kang, Liyuan Wang, Xingxing Zhang, and Karteek Alahari. Advancing prompt-based methods for replay-independent general continual learning. In *ICLR*, 2025.
- Gyuhak Kim, Changnan Xiao, Tatsuya Konishi, Zixuan Ke, and Bing Liu. A theoretical study on solving continual learning. *NeurIPS*, 35:5065–5079, 2022.
- Hyunseo Koh, Dahyun Kim, Jung-Woo Ha, and Jonghyun Choi. Online continual learning on class incremental blurry task configuration with anytime inference. *arXiv preprint arXiv:2110.10031*, 2021.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pp. 3519–3529. PMIR, 2019.
- Michael J Krashes, Alex C Keene, Benjamin Leung, J Douglas Armstrong, and Scott Waddell. Sequential use of mushroom body neuron subsets during drosophila odor memory processing. *Neuron*, 53(1):103–115, 2007.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Minh Le, An Nguyen, Huy Nguyen, Trang Nguyen, Trang Pham, Linh Van Ngo, and Nhat Ho. Mixture of experts meets prompt-based continual learning. *Advances in Neural Information Processing Systems*, 37:119025–119062, 2024.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Feng Li, Jack W Lindsey, Elizabeth C Marin, Nils Otto, Marisa Dreher, Georgia Dempsey, Ildiko Stark, Alexander S Bates, Markus William Pleijzier, Philipp Schlegel, et al. The connectome of the adult drosophila mushroom body provides insights into function. *Elife*, 9:e62576, 2020.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- Sen Lin, Li Yang, Deliang Fan, and Junshan Zhang. Beyond not-forgetting: Continual learning with backward knowledge transfer. *Advances in Neural Information Processing Systems*, 35:16165–16177, 2022.
- Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton van den Hengel. Ranpac: Random projections and pre-trained models for continual learning. *NeurIPS*, 36, 2024.
- Fei Mi, Lingjing Kong, Tao Lin, Kaicheng Yu, and Boi Faltings. Generalized class incremental learning. In *CVPR Workshops*, pp. 240–241, 2020.
- Mehrab N Modi, Yichun Shuai, and Glenn C Turner. The drosophila mushroom body: from architecture to algorithm in a learning circuit. *Annual review of neuroscience*, 43(1):465–484, 2020.

- Jun-Yeong Moon, Keon-Hee Park, Jung Uk Kim, and Gyeong-Moon Park. Online class incremental learning on stochastic blurry task boundary via mask and visual prompt tuning. In *ICCV*, pp. 11731–11741, 2023.
- David Oswald and Scott Waddell. Olfactory learning skews mushroom body output pathways to steer behavioral choice in drosophila. *Current opinion in neurobiology*, 35:178–184, 2015.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, pp. 524–540. Springer, 2020.
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. *NeurIPS*, 30, 2017.
- Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- 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*, pp. 11909–11919, 2023.
- Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- Glenn C Turner, Maxim Bazhenov, and Gilles Laurent. Olfactory representations by drosophila mushroom body neurons. *Journal of neurophysiology*, 99(2):734–746, 2008.
- Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, California Institute of Technology, 2011.
- Liyuan Wang and Qian Li. Convergent multi-modular architecture for adaptive learning in drosophila and artificial intelligence. *iScience*, 28(11), 2025.
- Liyuan Wang, Mingtian Zhang, Zhongfan Jia, Qian Li, Chenglong Bao, Kaisheng Ma, Jun Zhu, and Yi Zhong. Afec: Active forgetting of negative transfer in continual learning. *NeurIPS*, 2021.
- Liyuan Wang, Xingxing Zhang, Qian Li, Jun Zhu, and Yi Zhong. Coscl: Cooperation of small continual learners is stronger than a big one. In *ECCV*, 2022a.
- Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical decomposition of prompt-based continual learning: Rethinking obscured sub-optimality. *NeurIPS*, 2023a.
- Liyuan Wang, Xingxing Zhang, Qian Li, Mingtian Zhang, Hang Su, Jun Zhu, and Yi Zhong. Incorporating neuro-inspired adaptability for continual learning in artificial intelligence. *Nature Machine Intelligence*, 2023b.
- Liyuan Wang, Jingyi Xie, Xingxing Zhang, Hang Su, and Jun Zhu. Hide-pet: Continual learning via hierarchical decomposition of parameter-efficient tuning. *arXiv preprint arXiv:2407.05229*, 2024a.
- Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *PAMI*, 2024b.
- Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained transformers: An occam’s razor for domain incremental learning. *Advances in Neural Information Processing Systems*, 35:5682–5695, 2022b.

- 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*, pp. 631–648. Springer, 2022c.
- 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*, pp. 139–149, 2022d.
- Yichen Wu, Hongming Piao, Long-Kai Huang, Renzhen Wang, Wanhua Li, Hanspeter Pfister, Deyu Meng, Kede Ma, and Ying Wei. Sd-lora: Scalable decoupled low-rank adaptation for class incremental learning. In *ICLR*, 2025.
- Hongwei Yan, Liyuan Wang, Kaisheng Ma, and Yi Zhong. Orchestrate latent expertise: Advancing online continual learning with multi-level supervision and reverse self-distillation. In *CVPR*, pp. 23670–23680, 2024.
- Hongwei Yan, Guanglong Sun, Zhiqi Kang, Yi Zhong, and Liyuan Wang. Domain generalizable continual learning. *arXiv preprint arXiv:2510.16914*, 2025.
- Anthony Zador, Sean Escola, Blake Richards, Bence Ölveczky, Yoshua Bengio, Kwabena Boahen, Matthew Botvinick, Dmitri Chklovskii, Anne Churchland, Claudia Clopath, et al. Catalyzing next-generation artificial intelligence through neuroai. *Nature Communications*, 14(1):1597, 2023.
- Gengwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learner with classifier alignment for continual learning on a pre-trained model. *arXiv preprint arXiv:2303.05118*, 2023.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.
- Kanglei Zhou, Ruizhi Cai, Liyuan Wang, Hubert PH Shum, and Xiaohui Liang. A comprehensive survey of action quality assessment: Method and benchmark. *arXiv preprint arXiv:2412.11149*, 2024a.
- Kanglei Zhou, Liyuan Wang, Xingxing Zhang, Hubert PH Shum, Frederick WB Li, Jianguo Li, and Xiaohui Liang. Magr: Manifold-aligned graph regularization for continual action quality assessment. In *European Conference on Computer Vision*, pp. 375–392. Springer, 2024b.
- Kanglei Zhou, Zikai Hao, Liyuan Wang, and Xiaohui Liang. Adaptive score alignment learning for continual perceptual quality assessment of 360-degree videos in virtual reality. *IEEE Transactions on Visualization and Computer Graphics*, 31(5):2880–2890, 2025a.
- Kanglei Zhou, Qingyi Pan, Xingxing Zhang, Hubert PH Shum, Frederick WB Li, Xiaohui Liang, and Liyuan Wang. Continual action quality assessment via adaptive manifold-aligned graph regularization. *arXiv preprint arXiv:2510.06842*, 2025b.
- Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation and self-supervision for incremental learning. In *CVPR*, pp. 5867–5876, 2021. doi: 10.1109/CVPR46437.2021.00581.
- Huiping Zhuang, Zhenyu Weng, Hongxin Wei, Renchunzi Xie, Kar-Ann Toh, and Zhiping Lin. Acil: Analytic class-incremental learning with absolute memorization and privacy protection. *NeurIPS*, 35:11602–11614, 2022.
- Huiping Zhuang, Yizhu Chen, Di Fang, Run He, Kai Tong, Hongxin Wei, Ziqian Zeng, and Cen Chen. Gacil: Exemplar-free generalized analytic continual learning. *NeurIPS*, 37:83024–83047, 2024.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pp. 928–936, 2003.

A LARGE LANGUAGE MODELS ASSISTANCE

Large language models were used to polish the manuscript. The authors have thoroughly reviewed and edited all content and take full responsibility for the published work.

B RELATED WORK

Continual Learning (CL) aims to train models on task streams with evolving distributions (Wang et al., 2024b; 2021; 2022a; Parisi et al., 2019). Canonical CL settings are categorized into task-incremental learning (TIL), class-incremental learning (CIL), and domain-incremental learning (DIL) (Van de Ven & Tolias, 2019; Yan et al., 2024; 2025), depending on the structure of input and label spaces. TIL and CIL assume disjoint label spaces, with task identity provided only in TIL, while DIL shares label space but varies input domains. Theoretically, CL has been formalized by decoupling task identity prediction (TIP) and within-task prediction (WTP), which remain orthogonal under clearly segmented tasks and from-scratch training (Kim et al., 2022; Wang et al., 2023a).

The rise of pretrained models (PTMs) has shifted CL towards adapting frozen backbones via lightweight modules, known as parameter-efficient tuning (PET) (Lester et al., 2021; Li & Liang, 2021; Rebuffi et al., 2017; Hu et al., 2021). PET-based CL methods often employ either task-shared modules (Zhang et al., 2023; McDonnell et al., 2024; Zhou et al., 2024a; 2025b;a; 2024b) that require gradual updates, or task-specific experts (Wang et al., 2022d; 2023a) that demand effective expert selection (implicitly via external queries (Wang et al., 2022d) or explicitly via routing functions (Wang et al., 2023a)). Importantly, the strong priors embedded in PTMs blur the TIP-WTP decomposition, making classical CL theory less applicable (Wang et al., 2023a; 2024a).

General Continual Learning (GCL) extends CL to more practical scenarios by removing assumptions of clear task segmentation and offline data access (Buzzega et al., 2020; De Lange et al., 2021; Mi et al., 2020). Specifically, GCL emphasizes *online learning*, where each data point is seen only once; and *blurry or unknown task boundaries*, where task identities are absent or ill-defined (Aljundi et al., 2019a; Prabhu et al., 2020; Bang et al., 2021; Moon et al., 2023). These properties introduce unique challenges in expert selection, knowledge retention, and fast adaptation, without task identities or replay buffers. Additional constraints, such as constant memory budgets and anytime inference, further distinguish GCL from traditional CL (De Lange et al., 2021).

To implement the GCL challenges, benchmarks such as Task-Free CL (Aljundi et al., 2019a; Prabhu et al., 2020) and Si-Blurry (Moon et al., 2023) have been proposed, progressively relaxing task-awareness and enforcing stream-based learning. Correspondingly, GCL methods adapt replay-based sampling (Aljundi et al., 2019b; Bang et al., 2021), memory management (Koh et al., 2021), or PET-based designs (Moon et al., 2023; Kang et al., 2025). However, replay methods raise privacy and scalability concerns, while recent PET-based methods (e.g., MVP (Moon et al., 2023) and MISA (Kang et al., 2025)) still suffer from limited representation capacity and lack principled mechanisms for prompt expert selection under non-stationary inputs. Consequently, their improvements over naive PTM-based baselines remain modest.

C IMPLEMENTATION DETAILS

C.1 TRAINING SETUP

We follow the previous GCL studies (Moon et al., 2023; Kang et al., 2025) for a fair comparison. The standard ViT-B/16 transformer backbone has an embedding dimension of $d = 768$. For prompt-based methods, we unify the prompt length to 5 and the position to insert the prompt as the first five layers of ViT. All methods are trained with an Adam optimizer with a learning rate 0.005 and zero weight decay. We set the batch size to 64, the epoch number to 1 (online learning), and the online iteration of each batch to 3. All images are cropped and resized to 224×224 to fit the ViT format using standard data transformation operations (resize, random crop, random horizontal flip and normalization). Moreover, the logit mask m trick in Sec. 3.2 is generally applied to all methods to enhance training stability. All experiment jobs are performed on the same Linux server with Intel Xeon Silver 4316 2.3GHz CPUs (20 cores), 1 NVIDIA RTX 4090 GPU. For random seeds, we use the fixed values 1, 2, 3, 4, and 5 for all parallel runs.

C.2 BASELINES

Unless otherwise specified, all baselines share the common ViT-B/16 backbone, single-pass Si-Blurry streams, optimizer, and data preprocessing described at the beginning of this section; below, we highlight method-specific architectures and the GCL-specific adaptations.

Sequential fine-tuning (Seq FT / SL) and Linear probing. Seq FT fine-tunes all parameters of the ViT-B/16 backbone and classifier on the Si-Blurry streams without replay buffers or task-specific heads; SL is an otherwise identical variant with a smaller learning rate (i.e., 5×10^{-5} , 10 times smaller than 0.005 used by other baselines) to provide a more optimistic lower bound (Zhang et al., 2023). Linear probing instead freezes the backbone and trains only a linear classifier, with no prompt modules or expert structures; together, these methods serve as simple PTM-based lower bounds.

Prompt-based CL baselines (L2P, DualPrompt, CODA-Prompt). For L2P (Wang et al., 2022d), DualPrompt (Wang et al., 2022c), and CODA-Prompt (Smith et al., 2023), we keep their original prompt-controller designs (key-based prompt pool in L2P, global+task prompts in DualPrompt, and attention-based prompts in CODA-Prompt), but adapt them to GCL by freezing the ViT-B/16 backbone and unifying prompt length and insertion position as in Sec. C.1. All three methods are trained in a single online pass over the Si-Blurry streams with the same update schedule as FlyPrompt, without extra replay or offline fine-tuning, so that differences in performance come from their prompt mechanisms rather than from additional data passes.

GCL baselines (MVP, MISA). MVP (Moon et al., 2023) and MISA (Kang et al., 2025) are implemented on top of the same ViT-B/16 backbone and Si-Blurry streams. We follow their official configurations for expert/prompt structures and initialization, while enforcing the unified prompt configuration and online training protocol of Sec. C.1. As in prior work, they maintain session-wise experts or prompts, but do not use any privileged task oracle beyond the evolving data stream; their routing structures can be interpreted in the same way as FlyPrompt’s experts discussed in Sec. D. For fairness, MVP and MISA also use the batch-seen class logit mask in Sec. 3.2, whose effect is ablated alongside other mask types in Tab. 15.

Offline PTM-based CL methods (S-Prompt++, HiDe-Prompt/LoRA/Adapter, NoRGa, SD-LoRA). S-Prompt++ (Wang et al., 2022b) introduces prompt experts with a mixture-of-experts (MoE) structure with linear gating, while HiDe-Prompt/LoRA/Adapter (Wang et al., 2023a) and NoRGa (Le et al., 2024) build hierarchical decompositions and stronger MoE-based routing on top of S-Prompt++, and SD-LoRA Wu et al. (2025) leverages structured low-rank adapters by decomposing expert LoRA into learnable amplitude and fixed direction; all of these are originally designed for offline or task/class-incremental CL. Specifically, HiDe and NoRGa consist of a two-stage TIP+WTP (task-ID prediction then within-task prediction) pipeline. To make them compatible with GCL and Si-Blurry, we adapt their TIP step as follows: when the method predicts a class, it is allowed to activate *all* prompts corresponding to the candidate task IDs associated with that class, and we count the prediction as correct if *any* activated prompt outputs the true label. This is an intentionally favorable modification for these baselines. In addition, any feature statistics required by their alignment modules (e.g., for HiDe or NoRGa) are accumulated online from the stream, rather than being computed from stored per-task datasets. Quantitative comparisons of these adapted offline PTM-based methods with FlyPrompt are reported in Tab. 12.

Analytic random-projection baselines (RanPAC variants). RanPAC (McDonnell et al., 2024) was originally proposed for offline class-incremental learning, where the PTM is fine-tuned on the first task, frozen, and all task-1 features are recomputed and stored to form a stable Gram matrix for closed-form ridge regression; this protocol is incompatible with single-pass GCL and blurry task boundaries. To ensure a fair analytic baseline, we adapt RanPAC into three GCL-compliant variants: **RanPAC[†]** fine-tunes the PTM on the first Si-Blurry session without storing features and then solves a closed-form ridge classifier on the resulting random-feature representations, serving as the main analytic random-projection baseline in our ablations; **RanPAC[‡]** freezes the PTM during the first session and stores all features to approximate the original offline setting while still respecting the single-pass constraint on labels; **RanPAC^{*}** simultaneously fine-tunes the PTM and collects features during the first session, which yields an ill-conditioned Gram matrix due to representation drift but offers an optimistic upper bound on analytic-classifier performance under our setting. As summarized in Tab. 13, FlyPrompt consistently outperforms all RanPAC variants across datasets.

D DISCUSSION OF TASK BOUNDARY IN SI-BLURRY

Table 8: Performance comparison of different numbers of tasks and experts for GCL methods on CIFAR-100 dataset over Sup-21K. “task-correlated” indicates that the initialization and training of expert parameters are aligned with task/sessions. w denotes the sample budget (window size) of each expert when methods adopt a self-triggered expert allocation mechanism. All results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds.

# of Tasks	# of Experts	MVP		MISA		FlyPrompt	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
5	5 (task-correlated)	<u>67.74</u> ± 4.96	63.22 ± 0.69	<u>80.35</u> ± 2.39	80.75 ± 1.24	83.24 ± 2.23	86.76 ± 0.73
	5 ($w = 10000$)	67.75 ± 4.96	63.22 ± 0.76	80.60 ± 2.06	81.73 ± 1.17	83.67 ± 2.23	<u>85.78</u> ± 0.62
	10 ($w = 5000$)	67.23 ± 5.06	<u>63.47</u> ± 0.78	79.95 ± 1.86	<u>81.32</u> ± 1.15	<u>83.40</u> ± 2.28	85.43 ± 0.32
	20 ($w = 2500$)	67.19 ± 4.80	64.06 ± 1.48	79.94 ± 1.75	81.03 ± 0.92	82.26 ± 1.94	84.48 ± 0.64
10	10 (task-correlated)	58.23 ± 3.42	61.13 ± 6.00	75.71 ± 3.10	80.22 ± 0.47	77.65 ± 3.05	84.87 ± 0.50
	5 ($w = 10000$)	58.49 ± 3.57	<u>61.98</u> ± 6.19	76.25 ± 3.10	80.62 ± 0.62	<u>77.28</u> ± 2.79	<u>84.63</u> ± 0.50
	10 ($w = 5000$)	58.19 ± 3.42	60.80 ± 6.48	75.69 ± 3.11	80.18 ± 0.41	76.96 ± 3.23	84.03 ± 0.39
	20 ($w = 2500$)	<u>58.46</u> ± 3.34	62.15 ± 5.28	<u>75.73</u> ± 2.80	80.66 ± 0.61	76.47 ± 3.38	83.39 ± 0.55
20	20 (task-correlated)	<u>56.52</u> ± 3.20	56.87 ± 5.18	<u>73.96</u> ± 0.72	77.98 ± 1.17	75.87 ± 1.93	<u>81.98</u> ± 1.12
	5 ($w = 10000$)	56.59 ± 3.34	56.58 ± 6.76	74.27 ± 0.97	<u>77.89</u> ± 1.60	76.12 ± 1.57	81.90 ± 0.32
	10 ($w = 5000$)	56.35 ± 3.38	56.80 ± 6.57	73.86 ± 0.82	77.59 ± 1.42	76.12 ± 1.21	82.21 ± 0.80
	20 ($w = 2500$)	56.48 ± 3.15	<u>56.86</u> ± 5.39	73.66 ± 0.85	77.70 ± 1.24	75.36 ± 1.65	81.13 ± 0.55

GCL (Buzzega et al., 2020) is defined by a single-pass, non-stationary data stream without task boundaries during training and without a task oracle at test time. The Si-Blurry benchmark (Moon et al., 2023) that we adopt has been carefully analyzed in subsequent work (Kang et al., 2025): by controlling the disjoint-class ratio r_D and blurry-sample ratio r_B , it generates streams where (i) the number of active classes can vary across sessions, (ii) classes may reoccur across sessions, and (iii) the number of samples per class and per session is randomized (Mi et al., 2020). In particular, when r_D approaching 0, the nominal “task” or “session” index becomes decorrelated from distributional changes. These properties ensure that Si-Blurry conforms to the core GCL assumptions, rather than reducing to standard task-incremental CIL.

Within this setting, FlyPrompt does not assume any privileged boundary information beyond what is already used by prior GCL methods such as MVP and MISA. The “task” or “session” index provided by Si-Blurry is treated as a conceptual device to describe how the benchmark constructs streams, not as a supervision signal for the model. In our implementation, expert indices are aligned with nominal session identities purely for convenience: the same behavior can be reproduced by starting a new expert after a fixed number of observed samples or when a user-defined computational/storage budget is reached, without accessing the task index. Moreover, the total number of experts T is not a hard-coded prior; matrices such as $Q \in \mathbb{R}^{M \times T}$ and the router head can be dynamically extended from T to $T + 1$ via zero-padding, analogous to adding classes in a standard classifier. Implementation details for how both GCL methods and offline PTM-based baselines are instantiated under this regime are summarized in Appendix C.2.

To empirically validate that FlyPrompt and comparable GCL baselines do not gain an advantage from Si-Blurry’s session structure, we further compare task-aligned expert management with a self-triggered expert allocation mechanism. In the self-triggered setup, each method maintains a fixed sample budget and freezes the current expert while initializing a new one whenever the number of observed samples reaches a predefined threshold, fully decoupling expert updates from external task segmentation. We evaluate multiple combinations of nominal task counts (# of Tasks = 5, 10, 20) and expert budgets (# of Experts = 5, 10, 20) on CIFAR-100, with corresponding sample budgets chosen to cover the 50K training examples. As summarized in Tab. 8, self-triggered expert initialization achieves performance on par with, or slightly better than, session-aligned setups for MVP, MISA, and FlyPrompt across all tested budgets. This confirms that (i) task-switching signals offer no measurable benefit in this benchmark, and (ii) our expert management mechanism does not exploit any extra boundary information.

E THEORETICAL PROOFS

E.1 PROOFS FOR REAR (THEOREM 1)

NOTATION AND ASSUMPTIONS

We repeat and fix the notation used throughout the proof of Theorem 1:

1. $f_\theta : \mathcal{X} \rightarrow \mathbb{R}^d$ is the given pretrained backbone. For an input \mathbf{x} we write $\mathbf{h} = f_\theta(\mathbf{x}) \in \mathbb{R}^d$.
2. $\mathbf{R} \in \mathbb{R}^{d \times M}$ is a random matrix with i.i.d. $\mathcal{N}(0, 1)$ entries; the j -th column is $\mathbf{r}_j \in \mathbb{R}^d$.
3. The random-expanded feature map is

$$\boldsymbol{\varphi}(\mathbf{x}) = (\varphi_1(\mathbf{x}), \dots, \varphi_M(\mathbf{x}))^\top, \quad \varphi_j(\mathbf{x}) = \sigma(\mathbf{h}^\top \mathbf{r}_j).$$

4. Assume embedding-boundedness: $\|\mathbf{h}\|_2 \leq H$ for all \mathbf{x} . (This can be enforced in practice by layer-norm or clipping.)
5. Activation $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is L_σ -Lipschitz and has linear growth $|\sigma(z)| \leq C(1 + |z|)$. ReLU satisfies these with $L_\sigma = 1$ and linear growth $C = 1$.
6. For training, we accumulate batches (or singletons) to form

$$\mathbf{G} = \sum_{i=1}^N \boldsymbol{\varphi}(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x}_i)^\top \in \mathbb{R}^{M \times M}, \quad \mathbf{Q} = \sum_{i=1}^N \boldsymbol{\varphi}(\mathbf{x}_i) \mathbf{c}_i^\top \in \mathbb{R}^{M \times T},$$

where $\mathbf{c}_i \in \{0, 1\}^T$ is the one-hot indicator of the target expert; see Eq. (2)).

7. Ridge solution (router):

$$\widehat{\mathbf{U}}^\top = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{Q}, \quad \lambda > 0,$$

with regularization parameter λ as in Eq. (4).

8. For the theoretical analysis, we treat the training pairs (\mathbf{x}_i, y_i) as i.i.d. draws from an underlying distribution over $\mathcal{X} \times \{1, \dots, T\}$, with a fixed and finite number of experts T .
9. For the margin-based routing-accuracy corollary below, we additionally assume that there exists a margin $\gamma > 0$ such that, for the population minimizer U^* and almost every input \mathbf{x} , the score of the correct expert $t^*(\mathbf{x})$ satisfies $s_{U^*}(\mathbf{x})_{t^*(\mathbf{x})} \geq s_{U^*}(\mathbf{x})_t + \gamma$ for all $t \neq t^*(\mathbf{x})$, where $s_U(\mathbf{x}) := \boldsymbol{\varphi}(\mathbf{x}) U^\top$.²

At the population level we consider the regularized squared risk

$$\mathcal{R}(U) := \mathbb{E} \|\mathbf{s}_U(X) - C\|_2^2 + \lambda \|U\|_F^2,$$

where (X, C) denotes a random variable pair drawn from the same distribution as the training examples $(\mathbf{x}_i, \mathbf{c}_i)$, with $X \in \mathcal{X}$ and $C \in \{0, 1\}^T$ is the one-hot indicator of the target expert. We write $s_U(\mathbf{x}) := \boldsymbol{\varphi}(\mathbf{x}) U^\top$ for the router scores as in Eq. (5). The minimizer of \mathcal{R} in the kernel-induced feature space is precisely the U^* appearing below.

We then state the complete Theorem 1 here based on the above assumptions:

Theorem (REAR, full). *Under the standing assumptions above, form the online statistics \mathbf{G}, \mathbf{Q} as in Eq. (2) and the ridge solution $\widehat{\mathbf{U}}$ as in Eq. (4). Let N be the total number of samples used to form \mathbf{G}, \mathbf{Q} and let U^* denote the population regularized minimizer in the feature space induced by the kernel $k(\mathbf{h}, \mathbf{h}') = \mathbb{E}_{\mathbf{r}}[\sigma(\mathbf{h}^\top \mathbf{r}) \sigma(\mathbf{h}'^\top \mathbf{r})]$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ (over \mathbf{R} and the training samples), the excess (population) squared risk decomposes as*

$$\mathcal{R}(\widehat{\mathbf{U}}) - \mathcal{R}(U^*) \leq \mathcal{E}_{\text{feat}}(M, \delta) + \mathcal{E}_{\text{estim}}(N, \lambda, \delta) + \mathcal{E}_{\text{reg}}(\lambda),$$

where, for universal constants C_i (depending on H, L_σ, C),

$$\mathcal{E}_{\text{feat}}(M, \delta) \leq C_1 \sqrt{\frac{\log(N/\delta)}{M}}, \quad \mathcal{E}_{\text{estim}}(N, \lambda, \delta) \leq C_2 \frac{1}{\sqrt{N}} \cdot \frac{1}{\lambda}, \quad \mathcal{E}_{\text{reg}}(\lambda) \leq C_3 \lambda \|U^*\|_F^2.$$

²This assumption is not needed for the excess-risk decomposition in Theorem 1 itself, but acts as a bridge.

E.1.1 LEMMA: RANDOM-FEATURE CONCENTRATION

Lemma 1. Let $S = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ be the finite training set and write $\mathbf{h}_i = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$. Define the kernel

$$k(\mathbf{h}, \mathbf{h}') := \mathbb{E}_{\mathbf{r} \sim \mathcal{N}(0, \mathbf{I}_d)} [\sigma(\mathbf{h}^\top \mathbf{r}) \sigma(\mathbf{h}'^\top \mathbf{r})].$$

Then, under the standing assumptions, for any $\delta \in (0, 1)$ and any $\varepsilon \in (0, 1)$, if

$$M \geq C_{\text{rf}} \cdot \frac{\log(N^2/\delta)}{\varepsilon^2}$$

(with C_{rf} depending only on H, L_σ, C above), then with probability at least $1 - \delta$ over \mathbf{R} ,

$$\max_{1 \leq i, j \leq N} \left| \frac{1}{M} \boldsymbol{\varphi}(\mathbf{x}_i)^\top \boldsymbol{\varphi}(\mathbf{x}_j) - k(\mathbf{h}_i, \mathbf{h}_j) \right| \leq \varepsilon.$$

Proof. Fix a pair (i, j) . Write

$$Z_\ell := \varphi_\ell(\mathbf{x}_i) \varphi_\ell(\mathbf{x}_j) - \mathbb{E}[\varphi_\ell(\mathbf{x}_i) \varphi_\ell(\mathbf{x}_j)], \quad \ell = 1, \dots, M.$$

The Z_ℓ are independent (across ℓ) mean-zero random variables because columns \mathbf{r}_ℓ are independent. We will apply Bernstein's inequality for sums of independent sub-exponential variables; to do so, we need a variance proxy and a uniform tail bound.

From the growth assumption $|\sigma(z)| \leq C(1 + |z|)$ and $\mathbf{r}_\ell \sim \mathcal{N}(0, \mathbf{I}_d)$, the marginal $\varphi_\ell(\mathbf{x}) = \sigma(\mathbf{h}_i^\top \mathbf{r}_\ell)$ is sub-Gaussian / sub-exponential: more precisely, since $\mathbf{h}_i^\top \mathbf{r}_\ell \sim \mathcal{N}(0, \|\mathbf{h}_i\|^2) \leq \mathcal{N}(0, H^2)$, we have for some constants v, b (depending on H, L_σ, C) that $\mathbb{P}(|\varphi_\ell(\mathbf{x})| \geq t) \leq 2 \exp(-ct)$ for large t ; thus $\varphi_\ell(\mathbf{x}) \varphi_\ell(\mathbf{x}')$ is sub-exponential with parameters bounded by functions of H, L_σ, C . Concretely, one can verify

$$\|Z_\ell\|_{\psi_1} \leq \tilde{b}$$

for a finite constant \tilde{b} depending only on H, L_σ, C , where $\|\cdot\|_{\psi_1}$ denotes the standard sub-exponential Orlicz norm. Hence, applying Bernstein's inequality for sub-exponential variables yields, for any $\tau > 0$,

$$\mathbb{P}\left(\left|\sum_{\ell=1}^M Z_\ell\right| \geq \tau\right) \leq 2 \exp\left(-c \min\left(\frac{\tau^2}{M\tilde{v}^2}, \frac{\tau}{\tilde{b}}\right)\right),$$

with constants $\tilde{v}, \tilde{b}, c > 0$ determined by the sub-exponential parameters.

Choose $\tau = M\varepsilon$. Plugging $\tau = M\varepsilon$ and requiring the RHS to be $\leq \delta/N^2$ (to union bound over all $\leq N^2$ pairs) yields the condition

$$M \geq C_{\text{rf}} \frac{\log(N^2/\delta)}{\varepsilon^2}$$

for some C_{rf} (combining cases of Bernstein). This gives, for fixed pair (i, j) ,

$$\mathbb{P}\left(\left|\frac{1}{M} \sum_{\ell=1}^M Z_\ell\right| \geq \varepsilon\right) \leq \frac{\delta}{N^2}.$$

Apply union bound over all $\leq N^2$ ordered pairs (i, j) . This yields the claimed uniform bound with probability at least $1 - \delta$. \square

Remarks on applicability of Bernstein. We used Bernstein for independent sub-exponential summands. The summands are independent across random-feature index ℓ ; sub-exponentiality follows from (i) Gaussianity of \mathbf{r}_ℓ and (ii) Lipschitz + linear-growth of σ . For ReLU (which is Lipschitz with linear growth) the same argument applies (moments of Gaussian tails control tails of $\sigma(\cdot)$).

E.1.2 LEMMA: RIDGE PERTURBATION

We next show that when the empirical feature covariance concentrates around its population counterpart and the empirical cross-covariance concentrates, then the finite-sample ridge solution is close to the population ridge solution.

Lemma 2. *Let $\Phi \in \mathbb{R}^{N \times M}$ be the feature matrix with rows $\varphi(\mathbf{x}_i)^\top$, define empirical covariance*

$$\widehat{\Sigma} = \frac{1}{N} \Phi^\top \Phi \in \mathbb{R}^{M \times M}, \quad \widehat{b} = \frac{1}{N} \Phi^\top Y \in \mathbb{R}^{M \times T},$$

where $Y \in \mathbb{R}^{N \times T}$ is the one-hot label matrix (or soft labels). Let the population quantities be $\Sigma = \mathbb{E}[\varphi(\mathbf{x})\varphi(\mathbf{x})^\top]$ and $b = \mathbb{E}[\varphi(\mathbf{x})\mathbf{c}^\top]$. Denote population ridge solution

$$U_\lambda^* := (\Sigma + \lambda I)^{-1}b, \quad \widehat{U}_\lambda := (\widehat{\Sigma} + \lambda I)^{-1}\widehat{b}.$$

If $\|\widehat{\Sigma} - \Sigma\|_{\text{op}} \leq \frac{\lambda}{2}$ and $\|\widehat{b} - b\|_F \leq \epsilon_b$, then

$$\|\widehat{U}_\lambda - U_\lambda^*\|_F \leq \frac{2}{\lambda} \epsilon_b + \frac{2}{\lambda^2} \|b\|_F \|\widehat{\Sigma} - \Sigma\|_{\text{op}}.$$

Proof. We write

$$\widehat{U}_\lambda - U_\lambda^* = (\widehat{\Sigma} + \lambda I)^{-1}(\widehat{b} - b) + [(\widehat{\Sigma} + \lambda I)^{-1} - (\Sigma + \lambda I)^{-1}]b.$$

For the first term use operator norm bound $\|(\widehat{\Sigma} + \lambda I)^{-1}\|_{\text{op}} \leq 1/\lambda$ to get

$$\|(\widehat{\Sigma} + \lambda I)^{-1}(\widehat{b} - b)\|_F \leq \frac{1}{\lambda} \|\widehat{b} - b\|_F.$$

For the second term use the identity $A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$ with $A = \widehat{\Sigma} + \lambda I$, $B = \Sigma + \lambda I$. Hence

$$\|A^{-1} - B^{-1}\|_{\text{op}} \leq \|A^{-1}\|_{\text{op}} \|A - B\|_{\text{op}} \|B^{-1}\|_{\text{op}} \leq \frac{1}{\lambda} \|\widehat{\Sigma} - \Sigma\|_{\text{op}} \frac{1}{\lambda}.$$

Thus

$$\|[(\widehat{\Sigma} + \lambda I)^{-1} - (\Sigma + \lambda I)^{-1}]b\|_F \leq \frac{1}{\lambda^2} \|\widehat{\Sigma} - \Sigma\|_{\text{op}} \|b\|_F.$$

Combining the two terms and tightening constants when $\|\widehat{\Sigma} - \Sigma\|_{\text{op}} \leq \lambda/2$ gives the stated bound. \square

Concentration of $\widehat{\Sigma}$ and \widehat{b} . We next control the empirical covariance and cross-covariance. Recall

$$\widehat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_i)\varphi(\mathbf{x}_i)^\top, \quad \Sigma = \mathbb{E}[\varphi(\mathbf{x})\varphi(\mathbf{x})^\top].$$

It is convenient to write

$$\widehat{\Sigma} - \Sigma = \sum_{i=1}^N X_i, \quad X_i := \frac{1}{N} (\varphi(\mathbf{x}_i)\varphi(\mathbf{x}_i)^\top - \Sigma).$$

Each X_i is self-adjoint and satisfies $\mathbb{E}[X_i] = 0$. Under the bounded-embedding and activation assumptions (cf. Lemma 1), there exists a constant $C_\varphi > 0$ (depending only on (H, L_σ, C)) such that $\|\varphi(\mathbf{x})\|_2 \leq C_\varphi$ almost surely. Hence, for every i ,

$$\|X_i\|_{\text{op}} \leq \frac{1}{N} (\|\varphi(\mathbf{x}_i)\varphi(\mathbf{x}_i)^\top\|_{\text{op}} + \|\Sigma\|_{\text{op}}) \leq \frac{1}{N} (C_\varphi^2 + \|\Sigma\|_{\text{op}}) =: \frac{L_0}{N}.$$

Similarly, we can bound the “matrix variance” term

$$v^2 := \left\| \sum_{i=1}^N \mathbb{E}[X_i^2] \right\|_{\text{op}} = N \left\| \mathbb{E}[X_1^2] \right\|_{\text{op}} \leq \frac{N}{N^2} \left\| \mathbb{E}[(\varphi(\mathbf{x})\varphi(\mathbf{x})^\top - \Sigma)^2] \right\|_{\text{op}} \leq \frac{V_0}{N},$$

for some constant $V_0 > 0$ depending only on (H, L_σ, C) . For completeness, we recall a standard matrix Bernstein inequality (Theorem 6.1 in Tropp (2012)): if $\{X_i\}_{i=1}^N$ are independent, mean-zero, self-adjoint matrices with $\|X_i\|_{\text{op}} \leq L$ almost surely and

$$v^2 := \left\| \sum_{i=1}^N \mathbb{E}[X_i^2] \right\|_{\text{op}},$$

then for all $t > 0$:

$$\mathbb{P}\left(\left\| \sum_{i=1}^N X_i \right\|_{\text{op}} \geq t\right) \leq 2D \exp\left(-\frac{t^2/2}{v^2 + Lt/3}\right),$$

where D denotes the matrix dimension (in our case $D = M$, the feature dimension). Applying this with $t = \varepsilon$ and our bounds $L \leq L_0/N$ and $v^2 \leq V_0/N$ gives

$$\mathbb{P}(\|\hat{\Sigma} - \Sigma\|_{\text{op}} \geq \varepsilon) = \mathbb{P}\left(\left\| \sum_{i=1}^N X_i \right\|_{\text{op}} \geq \varepsilon\right) \leq 2M \exp\left(-\frac{N\varepsilon^2/2}{V_0 + L_0\varepsilon/3}\right).$$

For $\varepsilon \in (0, 1)$ the denominator in the exponent is bounded above by a constant multiple of V_0 , so there exists $C' > 0$ (depending only on (H, L_σ, C)) such that, for all $\delta \in (0, 1)$, taking

$$\varepsilon = C' \sqrt{\frac{\log(M/\delta)}{N}}$$

ensures

$$\mathbb{P}(\|\hat{\Sigma} - \Sigma\|_{\text{op}} \geq \varepsilon) \leq \delta.$$

Equivalently, with probability at least $1 - \delta$,

$$\|\hat{\Sigma} - \Sigma\|_{\text{op}} \leq C' \sqrt{\frac{\log(M/\delta)}{N}}.$$

For the empirical cross-covariance $\hat{b} = \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_i) \mathbf{c}_i^\top$ and its population counterpart $b = \mathbb{E}[\varphi(\mathbf{x}) \mathbf{c}^\top]$ we apply the same argument column-wise (each column is an average of bounded sub-exponential vectors of length M) and obtain

$$\|\hat{b} - b\|_F \leq C'' \sqrt{\frac{\log(T/\delta)}{N}},$$

for some constant $C'' > 0$ depending only on the same problem parameters. Adjusting constants to account for the two events and taking a union bound, we may assume that both inequalities hold simultaneously with probability at least $1 - \delta$. Choosing N large enough to make $\|\hat{\Sigma} - \Sigma\|_{\text{op}} \leq \lambda/2$ and to make $\|\hat{b} - b\|_F$ (denoted ϵ_b in Lemma 2) small then yields the desired estimation error term in Lemma 2.

E.1.3 LEMMA: ONLINE STATISTICS IMPLEMENT BATCH RIDGE

Lemma 3. *If \mathbf{G} and \mathbf{Q} are formed by accumulating per-example contributions*

$$\mathbf{G} = \sum_{i=1}^N \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^\top, \quad \mathbf{Q} = \sum_{i=1}^N \varphi(\mathbf{x}_i) \mathbf{c}_i^\top,$$

then the closed-form solution $\hat{\mathbf{U}}^\top = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{Q}$ equals the ridge regression solution computed in batch on features $\varphi(\mathbf{x}_i)$ and labels \mathbf{c}_i . Moreover, if the online implementation maintains $(\mathbf{G} + \lambda \mathbf{I})^{-1}$ via rank-1 updates, numerical equivalence holds up to floating-point precision.

Proof. This is algebraic: batch ridge with design matrix Φ and labels Y solves $\hat{\mathbf{U}}^\top = (\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top Y$. But $\Phi^\top \Phi = \sum_i \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^\top = \mathbf{G}$ and $\Phi^\top Y = \mathbf{Q}$. The equality follows. For incremental numerical maintenance of the inverse, standard Sherman–Morrison or Woodbury updates apply; also numerically stable Cholesky-updates are recommended when M is large. \square

E.1.4 COMBINING LEMMAS: PROOF OF THEOREM 1

Proof. The excess population risk decomposes as

$$\mathcal{R}(\widehat{U}) - \mathcal{R}(U^*) = \underbrace{\mathcal{R}(\widehat{U}) - \mathcal{R}(U_\lambda^*)}_{\text{estimation error}} + \underbrace{\mathcal{R}(U_\lambda^*) - \mathcal{R}(U^*)}_{\text{reg. bias}}.$$

The regularization bias is the usual ridge bias and yields the $\mathcal{E}_{\text{reg}}(\lambda)$ term; standard calculus shows it is bounded by $\lambda \|U^*\|_F^2$ up to constant factors.

For the estimation error, apply Lemma 2 to relate the empirical ridge solution (which equals the on-line \widehat{U} by Lemma 3) to the population ridge solution U_λ^* . The two perturbation terms are controlled by $\|\widehat{\Sigma} - \Sigma\|_{\text{op}}$ and $\|\widehat{b} - b\|_F$, which in turn are bounded by the matrix Bernstein concentration bounds for $\widehat{\Sigma}$ and \widehat{b} derived above. This yields the stated $O(\frac{1}{\sqrt{N}} \cdot \frac{1}{\lambda})$ behavior for $\mathcal{E}_{\text{estim}}$ (explicit constants follow from the above bounds).

Finally, the approximation error due to random features is precisely Lemma 1: replacing the kernel k by its Monte Carlo approximation using M independent features introduces a uniform $O(\sqrt{\log(N/\delta)/M})$ perturbation in inner products, which propagates to the excess risk as the term $\mathcal{E}_{\text{feat}}(M, \delta)$ displayed in Theorem 1. \square

Margin-based routing accuracy. Under the additional margin assumption in the REAR standing assumptions and a uniform bound $\|\varphi(\mathbf{x})\|_2 \leq C_\varphi$ (for some constant C_φ depending only on (H, L_σ, C)), the excess risk bound above can be converted into a bound on misrouting probability. Let $\hat{t}(\mathbf{x}) := \arg \max_t s_{\widehat{U}}(\mathbf{x})_t$ and $t^*(\mathbf{x}) := \arg \max_t s_{U^*}(\mathbf{x})_t$ denote the experts selected by \widehat{U} and U^* , respectively. Then

$$\mathbb{P}(\hat{t}(X) \neq t^*(X)) \leq \frac{8C_\varphi^2}{\lambda\gamma^2} (\mathcal{R}(\widehat{U}) - \mathcal{R}(U^*)).$$

Proof. On the event $\{\hat{t}(\mathbf{x}) \neq t^*(\mathbf{x})\}$, the margin condition implies

$$\gamma \leq s_{U^*}(\mathbf{x})_{t^*(\mathbf{x})} - s_{U^*}(\mathbf{x})_{\hat{t}(\mathbf{x})} \leq 2 \max_t |s_{U^*}(\mathbf{x})_t - s_{\widehat{U}}(\mathbf{x})_t|.$$

$$\text{Hence, } \mathbf{1}\{\hat{t}(\mathbf{x}) \neq t^*(\mathbf{x})\} \leq (2/\gamma)^2 \|s_{\widehat{U}}(\mathbf{x}) - s_{U^*}(\mathbf{x})\|_2^2.$$

Given $s_U(\mathbf{x}) = \varphi(\mathbf{x})U^\top$ and Cauchy–Schwarz yields:

$$\begin{aligned} \max_t |s_{U^*}(\mathbf{x})_t - s_{\widehat{U}}(\mathbf{x})_t| &\leq \|\varphi(\mathbf{x})\|_2 \|\widehat{U} - U^*\|_F \leq C_\varphi \|\widehat{U} - U^*\|_F, \\ \Rightarrow \mathbf{1}\{\hat{t}(\mathbf{x}) \neq t^*(\mathbf{x})\} &\leq (2C_\varphi/\gamma)^2 \|\widehat{U} - U^*\|_F^2. \end{aligned}$$

Since the regularized risk is defined as $\mathcal{R}(U) := \mathbb{E}\|s_U(X) - C\|_2^2 + \lambda\|U\|_F^2$, the quadratic ridge term $\lambda\|U\|_F^2$ makes \mathcal{R} (at least) λ -strongly convex in U . In particular, strong convexity implies:

$$\begin{aligned} \mathcal{R}(\widehat{U}) - \mathcal{R}(U^*) &\geq (\lambda/2) \|\widehat{U} - U^*\|_F^2, \\ \Rightarrow \|\widehat{U} - U^*\|_F^2 &\leq (2/\lambda) (\mathcal{R}(\widehat{U}) - \mathcal{R}(U^*)). \end{aligned}$$

Combining these inequalities and taking expectations over X gives the stated bound. \square

E.2 PROOFS FOR TE² (THEOREM 2)

E.2.1 NOTATION AND ASSUMPTIONS

We reuse notation from the main text. For a fixed expert (prompt/head), we denote:

- $\mathbf{W}_t^* \in \mathbb{R}^{|\mathcal{Y}| \times D}$: the (population) time- t optimal linear parameter (in the chosen feature space) for that expert.

- \mathbf{W}_t : the instantaneous (online) estimator after observing the t -th update; we model $\mathbf{W}_t = \mathbf{W}_t^* + \boldsymbol{\xi}_t$.
- EMA head with decay $\alpha \in (0, 1)$:

$$\widetilde{\mathbf{W}}_t^{(\alpha)} = (1 - \alpha) \sum_{k \geq 0} \alpha^k \mathbf{W}_{t-k}, \quad L(\alpha) = \frac{1}{1-\alpha}.$$

- Discounted path length (drift measure):

$$P_t := \sum_{j \geq 1} \gamma^{j-1} \|\mathbf{W}_{t-j+1}^* - \mathbf{W}_{t-j}^*\|,$$

where we define the discount factor to match the EMA decay, i.e., $\gamma = \alpha$ or comparable. This is a discounted analogue of the standard path length / variation measure used in dynamic regret (Zinkevich, 2003; Besbes et al., 2015).

- Standing conditions: finite P_t , zero-mean noise with bounded variance:

$$\mathbb{E}[\boldsymbol{\xi}_t] = \mathbf{0}, \quad \mathbb{E}\|\boldsymbol{\xi}_t\|^2 \leq \zeta^2.$$

- (Optional, for classification calibration) A margin $\Delta > 0$ and a Lipschitz map-to-logits with constant C_f imply a bound on 0/1 error from parameter MSE.

We then state the complete Theorem 2 here based on the above assumptions:

Theorem (TE², full). *Under the standing conditions above, fix t and an EMA decay α with $L = 1/(1 - \alpha)$. There exist constants $C_1, C_2 > 0$ such that*

$$\mathbb{E}\|\widetilde{\mathbf{W}}_t^{(\alpha)} - \mathbf{W}_t^*\|^2 \leq C_1 \frac{\zeta^2}{L} + C_2 (L P_t)^2.$$

Moreover, if we keep a geometric grid of windows $\{L_i\}_{i=1}^m$ with ratio $r > 1$ (e.g., $L_i = r^{i-1}$), then for every t there exists an index i_t with

$$\mathbb{E}\|\widetilde{\mathbf{W}}_t^{(\alpha_{i_t})} - \mathbf{W}_t^*\|^2 \leq c(r) \min_{L \geq 1} \left\{ C_1 \frac{\zeta^2}{L} + C_2 (L P_t)^2 \right\},$$

where $c(r)$ depends only on the grid ratio r (one can take, for example, $c(r) = \max(r^2, r)$ by the argument below). Additionally, if a margin $\Delta > 0$ holds and the logits map is Lipschitz with constant C_f , then the above parameter MSE implies a classification error bound $O((C_f \varepsilon / \Delta)^2)$ whenever $\mathbb{E}\|\widetilde{\mathbf{W}}_t^{(\alpha)} - \mathbf{W}_t^*\|^2 \leq \varepsilon^2$.

E.2.2 PROOF OF THEOREM 2

Proof. The proof proceeds via a variance–bias decomposition and a geometric grid selection argument, followed by a calibration from parameter MSE to classification error.

We start by decomposing the error into the variance term. Define the difference between the EMA and the population optimum

$$\widetilde{\mathbf{W}}_t^{(\alpha)} - \overline{\mathbf{W}}_t^* = (1 - \alpha) \sum_{k \geq 0} \alpha^k \boldsymbol{\xi}_{t-k},$$

$$\overline{\mathbf{W}}_t^* := (1 - \alpha) \sum_{k \geq 0} \alpha^k \mathbf{W}_{t-k}^*,$$

where $\overline{\mathbf{W}}_t^*$ is the EMA of the population optima.

Since we have $\mathbb{E}[\boldsymbol{\xi}_t] = \mathbf{0}$ and $\mathbb{E}\|\boldsymbol{\xi}_t\|^2 \leq \zeta^2$, and the EMA weights are $a_k = (1 - \alpha)\alpha^k$, we can compute their squared sum explicitly:

$$\sum_{k \geq 0} a_k^2 = (1 - \alpha)^2 \sum_{k \geq 0} \alpha^{2k} = \frac{(1 - \alpha)^2}{1 - \alpha^2} = \frac{1 - \alpha}{1 + \alpha} \leq \frac{1}{L},$$

Therefore,

$$\mathbb{E} \|\widetilde{\mathbf{W}}_t^{(\alpha)} - \overline{\mathbf{W}}_t^*\|^2 = \mathbb{E} \left\| \sum_{k \geq 0} a_k \boldsymbol{\xi}_{t-k} \right\|^2 \leq \zeta^2 \sum_{k \geq 0} a_k^2 \leq \frac{\zeta^2}{L},$$

which matches the variance term $C_1 \zeta^2/L$ in Theorem 2 (with C_1 absorbing the constant factor). Next, we address the bias term. The difference between the EMA of the population optima and the actual population optimum is given by:

$$\overline{\mathbf{W}}_t^* - \mathbf{W}_t^* = \sum_{k \geq 0} a_k (\mathbf{W}_{t-k}^* - \mathbf{W}_t^*).$$

Reordering sums yields:

$$\begin{aligned} \|\overline{\mathbf{W}}_t^* - \mathbf{W}_t^*\| &\leq \sum_{j \geq 1} \left(\sum_{k \geq j} a_k \right) \|\mathbf{W}_{t-j+1}^* - \mathbf{W}_{t-j}^*\| = \sum_{j \geq 1} \alpha^j \Delta_{t-j+1}, \\ \text{where } \Delta_u &= \|\mathbf{W}_u^* - \mathbf{W}_{u-1}^*\|. \end{aligned} \quad (11)$$

Using the discounted path length P_t with the same discount factor $\gamma = \alpha$ as in the EMA definition, we can simply rewrite the above bound as

$$\|\overline{\mathbf{W}}_t^* - \mathbf{W}_t^*\| = \sum_{j \geq 1} \alpha^j \Delta_{t-j+1} = \alpha \sum_{j \geq 1} \alpha^{j-1} \Delta_{t-j+1} = \alpha P_t \leq L P_t,$$

where we used $L = 1/(1 - \alpha) \geq \alpha$. Consequently, the squared bias admits the explicit bound

$$\|\overline{\mathbf{W}}_t^* - \mathbf{W}_t^*\|^2 \leq (L P_t)^2,$$

which corresponds to the term $C_2 (L P_t)^2$ in Theorem 2.

Then, by combining the variance and bias terms, we apply the inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$, which yields the claimed MSE bound:

$$\mathbb{E} \left[\|\widetilde{\mathbf{W}}_t^{(\alpha)} - \overline{\mathbf{W}}_t^*\|^2 \right] \leq 2 \frac{\zeta^2}{L} + 2 (L P_t)^2.$$

Therefore, in Theorem 2 we may take the explicit choice $C_1 = C_2 = 2$. If we allow γ to differ slightly from α , the constant C_2 becomes $(\alpha/\gamma)^2$ (bounded if we restrict $\gamma \in [(1 - \epsilon)\alpha, (1 + \epsilon)\alpha]$).

For the geometric bank, it is convenient to write the bound in the generic form

$$f(L) := \frac{A}{L} + B(L P_t)^2, \quad A \asymp \zeta^2, \quad B \asymp 1.$$

A direct derivative calculation of $f'(L)$ shows that the minimizer over $L > 0$ is

$$f'(L) = -\frac{A}{L^2} + 2B P_t^2 L \Rightarrow L^* = \left(\frac{A}{2B P_t^2} \right)^{1/3}.$$

If we maintain a geometric grid $L_i = r^{i-1}$ with ratio $r > 1$, then for any L^* there exists an index i_t such that $L_{i_t} \in [L^*/r, r L^*]$. Writing $L = L^* \eta$ with $\eta \in [1/r, r]$ and using the optimality condition $A/L^* = 2B(L^* P_t)^2$, we obtain

$$\frac{f(L)}{f(L^*)} = \frac{2B(L^* P_t)^2/\eta + \eta^2 B(L^* P_t)^2}{3B(L^* P_t)^2} = \frac{2/\eta + \eta^2}{3}.$$

The right-hand side is maximized over $\eta \in [1/r, r]$ at one of the endpoints; a simple bound yields

$$\sup_{\eta \in [1/r, r]} \frac{2/\eta + \eta^2}{3} \leq \max \left(\frac{2}{3r} + \frac{r^2}{3}, \frac{2r}{3} + \frac{1}{3r^2} \right) \leq \max(r^2, r).$$

Therefore $f(L_{i_t}) \leq c(r) f(L^*)$ with the explicit choice $c(r) = \max(r^2, r)$ used in Theorem 2 (obviously $c(r) = r^2$, given $r > 1$ in our case.)

From the general classification error's view, by the Lipschitz-to-logit condition, the induced logit error is at most:

$$C_f \|\widetilde{\mathbf{W}} - \mathbf{W}^*\|.$$

Under the margin condition ($\Delta > 0$), the standard margin-to-0/1 calibration gives an error bound $O((C_f \varepsilon / \Delta)^2)$ when $\mathbb{E} \|\widetilde{\mathbf{W}} - \mathbf{W}^*\|^2 \leq \varepsilon^2$. \square

F ADDITIONAL EXPERIMENT RESULTS

F.1 EXPERIMENT WITH DISJOINT AND BLURRY SETUP VARIANTS

Table 9: Performance comparison of r_D variants in FlyPrompt. We fixed $r_B = 10\%$ and use Sup-21K backbone PTM. All results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds.

$r_D(\%)$	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$	$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$	$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$
0 (pure blurry)	L2P	73.64 \pm 1.69	81.71 \pm 1.28	42.97 \pm 1.91	42.53 \pm 0.79	64.62 \pm 3.37	62.72 \pm 3.32
	DualPrompt	72.16 \pm 2.40	77.52 \pm 1.43	45.08 \pm 3.48	42.34 \pm 1.55	65.42 \pm 3.53	62.84 \pm 2.21
	CODA-P	72.33 \pm 5.84	78.56 \pm 4.46	50.31 \pm 4.42	48.88 \pm 3.46	66.64 \pm 3.42	63.70 \pm 2.48
	MVP	67.10 \pm 2.33	74.94 \pm 7.45	39.04 \pm 1.49	32.28 \pm 3.12	56.86 \pm 3.71	55.50 \pm 6.71
	MISA	78.38 \pm 1.49	83.04 \pm 1.44	51.78 \pm 3.22	47.64 \pm 0.64	67.38 \pm 3.25	64.49 \pm 2.55
	FlyPrompt (ours)	80.12 \pm 1.38	87.11 \pm 0.52	55.44 \pm 1.82	55.89 \pm 0.86	71.60 \pm 3.49	74.72 \pm 1.69
50 (mixed)	L2P	76.23 \pm 2.73	79.11 \pm 1.43	44.40 \pm 1.03	42.03 \pm 1.72	64.30 \pm 2.18	61.42 \pm 2.13
	DualPrompt	76.04 \pm 3.32	76.62 \pm 0.74	46.13 \pm 1.94	40.80 \pm 1.04	65.03 \pm 2.24	62.43 \pm 1.78
	CODA-P	79.13 \pm 3.06	80.91 \pm 0.70	51.87 \pm 2.81	48.09 \pm 2.75	66.01 \pm 2.20	62.90 \pm 2.46
	MVP	67.74 \pm 4.96	63.22 \pm 0.69	39.50 \pm 1.41	32.63 \pm 3.95	54.69 \pm 3.14	50.07 \pm 3.86
	MISA	80.35 \pm 2.39	80.75 \pm 1.24	51.52 \pm 2.09	45.08 \pm 1.43	65.40 \pm 3.01	60.20 \pm 1.82
	FlyPrompt (ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88
100 (disjoint)	L2P	82.98 \pm 0.72	78.79 \pm 0.95	45.48 \pm 0.71	43.12 \pm 0.75	71.74 \pm 3.58	61.52 \pm 1.82
	DualPrompt	81.12 \pm 2.10	75.94 \pm 0.37	46.79 \pm 2.00	41.42 \pm 0.63	72.81 \pm 3.80	62.46 \pm 0.97
	CODA-P	82.68 \pm 3.99	77.97 \pm 3.19	53.01 \pm 3.01	49.31 \pm 2.87	73.95 \pm 4.10	63.90 \pm 0.94
	MVP	74.92 \pm 1.10	56.17 \pm 2.98	40.43 \pm 0.51	28.32 \pm 5.55	63.41 \pm 2.52	43.80 \pm 2.60
	MISA	85.67 \pm 1.01	81.04 \pm 1.02	53.88 \pm 1.99	47.63 \pm 0.78	74.27 \pm 3.56	62.97 \pm 1.44
	FlyPrompt (ours)	88.25 \pm 0.90	85.51 \pm 0.64	57.41 \pm 0.95	55.69 \pm 0.33	78.14 \pm 3.62	74.34 \pm 0.74

Table 10: Performance comparison of r_B variants in FlyPrompt. We fixed $r_D = 50\%$ and use Sup-21K PTM backbone. All results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds.

$r_B(\%)$	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$	$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$	$A_{auc}(\%, \uparrow)$	$A_{last}(\%, \uparrow)$
10	L2P	76.23 \pm 2.73	79.11 \pm 1.43	44.40 \pm 1.03	42.03 \pm 1.72	64.30 \pm 2.18	61.42 \pm 2.13
	DualPrompt	76.04 \pm 3.32	76.62 \pm 0.74	46.13 \pm 1.94	40.80 \pm 1.04	65.03 \pm 2.24	62.43 \pm 1.78
	CODA-P	79.13 \pm 3.06	80.91 \pm 0.70	51.87 \pm 2.81	48.09 \pm 2.75	66.01 \pm 2.20	62.90 \pm 2.46
	MVP	67.74 \pm 4.96	63.22 \pm 0.69	39.50 \pm 1.41	32.63 \pm 3.95	54.69 \pm 3.14	50.07 \pm 3.86
	MISA	80.35 \pm 2.39	80.75 \pm 1.24	51.52 \pm 2.09	45.08 \pm 1.43	65.40 \pm 3.01	60.20 \pm 1.82
	FlyPrompt (ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88
30	L2P	78.48 \pm 0.92	80.13 \pm 0.87	43.32 \pm 0.79	42.27 \pm 1.40	63.67 \pm 2.03	63.60 \pm 3.09
	DualPrompt	77.76 \pm 1.65	77.50 \pm 0.49	45.11 \pm 1.09	41.01 \pm 0.70	64.36 \pm 1.98	63.72 \pm 2.22
	CODA-P	81.50 \pm 1.20	82.65 \pm 0.72	50.55 \pm 2.24	47.58 \pm 2.81	65.89 \pm 1.42	64.97 \pm 2.83
	MVP	71.01 \pm 1.70	65.71 \pm 4.60	38.62 \pm 1.35	32.43 \pm 4.75	54.16 \pm 4.56	51.04 \pm 6.89
	MISA	82.54 \pm 1.08	82.50 \pm 0.68	51.69 \pm 0.94	47.09 \pm 1.16	67.13 \pm 1.72	66.53 \pm 2.39
	FlyPrompt (ours)	84.61 \pm 1.25	86.89 \pm 0.38	55.30 \pm 0.86	55.43 \pm 1.04	70.19 \pm 2.01	74.25 \pm 1.27
50	L2P	77.44 \pm 2.42	80.31 \pm 0.69	44.39 \pm 1.72	43.66 \pm 1.04	65.42 \pm 2.71	64.62 \pm 1.63
	DualPrompt	77.44 \pm 2.64	77.13 \pm 1.08	46.23 \pm 1.83	41.99 \pm 0.72	66.40 \pm 2.72	65.34 \pm 3.05
	CODA-P	81.39 \pm 2.18	83.10 \pm 0.97	53.05 \pm 1.60	50.20 \pm 1.74	67.88 \pm 2.26	65.44 \pm 2.32
	MVP	67.97 \pm 4.78	58.11 \pm 1.26	40.68 \pm 1.59	31.87 \pm 6.56	57.25 \pm 3.76	53.86 \pm 3.41
	MISA	81.81 \pm 2.29	82.51 \pm 0.38	53.27 \pm 1.71	48.32 \pm 0.84	68.68 \pm 2.47	66.84 \pm 2.27
	FlyPrompt (ours)	83.69 \pm 1.81	86.31 \pm 0.73	56.50 \pm 1.87	55.59 \pm 0.83	71.95 \pm 1.92	74.03 \pm 1.02

As detailed in Tabs. 9 and 10, FlyPrompt consistently outperforms baselines from the purely blurry regime ($r_D = 0$) to fully disjoint tasks ($r_D = 1$, namely, online CIL) and across different blurry ratios r_B , demonstrating robustness under extreme GCL configurations and superior versatility.

F.2 GCL EVALUATION METRICS

Table 11: Average accuracy of 5 sessions and forgetting of different GCL methods over three datasets. All results are reported as an average of five runs (\pm standard deviation) with different random seeds.

PTM	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{avg}(\%, \uparrow)$	$F_{last}(\%, \downarrow)$	$A_{avg}(\%, \uparrow)$	$F_{last}(\%, \downarrow)$	$A_{avg}(\%, \uparrow)$	$F_{last}(\%, \downarrow)$
Sup-21K	L2P	75.41 \pm 2.75	11.53 \pm 1.44	47.82 \pm 0.95	19.16 \pm 1.92	65.44 \pm 3.30	29.15 \pm 3.19
	DualPrompt	75.96 \pm 2.56	11.42 \pm 0.91	49.90 \pm 2.34	18.24 \pm 4.34	66.54 \pm 3.05	27.18 \pm 2.91
	CODA-P	78.73 \pm 3.09	9.95 \pm 1.30	55.75 \pm 2.83	18.49 \pm 2.44	67.19 \pm 2.98	27.43 \pm 3.42
	MVP	64.70 \pm 4.14	33.19 \pm 2.33	40.14 \pm 1.39	43.39 \pm 4.23	52.42 \pm 5.06	47.61 \pm 5.73
	MISA	79.67 \pm 1.78	9.67 \pm 1.39	55.78 \pm 1.41	21.46 \pm 4.25	66.92 \pm 3.12	30.06 \pm 3.70
	FlyPrompt (ours)	82.72 \pm 2.69	5.03 \pm 1.16	60.38 \pm 1.76	13.42 \pm 1.74	72.22 \pm 4.35	10.97 \pm 1.52
Sup-21K/1K	L2P	60.98 \pm 10.04	14.88 \pm 6.27	50.21 \pm 3.09	35.88 \pm 4.51	43.72 \pm 5.19	35.05 \pm 9.31
	DualPrompt	66.13 \pm 4.34	19.18 \pm 5.13	56.53 \pm 1.22	30.73 \pm 7.45	47.37 \pm 5.07	35.40 \pm 8.16
	CODA-P	66.81 \pm 4.87	19.82 \pm 6.85	55.01 \pm 1.66	35.58 \pm 5.52	45.38 \pm 4.99	35.73 \pm 8.95
	MVP	63.22 \pm 1.67	46.98 \pm 8.15	50.91 \pm 3.09	51.11 \pm 2.86	46.13 \pm 1.76	62.93 \pm 8.05
	MISA	60.11 \pm 8.35	11.38 \pm 2.76	54.17 \pm 2.92	28.66 \pm 8.19	43.33 \pm 5.30	33.95 \pm 8.63
	FlyPrompt (ours)	76.86 \pm 2.29	9.50 \pm 2.33	64.41 \pm 1.80	21.59 \pm 4.28	55.12 \pm 5.49	21.71 \pm 3.56
iBOT-21K	L2P	52.73 \pm 7.78	12.58 \pm 6.24	38.41 \pm 6.38	34.17 \pm 7.54	14.93 \pm 2.30	19.09 \pm 8.05
	DualPrompt	62.86 \pm 7.45	19.22 \pm 4.22	46.07 \pm 2.16	37.96 \pm 9.50	21.96 \pm 5.20	30.18 \pm 8.79
	CODA-P	59.59 \pm 7.94	22.20 \pm 5.61	49.47 \pm 2.79	41.21 \pm 6.19	18.77 \pm 6.28	28.98 \pm 9.54
	MVP	61.48 \pm 2.55	50.23 \pm 11.27	44.01 \pm 4.18	62.17 \pm 6.10	30.92 \pm 2.26	62.20 \pm 5.81
	MISA	62.87 \pm 3.76	17.34 \pm 6.20	44.62 \pm 3.36	35.62 \pm 11.69	19.48 \pm 4.88	21.18 \pm 10.00
	FlyPrompt (ours)	73.08 \pm 2.18	8.38 \pm 1.52	59.40 \pm 1.84	23.40 \pm 3.20	30.78 \pm 8.41	18.85 \pm 4.15
iBOT-1K	L2P	48.97 \pm 6.18	16.47 \pm 7.80	41.23 \pm 7.31	33.01 \pm 8.96	19.75 \pm 3.57	21.26 \pm 11.09
	DualPrompt	50.16 \pm 4.91	20.45 \pm 3.88	49.51 \pm 1.57	35.14 \pm 7.72	30.29 \pm 3.97	32.70 \pm 9.00
	CODA-P	55.84 \pm 5.28	22.72 \pm 6.32	53.49 \pm 1.21	40.23 \pm 6.19	28.85 \pm 3.98	27.71 \pm 9.65
	MVP	56.41 \pm 2.00	53.47 \pm 12.87	46.96 \pm 3.97	56.26 \pm 5.49	34.97 \pm 1.44	62.61 \pm 6.88
	MISA	52.58 \pm 4.34	19.16 \pm 4.04	48.89 \pm 2.82	33.52 \pm 10.12	27.98 \pm 5.32	29.35 \pm 6.60
	FlyPrompt (ours)	64.94 \pm 2.57	11.52 \pm 3.71	63.77 \pm 1.42	20.89 \pm 3.75	38.05 \pm 6.98	21.62 \pm 3.82
DINO-1K	L2P	45.45 \pm 6.86	14.82 \pm 6.81	38.98 \pm 7.34	33.00 \pm 7.22	23.21 \pm 3.68	24.23 \pm 9.93
	DualPrompt	49.65 \pm 5.46	18.69 \pm 5.27	47.16 \pm 1.05	37.89 \pm 7.65	29.48 \pm 5.88	31.35 \pm 10.61
	CODA-P	50.76 \pm 5.65	19.43 \pm 5.83	48.68 \pm 2.79	41.26 \pm 6.20	29.43 \pm 4.89	32.69 \pm 10.65
	MVP	52.41 \pm 1.48	54.70 \pm 12.47	44.45 \pm 4.11	56.54 \pm 6.00	34.71 \pm 1.92	64.21 \pm 7.82
	MISA	49.81 \pm 3.67	17.76 \pm 4.28	46.61 \pm 2.24	34.38 \pm 10.60	27.46 \pm 5.11	25.77 \pm 10.77
	FlyPrompt (ours)	63.59 \pm 4.56	9.04 \pm 1.73	60.83 \pm 1.57	20.65 \pm 3.08	37.50 \pm 8.33	19.47 \pm 4.84
MoCo-1K	L2P	26.75 \pm 6.56	19.68 \pm 15.19	19.02 \pm 2.75	43.59 \pm 4.64	13.42 \pm 3.89	26.38 \pm 11.63
	DualPrompt	49.61 \pm 7.76	14.79 \pm 4.09	41.37 \pm 1.80	35.07 \pm 5.93	21.87 \pm 4.80	28.42 \pm 9.47
	CODA-P	48.26 \pm 5.59	16.97 \pm 6.15	44.71 \pm 2.43	43.27 \pm 2.27	22.33 \pm 3.94	30.90 \pm 10.77
	MVP	52.82 \pm 2.35	55.42 \pm 15.00	38.68 \pm 3.83	55.87 \pm 4.26	30.74 \pm 2.46	62.97 \pm 5.67
	MISA	53.46 \pm 6.41	16.19 \pm 5.26	45.41 \pm 4.62	34.42 \pm 8.60	26.92 \pm 4.88	35.50 \pm 7.97
	FlyPrompt (ours)	60.16 \pm 6.88	7.69 \pm 2.42	55.95 \pm 1.60	21.27 \pm 2.89	26.52 \pm 6.31	20.69 \pm 3.58

Following Moon et al. (2023); Kang et al. (2025), we consider the standard evaluation metrics A_{auc} , A_{last} , A_{avg} and F_{last} for GCL performance. We denote $R_{i,j}$ as the accuracy recorded right after session i with respect to the data in the session j . We then maintain a matrix \mathbf{R} whose j th column is the history of evaluation after each session with respect to the data in session j . Firstly we calculate the final average accuracy A_{last} as:

$$A_{last} = \frac{1}{T} \sum_{i=1}^T R_{T,i}, \quad (12)$$

average running accuracy A_{avg} :

$$A_{avg} = \frac{1}{T} \sum_{i=1}^T R_{i,i}, \quad (13)$$

and the final average forgetting F_{last} :

$$F_{last} = \frac{1}{T} \sum_{i=1}^T (\max(R_j) - R_{T,i}), \quad (14)$$

where A_{last} and A_{avg} are the higher the better and F_{last} is the lower the better. They are all conventional metrics used in classic CL research.

Moreover, A_{auc} (higher the better) is the anytime inference metric proposed by Koh et al. (2021) to better evaluate the GCL performance during the learning process. Specifically, the evaluation of GCL accuracy is performed every b batches. Throughout this work, we fix $b = 1000$. Denote the total number of the assessments performed as S , then given the history $\mathbf{a} \in \mathbb{R}^S$ (where a_s shows the accuracy of the model at time stamp s over the data it has observed so far), we calculate the area under curve of any-time accuracy A_{auc} as :

$$A_{\text{auc}} = \frac{1}{S} \sum_{s=1}^S a_s. \quad (15)$$

In addition to these metrics, we also report the backward transfer (BWT) metric (Lin et al., 2022) to quantify how learning later sessions influences performance on earlier ones. Using the same notation $R_{i,j}$ and T as above, we define

$$\text{BWT} = \frac{1}{T-1} \sum_{i=1}^{T-1} (R_{T,i} - R_{i,i}). \quad (16)$$

A positive BWT indicates that subsequent learning improves performance on previous sessions (positive backward transfer), while a negative value implies net forgetting on earlier sessions. BWT results of GCL methods are presented in Tab. 12.

F.3 COMPARISON WITH PROMINENT OFFLINE PTM-BASED METHODS

Table 12: Extended Comparison of performance, backward transfer, and computational cost across PTM-based CL methods. All performance results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds, over the CIFAR-100 dataset and Sup-21K backbone. Parameter counts are measured in millions. Time cost is reported in seconds per batch.

Method	$A_{\text{auc}}(\%, \uparrow)$	BWT ($\%, \uparrow$)	Total Param.	Trainable Param.	Training Time	Inference Time
L2P	76.23 \pm 2.73	0.10 \pm 2.65	86.01	0.22	5.57	0.95
DualPrompt	76.04 \pm 3.32	-2.93 \pm 2.42	86.35	0.55	4.78	0.90
CODA-P	79.13 \pm 3.06	-0.83 \pm 2.17	86.72	0.92	4.75	0.94
MVP	67.74 \pm 4.96	-18.09 \pm 3.24	86.12	0.32	5.35	1.27
MISA	80.35 \pm 2.39	-1.76 \pm 2.28	86.37	0.58	4.78	0.90
S-Prompt++	80.21 \pm 2.55	0.81 \pm 1.86	86.26	0.46	6.03	1.18
HiDe-Prompt	77.10 \pm 3.81	3.35 \pm 2.71	86.81	0.94	6.13	1.27
HiDe-LoRA	80.07 \pm 2.41	0.36 \pm 0.94	87.39	1.51	6.80	1.17
HiDe-Adapter	79.52 \pm 2.81	-2.05 \pm 1.95	87.41	1.53	6.68	1.04
NoRGa	78.89 \pm 3.33	2.72 \pm 1.98	86.81	0.94	6.69	1.05
SD-LoRA	79.26 \pm 2.21	-6.66 \pm 3.22	87.72	1.92	7.24	0.82
FlyPrompt (ours)	83.24 \pm 2.23	4.35 \pm 1.19	87.08	0.46	4.96	0.92

Table 13: Comparison between RanPAC variants and FlyPrompt over three GCL benchmarks. All results are reported as an average of five parallel runs (\pm standard deviation) with different seeds.

Method	CIFAR-100		ImageNet-R		CUB-200	
	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
RanPAC [†]	69.91 \pm 3.88	79.92 \pm 0.07	47.14 \pm 2.18	50.75 \pm 2.15	60.18 \pm 5.52	66.21 \pm 6.15
RanPAC [‡]	57.35 \pm 8.23	77.65 \pm 0.21	36.90 \pm 4.17	44.39 \pm 0.11	64.52 \pm 8.23	71.65 \pm 0.17
RanPAC*	77.88 \pm 4.28	86.52 \pm 1.15	53.18 \pm 2.22	54.71 \pm 2.48	69.64 \pm 3.89	72.30 \pm 1.09
FlyPrompt (ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88

F.4 CKA ANALYSIS OF EXPERTS’ REPRESENTATION SIMILARITY

We use centered kernel alignment (CKA) (Kornblith et al., 2019) to quantify the similarity between expert-specific representations while factoring out the shared contribution of the frozen PTM backbone. For a given method and dataset, we first fix the backbone f_θ and, for each expert E_t , apply its prompt (or expert-specific encoder parameters) to obtain a matrix of CLS features $\mathbf{Z}_t \in \mathbb{R}^{n \times d}$

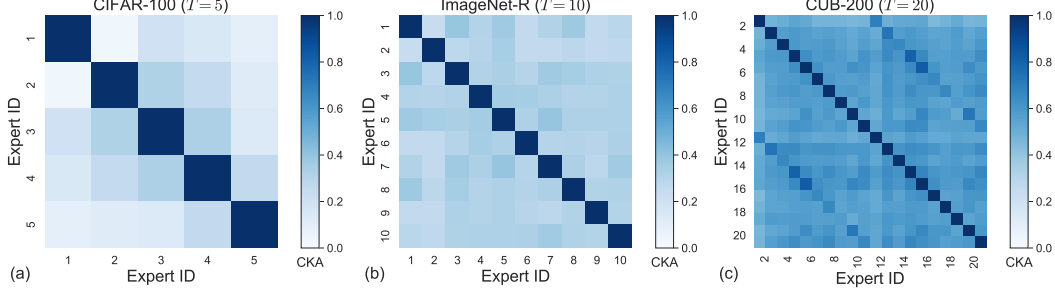


Figure 6: CKA similarity of feature representations between experts of MISA on three datasets.

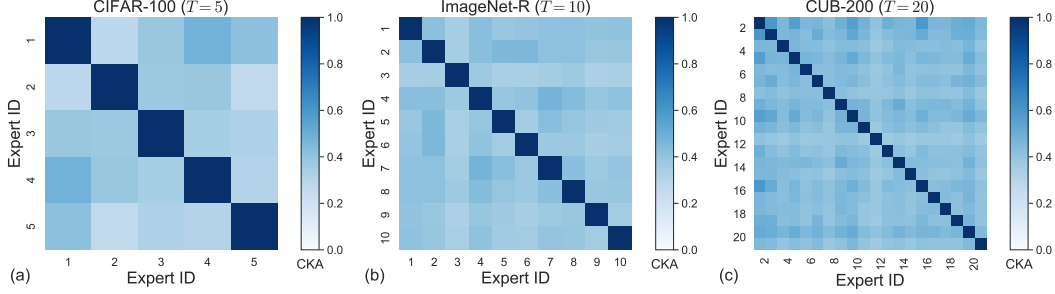


Figure 7: CKA similarity of feature representations between FlyPrompt experts on three datasets.

over a common set of n samples. Using the same backbone together with the expert-shared parameters (e.g., the global prompt in DualPrompt), we also compute a common representation matrix $\mathbf{Z}^{\text{com}} \in \mathbb{R}^{n \times d}$. We then define the residual features of expert E_t as $\tilde{\mathbf{Z}}_t = \mathbf{Z}_t - \mathbf{Z}^{\text{com}}$, which remove the largely stable PTM-driven component and highlight the expert-specific modulation that emerges during online GCL. Based on these residual features, we measure the (linear) CKA between experts E_t and $E_{t'}$ as

$$\text{CKA}(\tilde{\mathbf{Z}}_t, \tilde{\mathbf{Z}}_{t'}) = \frac{\|\tilde{\mathbf{Z}}_t^\top \tilde{\mathbf{Z}}_{t'}\|_F^2}{\|\tilde{\mathbf{Z}}_t^\top \tilde{\mathbf{Z}}_t\|_F \|\tilde{\mathbf{Z}}_{t'}^\top \tilde{\mathbf{Z}}_{t'}\|_F}. \quad (17)$$

This similarity measure is invariant to isotropic rescaling and orthogonal transformations of the features, and is well-suited for comparing representations across experts. We report the pairwise CKA scores between all experts as a heatmap: diagonal entries capture self-similarity, whereas off-diagonal values reveal the degree of specialization or redundancy among experts after removing the common PTM-induced component.

F.5 EXTRA HYPERPARAMETER SENSITIVITY TEST RESULTS

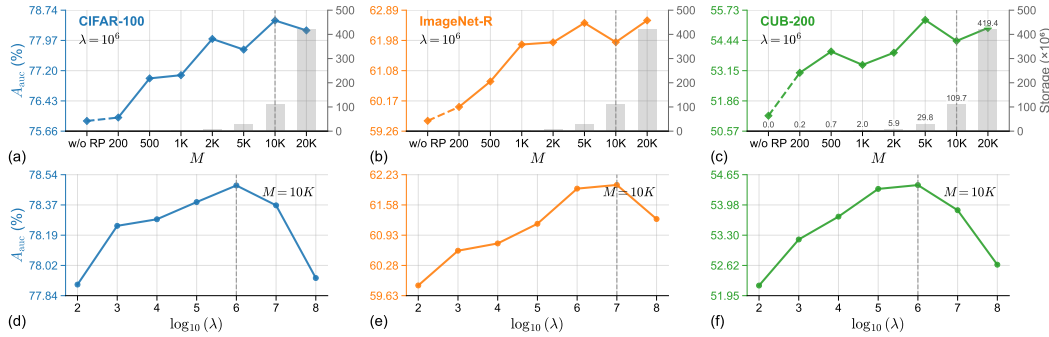


Figure 8: Analysis of hyperparameters in REAR [Backbone: Sup-21K/1K]. (a-c) Different random projection dimension M with fixed $\lambda = 10^6$: we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$.

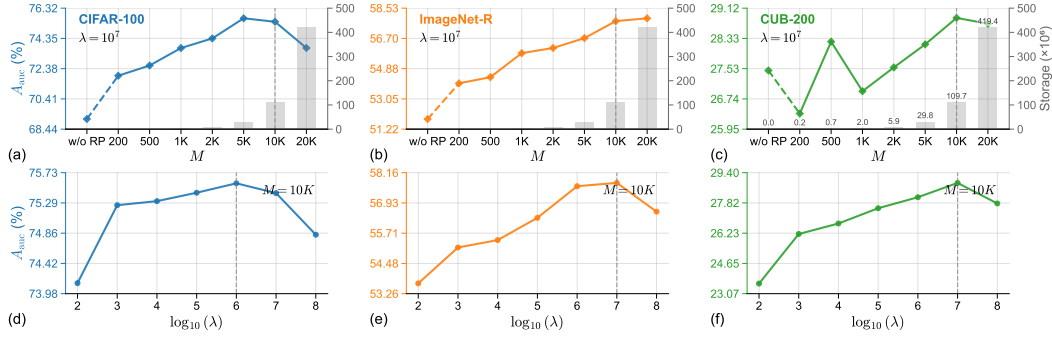


Figure 9: Analysis of hyperparameters in REAR [Backbone: iBOT-21K]. (a-c) Different random projection dimension M with fixed $\lambda = 10^7$: we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$.

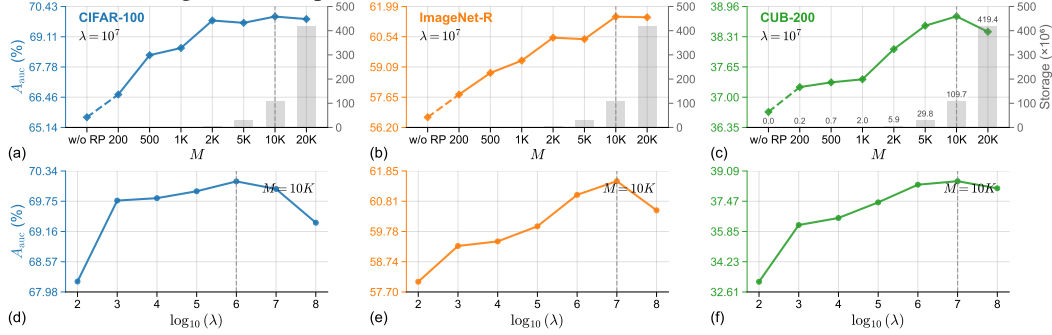


Figure 10: Analysis of hyperparameters in REAR [Backbone: iBOT-1K]. (a-c) Different random projection dimension M with fixed $\lambda = 10^7$: we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$.

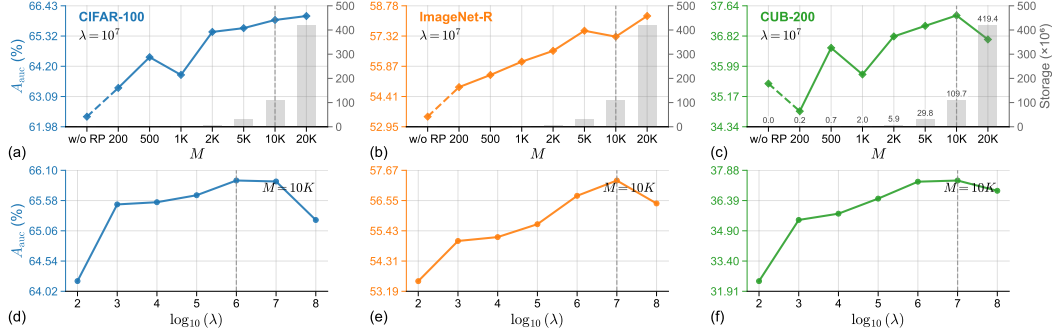


Figure 11: Analysis of hyperparameters in REAR [Backbone: DINO-1K]. (a-c) Different random projection dimension M with fixed $\lambda = 10^7$: we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$.

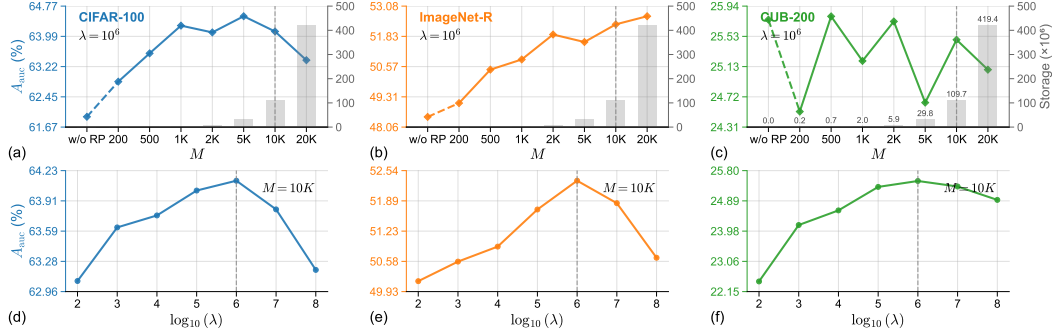


Figure 12: Analysis of hyperparameters in REAR [Backbone: MoCo v3-1K]. (a-c) Different random projection dimension M with fixed $\lambda = 10^6$: we report A_{auc} and extra storage cost (bar) given M . (d-f) Different regularization parameter λ with fixed $M = 10^4$.

F.6 ADDITIONAL ABLATION STUDY

Table 14: Effect of REAR and TE² on various PTM-based CL methods over three GCL benchmarks. All results are reported as an average of five runs (\pm standard deviation) over Sup-21K backbone.

Setup	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Baseline	DualPrompt	76.04 \pm 3.32	76.62 \pm 0.74	46.13 \pm 1.94	40.80 \pm 1.04	65.03 \pm 2.24	62.43 \pm 1.78
	MVP	67.74 \pm 4.96	63.22 \pm 0.69	39.50 \pm 1.41	32.63 \pm 3.95	54.69 \pm 3.14	50.07 \pm 3.86
	MISA	80.35 \pm 2.39	80.75 \pm 1.24	51.52 \pm 2.09	45.08 \pm 1.43	65.46 \pm 3.01	60.20 \pm 1.82
	S-Prompt++	80.21 \pm 2.55	83.48 \pm 1.20	52.14 \pm 1.65	49.13 \pm 1.60	66.61 \pm 2.21	64.73 \pm 2.25
	HiDe-Prompt	77.10 \pm 3.81	81.77 \pm 2.00	53.77 \pm 1.09	49.87 \pm 3.01	67.05 \pm 2.37	67.12 \pm 0.50
	NoRGa	78.89 \pm 3.33	83.03 \pm 1.20	54.12 \pm 1.37	50.09 \pm 3.66	67.16 \pm 2.44	67.06 \pm 0.58
w/ REAR	DualPrompt	80.63 \pm 2.25	83.65 \pm 1.25	53.16 \pm 1.26	51.11 \pm 0.91	65.96 \pm 2.50	63.81 \pm 1.74
	MVP	67.44 \pm 4.89	62.33 \pm 1.62	38.87 \pm 1.27	31.59 \pm 4.19	53.65 \pm 3.17	48.25 \pm 3.38
	MISA	82.03 \pm 1.97	83.82 \pm 1.04	<u>57.30</u> \pm 1.12	54.02 \pm 0.66	68.04 \pm 2.34	65.59 \pm 2.69
	S-Prompt++	81.43 \pm 2.45	83.93 \pm 0.84	54.74 \pm 1.53	52.30 \pm 1.05	66.80 \pm 2.48	64.72 \pm 1.97
	HiDe-Prompt	78.41 \pm 2.64	83.46 \pm 1.14	53.61 \pm 1.16	49.26 \pm 2.56	67.05 \pm 2.36	66.99 \pm 1.01
	NoRGa	79.37 \pm 2.71	83.79 \pm 1.28	54.78 \pm 1.05	50.24 \pm 3.39	67.26 \pm 2.50	67.20 \pm 0.85
w/ TE ²	DualPrompt	76.83 \pm 3.44	78.00 \pm 0.61	47.11 \pm 2.19	42.15 \pm 0.81	66.47 \pm 2.72	65.42 \pm 1.55
	MVP	68.91 \pm 4.86	64.28 \pm 1.00	42.06 \pm 1.11	35.94 \pm 0.92	56.98 \pm 2.79	54.14 \pm 2.89
	MISA	81.65 \pm 2.24	82.80 \pm 1.06	54.05 \pm 1.70	48.46 \pm 1.15	69.30 \pm 2.43	67.29 \pm 2.44
	S-Prompt++	81.93 \pm 2.21	83.98 \pm 0.65	55.37 \pm 1.64	52.91 \pm 1.53	67.97 \pm 2.51	67.68 \pm 1.49
	HiDe-Prompt	77.46 \pm 3.56	82.09 \pm 1.92	54.83 \pm 1.08	50.26 \pm 2.78	67.77 \pm 2.60	69.64 \pm 0.76
	NoRGa	79.16 \pm 3.28	83.01 \pm 1.45	54.08 \pm 1.58	51.81 \pm 3.51	67.97 \pm 2.70	69.58 \pm 0.90
w/ both	DualPrompt	82.33 \pm 2.17	86.14 \pm 0.90	54.72 \pm 1.29	54.00 \pm 0.76	69.65 \pm 2.93	72.75 \pm 2.03
	MVP	68.93 \pm 4.60	64.52 \pm 1.45	41.59 \pm 1.34	35.45 \pm 1.44	55.65 \pm 2.79	51.84 \pm 2.44
	MISA	83.60 \pm 2.08	86.66 \pm 0.55	59.12 \pm 1.02	56.62 \pm 0.82	72.38 \pm 2.98	74.68 \pm 2.11
	S-Prompt++	83.11 \pm 2.30	<u>86.67</u> \pm 0.45	56.57 \pm 1.48	55.24 \pm 1.24	<u>70.65</u> \pm 2.84	<u>73.42</u> \pm 1.88
	HiDe-Prompt	78.60 \pm 2.53	83.14 \pm 1.12	54.79 \pm 1.13	51.30 \pm 2.81	68.27 \pm 2.57	69.61 \pm 0.84
	NoRGa	79.37 \pm 2.74	83.79 \pm 0.96	55.75 \pm 1.31	52.50 \pm 3.33	68.32 \pm 2.64	70.03 \pm 0.67
	FlyPrompt (ours)	<u>83.24</u> \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	<u>55.27</u> \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88

F.7 DIFFERENT LOGIT MASK STRATEGIES FOR GCL MASKS

Table 15: Performance of different logit mask strategies for GCL methods. All results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds, over Sup-21K.

Mask Type	Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
No Mask	L2P	62.74 \pm 4.39	56.08 \pm 2.10	34.58 \pm 1.23	26.18 \pm 4.69	54.83 \pm 2.65	47.94 \pm 4.64
	DualPrompt	66.68 \pm 5.25	61.98 \pm 4.09	41.62 \pm 1.43	<u>36.08</u> \pm 1.39	56.68 \pm 2.57	50.71 \pm 4.21
	CODA-P	66.15 \pm 5.22	58.42 \pm 1.39	40.71 \pm 3.50	30.56 \pm 5.15	56.44 \pm 2.84	49.50 \pm 6.04
	MVP	68.18 \pm 4.85	63.96 \pm 1.48	38.79 \pm 1.12	32.01 \pm 2.80	54.74 \pm 2.01	<u>52.88</u> \pm 3.08
	MISA	<u>69.85</u> \pm 3.73	<u>65.13</u> \pm 1.70	<u>45.15</u> \pm 1.75	35.91 \pm 3.48	<u>57.87</u> \pm 2.49	52.15 \pm 4.27
	FlyPrompt (ours)	78.73 \pm 3.55	83.62 \pm 0.50	51.39 \pm 1.80	48.72 \pm 1.06	69.22 \pm 3.04	73.07 \pm 2.34
Random Mask	L2P	62.46 \pm 4.54	54.67 \pm 1.39	33.10 \pm 1.34	24.73 \pm 4.52	52.92 \pm 2.51	45.05 \pm 4.89
	DualPrompt	65.48 \pm 4.45	59.66 \pm 1.70	36.94 \pm 1.77	29.22 \pm 1.66	55.07 \pm 2.21	47.83 \pm 5.00
	CODA-P	65.58 \pm 5.23	57.05 \pm 1.97	39.23 \pm 2.80	28.91 \pm 5.22	55.64 \pm 2.27	48.00 \pm 5.50
	MVP	67.75 \pm 4.95	<u>63.21</u> \pm 0.78	39.45 \pm 1.42	<u>32.70</u> \pm 3.96	54.72 \pm 3.14	<u>50.01</u> \pm 3.81
	MISA	<u>68.23</u> \pm 3.82	61.55 \pm 2.02	<u>40.67</u> \pm 1.93	29.48 \pm 4.32	<u>56.06</u> \pm 2.19	48.12 \pm 4.86
	FlyPrompt (ours)	78.32 \pm 3.48	81.88 \pm 0.88	51.67 \pm 1.30	47.26 \pm 1.25	68.63 \pm 2.82	69.52 \pm 4.03
Seen-Class Mask	L2P	62.22 \pm 4.39	53.36 \pm 2.03	33.80 \pm 1.22	25.20 \pm 4.65	53.10 \pm 2.60	45.55 \pm 4.96
	DualPrompt	65.29 \pm 4.62	57.74 \pm 2.53	37.31 \pm 1.80	29.72 \pm 1.49	55.25 \pm 2.43	47.67 \pm 4.65
	CODA-P	65.63 \pm 5.40	56.75 \pm 1.51	40.13 \pm 2.46	29.35 \pm 5.01	55.80 \pm 2.58	47.71 \pm 5.55
	MVP	67.72 \pm 4.87	<u>62.99</u> \pm 0.95	39.57 \pm 1.44	<u>32.72</u> \pm 4.00	54.72 \pm 3.14	<u>50.14</u> \pm 3.80
	MISA	<u>68.34</u> \pm 3.90	61.44 \pm 2.54	<u>41.17</u> \pm 1.85	29.97 \pm 4.14	<u>56.44</u> \pm 2.42	48.58 \pm 4.74
	FlyPrompt (ours)	78.75 \pm 3.52	82.87 \pm 0.82	52.39 \pm 1.46	48.02 \pm 1.00	69.28 \pm 2.95	70.91 \pm 2.83
Batch Seen-Class Mask	L2P	76.23 \pm 2.73	79.11 \pm 1.43	44.40 \pm 1.03	42.03 \pm 1.72	64.30 \pm 2.18	61.42 \pm 2.13
	DualPrompt	76.04 \pm 3.32	76.62 \pm 0.74	46.13 \pm 1.94	40.80 \pm 1.04	65.03 \pm 2.24	62.43 \pm 1.78
	CODA-P	79.13 \pm 3.06	<u>80.91</u> \pm 0.70	<u>51.87</u> \pm 2.81	<u>48.09</u> \pm 2.75	<u>66.01</u> \pm 2.20	<u>62.90</u> \pm 2.46
	MVP	67.74 \pm 4.96	63.22 \pm 0.69	39.50 \pm 1.41	32.63 \pm 3.95	54.69 \pm 3.14	50.07 \pm 3.86
	MISA	<u>80.35</u> \pm 2.39	80.75 \pm 1.24	51.52 \pm 2.09	45.08 \pm 1.43	65.40 \pm 3.01	60.20 \pm 1.82
	FlyPrompt (ours)	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88

(1) **No Mask**, standard softmax over all output classes. (2) **Random Mask**, for each sample (x, y) , set $m_y = 0$ and assign $m_c = 0$ or $-\infty$ randomly with 0.5 probability for each previously seen class $c \neq y$. (3) **Seen-Class Mask**, setting $m_c = 0$ for all previously seen classes, $-\infty$ otherwise. (4) **Batch Seen-Class Mask** (used), setting $m_c = 0$ only for the classes present in the current batch y .

F.8 COMPLETE RESULTS OF DIFFERENT EMA DECAY RATES FOR TEMPORAL ENSEMBLE

Table 16: Performance comparison of different EMA decay rates for TE² across all PTMs. All results are reported as an average of five parallel runs (\pm standard deviation) with different seeds.

PTM	EMA Decay Rate	CIFAR-100		ImageNet-R	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Sup-21K	online	81.90 \pm 2.20	84.23 \pm 1.32	54.91 \pm 1.32	52.58 \pm 1.36
	0.9	82.81 \pm 2.28	86.36 \pm 0.54	56.36 \pm 1.52	55.09 \pm 0.89
	0.99	82.84 \pm 2.51	86.41 \pm 0.39	55.94 \pm 1.65	54.67 \pm 0.89
	0.999	81.80 \pm 2.37	84.39 \pm 0.83	55.15 \pm 1.39	53.52 \pm 0.80
	0.9,0.99	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91
	0.9,0.99,0.999	82.99 \pm 2.22	86.24 \pm 0.79	56.35 \pm 1.72	55.50 \pm 0.77
Sup-21K/1K	online	71.28 \pm 2.58	69.73 \pm 5.78	53.12 \pm 2.19	44.69 \pm 3.65
	0.9	75.59 \pm 2.93	77.39 \pm 6.14	61.56 \pm 1.48	57.35 \pm 1.63
	0.99	77.96 \pm 2.15	79.71 \pm 3.67	60.96 \pm 1.94	56.32 \pm 1.87
	0.999	74.13 \pm 1.64	74.68 \pm 2.93	53.96 \pm 1.30	46.55 \pm 1.17
	0.9,0.99	78.48 \pm 1.31	80.39 \pm 3.54	62.01 \pm 2.32	56.55 \pm 3.94
	0.9,0.99,0.999	78.44 \pm 1.38	80.55 \pm 4.12	62.59 \pm 2.22	58.00 \pm 1.79
iBOT-21K	online	67.01 \pm 3.85	65.38 \pm 7.24	45.32 \pm 1.46	35.45 \pm 4.60
	0.9	72.43 \pm 1.56	75.46 \pm 4.96	56.37 \pm 2.06	52.37 \pm 0.87
	0.99	75.03 \pm 0.78	78.08 \pm 3.35	55.93 \pm 2.35	51.64 \pm 0.53
	0.999	68.96 \pm 3.22	69.30 \pm 2.94	43.78 \pm 2.21	34.64 \pm 2.90
	0.9,0.99	75.58 \pm 1.70	79.36 \pm 3.47	57.75 \pm 2.12	54.39 \pm 1.29
	0.9,0.99,0.999	74.16 \pm 2.47	76.61 \pm 2.26	55.94 \pm 3.12	52.12 \pm 0.79
iBOT-1K	online	61.38 \pm 2.32	60.58 \pm 7.91	50.79 \pm 1.43	41.92 \pm 1.76
	0.9	68.01 \pm 1.18	71.99 \pm 4.62	60.52 \pm 1.51	56.98 \pm 1.17
	0.99	71.50 \pm 1.00	75.20 \pm 3.10	60.66 \pm 1.56	56.57 \pm 0.70
	0.999	65.88 \pm 3.51	67.07 \pm 1.95	50.27 \pm 1.40	42.71 \pm 2.03
	0.9,0.99	70.14 \pm 1.76	74.84 \pm 4.26	61.50 \pm 1.66	57.18 \pm 1.36
	0.9,0.99,0.999	67.93 \pm 3.07	70.69 \pm 3.50	59.60 \pm 1.88	55.35 \pm 1.34
DINO-1K	online	58.61 \pm 3.26	60.76 \pm 6.77	47.35 \pm 2.24	41.33 \pm 2.04
	0.9	62.95 \pm 4.13	69.65 \pm 6.54	56.83 \pm 1.47	53.68 \pm 0.83
	0.99	66.42 \pm 2.51	73.03 \pm 3.61	56.67 \pm 1.74	53.23 \pm 0.77
	0.999	60.37 \pm 4.38	64.06 \pm 2.37	46.38 \pm 1.51	39.69 \pm 2.07
	0.9,0.99	65.92 \pm 2.74	72.66 \pm 4.52	57.29 \pm 2.40	54.72 \pm 1.89
	0.9,0.99,0.999	65.27 \pm 2.98	70.83 \pm 4.32	55.66 \pm 1.57	51.91 \pm 1.73
MoCo-1K	online	57.90 \pm 5.29	62.20 \pm 10.03	42.81 \pm 0.83	35.46 \pm 3.32
	0.9	61.96 \pm 6.50	69.83 \pm 10.05	51.47 \pm 1.64	47.88 \pm 1.49
	0.99	65.95 \pm 4.40	73.28 \pm 6.96	50.75 \pm 1.89	47.42 \pm 1.29
	0.999	61.26 \pm 3.27	66.42 \pm 5.07	42.33 \pm 1.54	36.82 \pm 2.14
	0.9,0.99	64.12 \pm 5.18	71.51 \pm 8.48	52.32 \pm 1.50	49.06 \pm 1.35
	0.9,0.99,0.999	64.03 \pm 4.47	69.92 \pm 6.13	51.64 \pm 2.59	48.69 \pm 0.68

F.9 COMPARISON OF DIFFERENT ROUTING ALGORITHMS ON RANDOM EXPANDED FEATURES

Table 17: Comparison of routing algorithms based on random expanded features in FlyPrompt. M : expansion dimension (default 10,000); T : number of experts (default 5); H : hidden dimension of MLP ($H = 512$ is used); K : number of nearest neighbors ($K = 10$ is used). Time cost is reported in seconds per batch on CIFAR-100. All performance results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds over Sup-21K.

Routing Algorithm	Train Time	Inference Time	Inference Complexity	CIFAR-100		ImageNet-R		CUB-200	
				$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Prototype Similarity	5.58	0.90	$O(MT)$	80.67 \pm 2.48	83.80 \pm 1.15	54.29 \pm 1.72	52.36 \pm 1.12	67.00 \pm 2.77	66.66 \pm 1.56
Naïve Bayes	5.30	0.93	$O(MT)$	82.73 \pm 2.17	85.51 \pm 0.97	55.85 \pm 1.83	53.84 \pm 1.40	69.08 \pm 2.91	69.63 \pm 1.18
MLP	7.03	1.00	$O(MH + HT)$	81.75 \pm 2.09	82.76 \pm 1.98	56.31 \pm 1.29	53.70 \pm 0.96	68.53 \pm 2.39	67.92 \pm 1.91
K-Means	6.11	1.49	$O(KMT)$	82.22 \pm 2.04	85.27 \pm 0.83	54.93 \pm 1.94	53.08 \pm 1.43	68.33 \pm 2.71	68.24 \pm 2.59
Ridge Regression (ours)	4.96	0.92	$O(MT)$	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88

(1) **Prototype Similarity**, cosine similarity to each expert’s mean feature. (2) **Naïve Bayes**, assuming Gaussian-distributed features per expert. (3) **MLP**, a two-layer MLP router, as in HiDe (Wang et al., 2023a). (4) **K-Means**, clusters each expert’s features and routes based on the nearest center. (5) **Ridge Regression (Ours)**, the REAR analytic router trained once over accumulated statistics.

F.10 COMPLETE RESULTS OF DIFFERENT AGGREGATION METHODS FOR TEMPORAL ENSEMBLE

Table 18: Performance comparison of different aggregation choices for TE² across all PTMs. All results are reported as an average of five parallel runs (\pm standard deviation) with different random seeds.

PTM	Ensemble Method	CIFAR-100		ImageNet-R		CUB-200	
		$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$	$A_{\text{auc}}(\%, \uparrow)$	$A_{\text{last}}(\%, \uparrow)$
Sup-21K	Mean	81.34 \pm 1.64	85.11 \pm 1.03	52.71 \pm 1.36	53.24 \pm 1.22	68.49 \pm 2.57	73.95 \pm 1.90
	Max Prob	82.29 \pm 2.25	84.95 \pm 1.20	55.56 \pm 1.38	53.53 \pm 1.40	68.00 \pm 2.50	66.56 \pm 1.60
	Min Entropy	81.92 \pm 2.19	84.23 \pm 1.32	55.05 \pm 1.31	52.88 \pm 1.38	66.78 \pm 2.53	64.73 \pm 1.36
	SoftMax+Mean	82.30 \pm 1.82	85.98 \pm 0.80	56.16 \pm 1.56	55.53 \pm 0.89	70.77 \pm 3.00	74.86 \pm 1.54
	SoftMax+Max	83.24 \pm 2.23	86.76 \pm 0.73	56.58 \pm 1.47	55.27 \pm 0.91	70.64 \pm 2.85	73.40 \pm 1.88
	SoftMax+Min	83.11 \pm 2.34	86.50 \pm 0.64	55.94 \pm 1.41	54.24 \pm 1.34	69.86 \pm 2.80	71.51 \pm 1.79
Sup-21K/1K	Mean	79.44 \pm 1.14	82.82 \pm 1.63	60.84 \pm 2.09	56.97 \pm 2.72	56.64 \pm 4.59	60.39 \pm 3.95
	Max Prob	72.58 \pm 2.99	71.31 \pm 7.52	54.28 \pm 2.06	46.44 \pm 3.23	47.05 \pm 3.41	44.42 \pm 3.59
	Min Entropy	71.27 \pm 2.58	69.73 \pm 5.82	53.05 \pm 2.11	44.80 \pm 3.42	45.36 \pm 3.15	42.82 \pm 4.31
	SoftMax+Mean	77.72 \pm 1.48	81.77 \pm 3.15	60.35 \pm 1.86	56.63 \pm 2.05	57.68 \pm 5.02	62.60 \pm 4.05
	SoftMax+Max	78.48 \pm 1.31	80.39 \pm 3.54	62.01 \pm 2.32	56.55 \pm 3.94	54.42 \pm 4.67	55.50 \pm 3.55
	SoftMax+Min	78.30 \pm 1.25	80.07 \pm 3.68	60.89 \pm 2.15	55.59 \pm 3.37	54.12 \pm 4.61	54.75 \pm 3.57
iBOT-21K	Mean	76.59 \pm 1.33	81.40 \pm 2.48	54.92 \pm 1.86	51.82 \pm 1.97	29.31 \pm 5.17	37.60 \pm 4.77
	Max Prob	68.51 \pm 4.08	67.60 \pm 8.11	46.74 \pm 1.39	37.64 \pm 4.20	24.14 \pm 3.53	28.53 \pm 3.46
	Min Entropy	67.03 \pm 3.85	65.44 \pm 7.22	45.33 \pm 1.36	35.86 \pm 4.55	23.44 \pm 3.39	27.40 \pm 3.60
	SoftMax+Mean	74.69 \pm 1.77	80.58 \pm 2.75	54.41 \pm 1.95	52.58 \pm 1.10	29.09 \pm 6.33	36.28 \pm 7.54
	SoftMax+Max	75.58 \pm 1.70	79.36 \pm 3.47	57.75 \pm 2.12	54.39 \pm 1.29	28.86 \pm 5.84	36.79 \pm 7.58
	SoftMax+Min	74.87 \pm 1.89	77.60 \pm 4.71	55.98 \pm 1.90	52.03 \pm 1.61	27.47 \pm 5.43	34.57 \pm 5.12
iBOT-1K	Mean	70.96 \pm 1.13	76.38 \pm 4.03	58.15 \pm 1.53	54.24 \pm 1.39	37.64 \pm 5.03	44.46 \pm 3.01
	Max Prob	62.92 \pm 2.41	63.09 \pm 8.29	52.02 \pm 1.32	44.15 \pm 1.01	31.53 \pm 3.46	34.27 \pm 4.92
	Min Entropy	61.36 \pm 2.30	60.73 \pm 7.76	50.74 \pm 1.30	42.27 \pm 1.46	30.53 \pm 3.62	32.82 \pm 4.98
	SoftMax+Mean	67.24 \pm 1.92	72.80 \pm 4.63	56.14 \pm 1.33	53.35 \pm 0.95	38.08 \pm 5.57	44.13 \pm 4.27
	SoftMax+Max	70.14 \pm 1.76	74.84 \pm 4.26	61.50 \pm 1.66	57.18 \pm 1.36	38.54 \pm 5.72	45.00 \pm 4.19
	SoftMax+Min	69.72 \pm 1.64	73.76 \pm 5.07	59.87 \pm 1.52	55.39 \pm 0.90	37.31 \pm 5.53	41.86 \pm 4.01
DINO-1K	Mean	66.72 \pm 1.89	73.94 \pm 3.66	54.00 \pm 2.34	51.60 \pm 2.21	37.91 \pm 6.38	44.02 \pm 2.34
	Max Prob	59.76 \pm 3.26	62.61 \pm 7.14	48.12 \pm 2.33	42.17 \pm 2.21	31.36 \pm 3.40	34.50 \pm 4.80
	Min Entropy	58.51 \pm 3.22	60.79 \pm 6.67	47.14 \pm 2.31	40.95 \pm 2.13	30.15 \pm 3.34	32.83 \pm 4.91
	SoftMax+Mean	64.79 \pm 3.29	72.87 \pm 4.57	52.98 \pm 1.49	51.38 \pm 0.95	37.22 \pm 7.26	43.74 \pm 5.34
	SoftMax+Max	65.92 \pm 2.74	72.66 \pm 4.52	57.29 \pm 2.40	54.72 \pm 1.89	37.38 \pm 5.86	44.66 \pm 2.35
	SoftMax+Min	65.48 \pm 2.69	71.88 \pm 5.38	55.69 \pm 2.51	52.55 \pm 1.90	36.48 \pm 5.73	41.43 \pm 2.14
MoCo-1K	Mean	64.29 \pm 6.09	72.17 \pm 8.65	49.69 \pm 1.28	46.86 \pm 0.90	27.92 \pm 5.19	33.32 \pm 3.58
	Max Prob	58.70 \pm 5.09	63.89 \pm 10.06	44.07 \pm 0.90	37.06 \pm 3.05	21.35 \pm 3.06	24.21 \pm 4.51
	Min Entropy	57.84 \pm 5.27	62.21 \pm 9.97	42.79 \pm 0.82	35.51 \pm 3.34	20.59 \pm 2.95	22.75 \pm 4.42
	SoftMax+Mean	62.90 \pm 6.00	71.71 \pm 8.20	50.56 \pm 1.57	47.95 \pm 0.70	26.95 \pm 5.22	31.50 \pm 3.85
	SoftMax+Max	64.12 \pm 5.18	71.51 \pm 8.48	52.32 \pm 1.50	49.06 \pm 1.35	25.49 \pm 4.53	30.44 \pm 4.50
	SoftMax+Min	63.92 \pm 5.03	71.29 \pm 8.50	51.46 \pm 1.44	47.93 \pm 1.27	25.26 \pm 4.52	29.33 \pm 3.88

F.11 MORE RESULTS ON SCALABILITY AND EFFICIENCY OF FLYPROMPT

Table 19: Parameter counts, storage and computational complexity breakdown of FlyPrompt components. M : expansion dimension (default 10,000); T : number of experts (default 5); l : prompt length (default 20); d : embedding dimension (default 768). Results are reported in millions.

Components	Total Param.	Trainable Param.	Storage	Storage Cost	Computation Cost
G matrix	0.00	0.00	100	$O(M^2)$	$O(M^3)$
Q matrix	0.00	0.00	0.05	$O(MT)$	$O(MT)$
Router Head	0.05	0.00	0.05	$O(MT)$	$O(MT)$
Prompts	0.38	0.38	0.38	$O(ld)$	$O(l^2d)$
TE ² heads	0.77	0.08	0.77	$O(dT)$	$O(dT)$

G PSEUDO-CODE OF FLYPROMPT

G.1 PSEUDO-CODE: ONLINE REAR UPDATES AND TE² AGGREGATION

Algorithm 1 FlyPrompt: online REAR maintenance and TE² inference

```

1: Inputs: sessions  $\{\mathcal{D}_t\}_{t=1}^T$ , backbone  $f_\theta$ , random matrix  $\mathbf{R} \in \mathbb{R}^{d \times M}$ , ridge  $\lambda$ , EMA decays  $\{\alpha_j\}_{j=1}^n$ , online iterations  $k$ .
2: Initialize:  $\mathbf{G} \leftarrow \mathbf{0}_{M \times M}$ ,  $\mathbf{Q} \leftarrow \mathbf{0}_{M \times T}$ , prompt set  $\mathcal{P} \leftarrow \emptyset$ , online head  $(\mathbf{W}, \mathbf{b})$ , logit mask  $\mathbf{m} \leftarrow \mathbf{0}$ 
3:
4: (1) Online Training Phase
5: for  $t=1$  to  $T$  do
6:   Set expert  $E_t$ :  $\mathbf{p}_t \leftarrow \frac{1}{t-1} \sum_{i=1}^{t-1} \mathbf{p}_i$  if  $t > 1$  else random
7:    $\mathcal{P} \leftarrow \mathcal{P} \cup \{\mathbf{p}_t\}$ 
8:   Initialize EMA heads:  $(\mathbf{W}_t^{(j)}, \mathbf{b}_t^{(j)}) \leftarrow (\mathbf{W}, \mathbf{b})$ ,  $\forall j \in \{1, \dots, n\}$ 
9:   for each batch  $(\mathbf{X}, \mathbf{y}) \in \mathcal{D}_t$  do
10:    Set logit mask  $\mathbf{m}$ : for any class  $c \in \mathbf{y}$ ,  $m_c \leftarrow 0$ , and for  $c' \notin \mathbf{y}$ ,  $m_{c'} \leftarrow -\infty$ 
11:    for 1 to  $k$  do
12:      Update  $(\mathbf{W}, \mathbf{b})$  and  $\mathbf{p}_t$  by minimizing  $\text{CE}(f_\theta(\mathbf{X}; \mathbf{p}_t) \mathbf{W}^\top + \mathbf{b} + \mathbf{m}, \mathbf{y})$ 
13:      EMA for  $E_t$ :  $\mathbf{W}_t^{(j)} \leftarrow \alpha_j \mathbf{W}_t^{(j)} + (1 - \alpha_j) \mathbf{W}$ ,  $\mathbf{b}_t^{(j)} \leftarrow \alpha_j \mathbf{b}_t^{(j)} + (1 - \alpha_j) \mathbf{b}$ ,  $\forall j$ 
14:    end for
15:     $\mathbf{H} \leftarrow f_\theta(\mathbf{X}; \mathbf{p}_t)$ ;  $\Phi \leftarrow \sigma(\mathbf{H} \mathbf{R})$ 
16:    Update REAR stats:  $\mathbf{G} \leftarrow \mathbf{G} + \Phi^\top \Phi$ ;  $\mathbf{Q} \leftarrow \mathbf{Q} + \Phi^\top \mathbf{C}_t$  (one-hot  $\mathbf{C}_t$  for expert  $E_t$ )
17:  end for
18: end for
19: Update closed-form router offline:  $\hat{\mathbf{U}}^\top \leftarrow (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{Q}$  (only once after the training phase)
20:
21: (2) Inference Phase
22: for  $\mathbf{x}$  from the test dataset do
23:   Get routing score:  $\mathbf{s}(\mathbf{x}) \leftarrow \sigma(f_\theta(\mathbf{x}) \mathbf{R}) \hat{\mathbf{U}}^\top$ 
24:   Select the expert:  $e \leftarrow \arg \max_{t \leq T} s_t(\mathbf{x})$ 
25:   Get online head outputs:  $\mathbf{z}^{(0)} \leftarrow f_\theta(\mathbf{x}; \mathbf{p}_e) \mathbf{W}^\top + \mathbf{b}$ 
26:   Get EMA heads output:  $\mathbf{z}^{(j)} \leftarrow f_\theta(\mathbf{x}; \mathbf{p}_e) \mathbf{W}_e^{(j)\top} + \mathbf{b}_e^{(j)}$ 
27:   Get aggregated ensemble output:  $\hat{\mathbf{z}}(\mathbf{x}) \leftarrow \max_{j \in \{0, \dots, n\}} \text{softmax}(\mathbf{z}^{(j)} + \mathbf{m})$ 
28:   Final output:  $\hat{y}(\mathbf{x}) = \arg \max_c \hat{z}_c(\mathbf{x})$ 
29: end for

```

G.1.1 COMPLEXITY

- Memory for REAR: storing $\mathbf{G} \in \mathbb{R}^{M \times M}$ and $\mathbf{Q} \in \mathbb{R}^{M \times T}$ is $O(M^2)$.
- Time complexity of solving the analytic router: $O(M^3)$ (matrix inverse and multiply).
- Per-sample cost: forming $\varphi(\mathbf{x}) = \sigma(\mathbf{h}^\top \mathbf{R})$ costs $O(dM)$ (matrix multiply).