REVISITING PROMPT-BASED METHODS IN CLASS IN CREMENTAL LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

In recent years, prompt-based methods have emerged as a promising direction for continual learning, demonstrating impressive performance across various benchmarks. These methods create learnable prompts to infer task identity, then select and integrate specific prompts into the pretrained model to generate instructed features for prediction. In this paper, we first analyze the working patterns of such method across different distribution scenarios through extensive empirical analysis. Our analysis exposes the limitations of existing methods: first, two-stage inference can make mistakes even when the first stage has already provided reliable predictions; second, enforcing identical architectures for both stages hampers performance gains. To address these issues, we incorporated a self-supervised learning objective to learn discriminative features, thereby boosting the plasticity of the model. During inference, we implemented a simple yet effective threshold filtering strategy to selectively pass data to the second stage. This approach prevents errors in the second stage when the first stage has already made reliable predictions, while also conserving computational resources. Ultimately, we explore utilizing self-supervised pretrained models as a unified task identity provider. Comparing to state-of-the-art methods, our method achieves comparable results under indistribution scenarios and demonstrates substantial gains under out-of-distribution scenarios (e.g., up to 6.34% and 5.15% improvements on Split Aircrafts and Split Cars-196, respectively).

032

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

1 INTRODUCTION

033 Continual learning (CL) aims to equip models with the ability to constantly learn new knowledge without forgetting previously acquired knowledge, with the main challenge being how to mitigate catastrophic forgetting McCloskey & Cohen (1989); Goodfellow et al. (2013) occurs in deep neural 035 networks. Catastrophic forgetting refers to a phenomenon in which introducing new information causes the model to forget old knowledge, thereby drastically reducing its performance on previous 037 learned tasks. Early studies Zenke et al. (2017); Aljundi et al. (2018); Chaudhry et al. (2018) generally start with neural networks that are randomly initialized, focusing primarily on class incremental learning (CIL) scenario Van de Ven & Tolias (2019); De Lange et al. (2021); Masana et al. (2022); 040 Wang et al. (2024b), where each incremental stage involves non-overlapping categories. With the 041 emergence of large-scale pretrained models Dosovitskiy (2020); Kolesnikov et al. (2020); Yalniz 042 et al. (2019); Caron et al. (2021); Oquab et al. (2023), extensive efforts have started to apply them to 043 CL. Among these, a series of works Wang et al. (2022b;a); Smith et al. (2023); Chen et al. (2023); 044 Wang et al. (2024a) are based on prompt learning, which learns and retrieves task-related prompts during training and inference, demonstrating remarkably superior performance.

Prompt-based CL methods typically involve creating a set of learnable prompts, which are then optimized throughout the sequential learning of tasks. The core idea of these methods is to establish a unified query mechanism during training and inference, where the most relevant prompt to the current input are identified and incorporated into the pretrained model to generate instructed features for prediction. As shown in Figure 1, some methods Wang et al. (2022a;b) utilize a frozen vision transformer (ViT) to obtain uninstructed features and select prompts based on the cosine similarity between these features and learnable keys, which are associated with the prompts. The matching result is implicitly treated as inferred task identity (in gray). While the newly proposed HiDe-Prompt Wang et al. (2024a) trains a separate model to predict labels, explicitly incorporating task

069





Figure 2: Performance comparison of different prompt-based CL methods on Split CIFAR-100 (ID) and Split Aircrafts (OOD) using iBOT model pretrained on ImageNet-21K, all datasets were splitted into ten incremental tasks.

Figure 1: Illustration of prompt-based CL methods.

070 identity information through a label-to-task mapping to choose prompts (in beige).¹ It significantly 071 bridges the performance gap between self-supervised and supervised pretrained models. However, 072 our empirical analysis reveals that these methods still fall short under out-of-distribution (OOD) 073 scenarios, especially in comparison to their remarkably high performance under in-distribution (ID) 074 scenarios, as shown in Figure 2.

075 In this paper, we first conducted a series of empirical analyses to thoroughly investigate the behav-076 ior of the existing prompt-based CL methods across two different scenarios (i.e. ID and OOD). 077 The main observations are as follows: (1) self-supervised and supervised pretrained models per-078 form similarly in ID scenarios, but self-supervised models excel in leveraging task identity in OOD 079 settings. (2) A reliable task identity, which is model-agnostic, benefits the second stage without requiring identical architectures for both stages. Based on these observations, we first incorporated self-distillation loss while training first-stage model to enhance the accuracy of its predicted task 081 identities. During inference, we found not all data requires two-stage inference, especially when the first-stage model has already provided reliable predictions. We first modeled the prediction 083 confidence of the first-stage model using a β distribution, and then proposed a simple yet effec-084 tive threshold filtering strategy. This approach mitigates issues encountered under OOD scenarios, 085 where second-stage model predicts incorrectly even with the correct task identity. Concretely, we 086 chose the lowest boundary of a specific highest density interval of the distribution as the threshold 087 for making decisions. We use this threshold to selectively send data to the second-stage model. Ul-880 timately, we explored leveraging self-supervised pretrained models to improve the accuracy of task 089 identity inference, thereby enhancing the models' continual learning capabilities. 090

The main contributions of this paper are summarized as follows: (1) We revisited existing prompt-091 based CL methods, analyzing the performance differences these methods exhibit across different 092 scenarios and their main limitations. (2) We improved the first-stage model's adaptability with 093 self-supervised learning, implemented a threshold filtering strategy to reduce second-stage errors 094 from reliable predictions, and used self-supervised pretrained models as a unified task identifier, 095 eliminating the need for identical architectures. (3) We achieved substantial improvements over 096 state-of-the-art methods in OOD benchmarks and obtained comparable results in ID benchmarks.

098

099

2 **RELATED WORK**

100 Continual Learning. Continual learning has received increasing attention from researchers. Early 101 approaches Zhu et al. (2021); Yu et al. (2020; 2022); Liu et al. (2022); Tao et al. (2024) tended to 102 train a deep neural network from scratch, broadly divided into three categories. The first is replay-103 based methods, which alleviate forgetting by replaying stored previous exemplars Rebuffi et al. 104 (2017); Hou et al. (2019) or generated pseudo samples Shin et al. (2017); Ostapenko et al. (2019).

¹⁰⁵

¹⁰⁶ ¹For clarity, we refer to the model inferring task identity as the "first-stage model" and the model that uses 107 selected prompts to generate instructed features as the "second-stage model." The following text will follow this terminology.

The second is regularization-based methods, which typically apply constraints to the model Li & Hoiem (2017) by combining knowledge distillation Hinton (2015); Yu et al. (2019) techniques or restricting changes of important parameters Kirkpatrick et al. (2017). The third is architecture-based methods, which either expand the model when learning new tasks Zhu et al. (2022); Zhou et al. (2022) or assign different parts of the model's parameters to different tasks through masking Serra et al. (2018). Some methods Douillard et al. (2022); Zhai et al. (2023a; 2024) employ a combination of these techniques to achieve superior performance.

115 Approaches equipped with large-scale pretrained models (PTMs) have shifted the traditional CL 116 paradigms. By extracting robust features from PTMs, some works directly take advantage of this to 117 construct classifiers based on PTMs. SLCA Zhang et al. (2023) applied varying learning rates for 118 representation layer and classifier to mitigate progressive overfitting, and reduce classification bias through prototype replay. RanPAC McDonnell et al. (2024) employed a frozen random projection 119 layer to project pretrained features into a higher-dimensional space for better linear separability. 120 Some works applied mixture of experts (MoE) methods, LAE Gao et al. (2023) employed an on-121 line module for learning new tasks and an offline module for preserving learned knowledge, with 122 final predictions during the inference stage derived from the maximum logit of the two modules. 123 ESN Wang et al. (2023) trained a separate classifier for each new task and introduced an anchor-124 based energy self-normalization strategy, with a voting-based strategy during inference stage to unify 125 the classifiers. The recent proposed Yu et al. (2024) expanded the vision-language model Radford 126 et al. (2021) through MoE and designed a distribution discriminator to dynamically allocate test 127 samples to either MoE adapters or the original CLIP during inference.

128 Prompt-based CL Methods. Such approaches typically create a set of learnable prompts and pre-129 dict by incorporating the prompts most relevant to the current input. L2P Wang et al. (2022b) selects 130 prompts based on the cosine similarity between pretrained features and learnable keys, and integrates 131 the selected prompts into the token sequence after the image is patchified by ViT, before feeding it to 132 the transformer encoder. DualPrompt Wang et al. (2022a) categorizes prompts into general prompts, 133 which are shared by all samples, and expert prompts that are attached to the key and value following 134 L2P. CODA-Prompt Smith et al. (2023) proposed learning a set of prompt components to gener-135 ate attention scores for weighting the prompts. HiDe-Prompt Wang et al. (2024a) adds a learnable MLP to the frozen ViT and uses prototype replay for sequential training first-stage model to directly 136 predict class labels. It retrieves task identities by mapping them with class labels, allowing for the 137 explicit selection of prompts based on the predicted task identities to obtain instructed features. 138

139 140

141

3 PRELIMINARY

142 Prompt-based CL methods typically create a set of learnable prompts to generate instructed features. In the case of vision tasks, for input image $x \in \mathbb{R}^{H \times W \times C}$, a pretrained ViT f_{θ} first divides it into 143 144 N non-overlapping patches. It then attaches a class token and incorporated position encodings into 145 these patch embeddings to form a token sequence. This sequence is then processed through L-146 layers of stacked multi-head self-attention (MSA) blocks. We denote the input sequence of the *i*-th MSA layer as $x_e^i \in \mathbb{R}^{(N+1) \times D}$. For the *i*-th MSA layer, Query x_q^i , Key x_k^i , and Value x_v^i are first 147 148 created by multiplying the identical input sequence x_e^i with projection matrices W_Q^i , W_K^i , and W_V^i , 149 respectively. Then it calculates the self-attention scores and produces the output sequence through 150 the projection matrix W_O^i as:

151

$$x_{j} = Attention(x_{q}^{i}W_{Q}^{i,j}, x_{k}^{i}W_{K}^{i,j}, x_{v}^{i}W_{V}^{i,j}), j = 1, 2, ..., m$$
(1)

$$x_e^{i+1} = MSA(x_q^i, x_k^i, x_v^i) = Concat(x_1, ..., x_m)W_O^i$$
(2)

where $Attention(Q, K, V) = \frac{QK^T}{\sqrt{D}}V$, m is the number of attention heads, x_j is the output of the *j*th head.

Existing prompt-based CL methods generally adopt two techniques: prompt tuning Lester et al. (2021) and prefix tuning Li & Liang (2021). Prompt tuning appends learnable tokens to the sequence x_e^i before it proceeds to the next MSA layer, whereas prefix tuning involves appending them to the x_k^i and x_v^i sequences. 164 Split CIFAR-100 165 Split Aircrafts (In-distribution) (Out-of-distribution) 166 Task ID Final Label **Proportion**(%) Task ID Final Label | Proportion(%) 167 ViT-B-16 78.68 27.57 168 supervised pretrained X 1.89 1 X 12.90 1 169 X on ImageNet-21K X 17.49 1 14.28 1 170 X X X Х 5.15 42.04 171 Task ID **Final Label Proportion**(%) Task ID Final Label **Proportion(%)** 172 ViT-B-16-DINO 79.42 49.27 173 self-supervised pretrained Х 1 х 9.93 1 2.51 X on ImageNet-1K X 1 10.68 1 15.27 174 X X 7.39 Х Х 25.08 175

Table 1: Performance of supervised and self-supervised pretrained ViT models on Split CIFAR-100
 and Split Aircrafts datasets, both datasets were splitted into ten incremental tasks.

Table 2: Various combinations and effects between the first-stage model and second-stage model, experiments were conducted on Split Aircrafts, splitted into ten incremental tasks.

First-stage model	Second-stage model	Task accuracy	Final accuracy
ViT-B-16-IN21K	ViT-B-16-IN21K	40.47	45.06
ViT-B-14-DINOv2	ViT-B-16-IN21K	73.39	58.18

As illustrated in Figure 1, some methods Wang et al. (2022a;b); Smith et al. (2023) calculate the cosine similarity between uninstructed features and keys, implicitly treating the matching results as inferred task identity information (in gray). In contrast, a recent state-of-the-art method, HiDe-Prompt Wang et al. (2024b), trains a model and explicitly uses its prediction as task identity to select prompts (in beige). For prompt-based CL methods, an accurate task identity helps integrate relevant information for correct predictions, while an incorrect identity typically leads to misclassification by incorporating wrong information. To address the limitations of current approaches, we propose a novel framework (in brown) that incorporates self-distillation loss and a threshold filtering strategy.

190 191

176

4 EMPIRICAL ANALYSIS

192 193

We will investigate two key questions regarding prompt-based CL methods: 1. *How do the super*vised pretrained and self-supervised models perform under ID and OOD scenarios? 2. Does the first-stage model, which infers task identity, need to be identical to the second-stage model, which incorporates prompts and provides the final predictions? Our experiments are conducted using the HiDe-Prompt method Wang et al. (2024a), which has shown impressive performance across various benchmarks. We will present our observations based on these findings.

200 Self-supervised models excel under OOD scenarios. To demonstrate the performance of prompt-201 based CL methods across different distribution scenarios, we conducted a comparison of two ViT-B-16 models: one supervised pretrained on ImageNet-21K and the other self-supervised pretrained 202 on ImageNet-1K using DINO Caron et al. (2021), across in-distribution dataset Split CIFAR-100 203 and out-of-distribution dataset Split Aircrafts. The results shown in Table 1 indicate that both mod-204 els effectively utilized task identity under ID scenario, rarely making errors when the task identity 205 is correct, with incorrect predictions on only 1.89% and 2.51% of test samples under this condi-206 tion. The proportion of correctly predicted samples was as high as 78.68% and 79.42%. Under 207 OOD scenarios, these two models exhibit significantly different performance. Specifically, the self-208 supervised model was much better at leveraging the task identity, with 21.7% higher than supervised 209 pretrained model when both the task identity and the final prediction were correct. 210

Observation 1: Under ID scenarios, both self-supervised and supervised pretrained models perform similarly, rarely making errors. However, under OOD scenarios, self-supervised models significantly outperform supervised ones in leveraging task identity.

Task identity provision is model-agnostic. Previous methods typically restricted the first-stage model and second-stage model to use the same architecture and pretrained weights. The accuracy of the task identity provided by the first-stage model is critical for the subsequent task-specific

prompt selection by the second-stage model. If the first-stage task predictions are unreliable, it will
inevitably affect the final performance outputed from the second-stage model. The results presented
in Table 2 for the Aircraft dataset indicate that when using the same model, ViT-B-16-IN21K, for
both the first and second stages, the final accuracy reaches 45.06%. we then replaced the first-stage
model with ViT-B-14-DINOv2, which significantly improved task accuracy by 32.92%, leading to
an overall gain of 13.12% in final accuracy.

Observation 2: Providing task identity is a model-agnostic behavior, a reliable task identity benefits the second stage without requiring identical architecture for both stages.

5 Method

5.1 BOOSTING CLASS INCREMENTAL LEARNING VIA SELF-SUPERVISED LEARNING

Our empirical analysis in Section 4 suggests that the first-stage and second-stage models can be trained using different architectures. We propose incorporating self-supervised learning into our framework from two angles: utilizing self-supervised pretrained models and applying selfsupervised loss to enhance task identity prediction.

232 233 234

222

223

224 225

226 227

228 229

230

231

5.1.1 WITH SELF-SUPERVISED PRETRAINED MODEL

235 Supervised pretrain often leads to neural collapse Galanti et al. (2021); Papyan et al. (2020); Fang 236 et al. (2021); Zhai et al. (2023b), where features of the same class cluster around their mean, making 237 it hard to generalize to OOD scenarios where more fine-grained discriminative ability is required. 238 Self-supervised pretrain avoids this issue by learning general representations capable of generalizing 239 to novel or unseen data, thereby enhancing OOD performance, which is especially advantageous for 240 CL. Considering the high transferability of features in self-supervised pretrained models and their 241 inherent advantages under OOD scenarios, using them as the first-stage models to provide task 242 identity proves more effective. Specifically, in this work, we explore models that have been self-243 supervised pretrained at different scales with various pretrain methods, including iBOT Zhou et al. (2021) pretrained on ImageNet-21K/1K, DINO Caron et al. (2021) and MoCo v3 Chen et al. (2021) 244 pre-trained on ImageNet-1K, and DINOv2 Oquab et al. (2023) pre-trained on LVD-142M. 245

246 247

260

261 262 263

5.1.2 WITH SELF-SUPERVISED LOSS

248 We incorporated self-distillation loss into our approach, a self-supervised learning objective derived 249 from the DINO (self-distillation with no labels) framework Caron et al. (2021), which is effective 250 in enhancing feature extraction capabilities of vision transformers. During training the first-stage 251 model, we consider the current network as f_{θ_s} and replicate it as the teacher network, denoted as 252 f_{θ_t} , parameterized by θ_s and θ_t respectively. Initially, we generate a set of augmented views V from 253 the input image x, including two global views $\{x_1^g, x_2^g\}$ and several local views. Then all views are fed into the student network, while the global views are input only into the teacher network. Cross-254 entropy loss is minimized between the outputs of the two networks to match their distributions: 255

$$\min_{\theta_s} \sum_{x \in \{x_g^1, x_g^2\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(P_t(x), P_s(x'))$$
(3)

Where $H(a,b) = -a \log b$, P(x) represents the network's probability distribution over K dimensions, calculated by the following equation:

$$P(x)^{(i)} = \frac{\exp(f_{\theta}(x)^{(i)}/\tau)}{\sum_{k=1}^{K} \exp(f_{\theta}(x)^{(k)}/\tau)}$$
(4)

with f_{θ_s} , τ_s for student network, and f_{θ_t} , τ_t for teacher network. The parameters of the student network θ_s are optimized via stochastic gradient descent during the minimization of 3, while the parameters of the teacher network θ_t are updated through the EMA (Exponential Moving Average) algorithm Morales-Brotons et al. (2024).

Applying self-distillation loss (SDL) for CL on downstream data streams offers several advantages. Firstly, the computation of self-distillation loss effectively acts like traditional knowledge 270 distillation-based regularization methods for preventing forgetting, but it relaxes constraints to main-271 tain model plasticity. By promoting consistency between multiple views (global and local) of the 272 same input, it facilitates the extraction of generalized features from images. The learned features are 273 task-invariant, making these generalized features favorable for the model's improved understanding 274 and adaptation to new and unseen data, which is crucial for CL. Secondly, by updating the teacher network with EMA, which is more stable and updates slowly, old knowledge is preserved. By con-275 tinuously distilling knowledge from the teacher to the student and utilizing different views of the 276 data, SDL can potentially mitigate catastrophic forgetting. 277

278

297

303

308

309 310 311

312 313

314

318

280 Our analysis in Section 4 reveals a key issue with the existing prompt-based CL methods: even if 281 the first-stage model correctly predicts the class label, the second-stage model may still misclassify 282 the category, despite receiving the correct task identity. This situation occurs frequently in OOD 283 scenarios, with an occurrence rate of around 10%. To mitigate this issue, we propose a simple yet 284 effective threshold filtering strategy to determine which samples require two-stage inference and 285 which do not. We utilize β distribution for modeling confidence scores and choosing the lowest 286 boundary of a specific highest density interval as the decision threshold. Predictions from the first-287 stage model are used directly for samples above this threshold, while samples below this threshold are fed to the second-stage model. Specially, given an input image x, the output probability for class 288 *i* of this image is $P_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$, where z_i is its logits, *j* is the number of classes encountered so 289 290 far. The predicted class $\hat{y} = \arg \max P_i$ is the class with the highest probability. The confidence 291 $c = max P_i$ that the model assigns to x is the probability of the predicted class. 292

Assume that the confidence score follows β distribution, then the prior distribution can be represented as: $\theta \sim Beta(\alpha_0, \beta_0)$, where α_0 and β_0 are two parameters of the β distribution, which control the shape of the distribution.

For *n* observed confidence scores from all test samples so far $c_1, c_2, ..., c_n$, assuming they are also sampled from β distribution, the likelihood function is:

$$L(\theta) = \prod_{i=1}^{n} c_i^{\alpha - 1} \cdot (1 - c_i)^{\beta - 1}$$
(5)

According to Bayes theorem, combining the prior distribution and the likelihood function, the posterior distribution is:

$$P(\theta|c_1,\ldots,c_n) \propto \left(\prod_{i=1}^n c_i^{\alpha-1} \cdot (1-c_i)^{\beta-1}\right) \cdot \operatorname{Beta}(\alpha_0,\beta_0)$$
(6)

Since both the prior and likelihood functions are β distributions, the posterior distribution is still β distribution, and the updated parameters are:

$$\alpha_{post} = \alpha_0 + \sum_{i=1}^{n} c_i, \beta_{post} = \beta_0 + n - \sum_{i=1}^{n} c_i \tag{7}$$

$$\theta_{post} \sim Beta(\alpha_{post}, \beta_{post})$$
 (8)

After obtaining the posterior distribution, the confidence threshold can be determined by calculating the quantile of the posterior distribution. For example, the lower limit of the d% highest density interval (HDI) is selected as the confidence threshold for classification:

$$\tau = Beta^{-1}\left(d|\alpha_{post},\beta_{post}\right) \tag{9}$$

where $Beta^{-1}$ represents the quantile function (or inverse cumulative distribution function) of the β distribution.

For samples that exceed the threshold τ , we consider that the predictions from the first-stage model are reliable and adopt them. For samples below τ , we fed them to the second-stage model for a secondary inference phase, guided by the inferred task identity.

			Split Aircrafts			Split Cars-196		
PTM (A-B)	Method	FAA (†)	CAA (†)	FFM (\downarrow)	FAA (†)	CAA (†)	FFM (\downarrow)	
	L2P	22.76 ± 0.66	35.99 ± 1.18	21.22 ± 3.79	34.49 ± 0.19	45.97 ± 1.40	12.39 ±2.18	
	DualPrompt	23.82 ± 1.76	35.88 ± 1.85	16.76 ± 1.87	43.21 ± 0.50	$51.72 \pm\!\! 1.46$	12.00 ± 1.76	
Sup-21K*	CODA-Prompt	19.02 ± 0.99	35.69 ± 1.14	34.04 ± 4.26	33.12 ± 0.13	50.16 ± 1.05	34.54 ± 2.19	
	HiDe-Prompt	$\underline{44.21} \pm 0.79$	$\underline{52.75 \pm 1.78}$	10.04 ± 0.54	$\underline{49.75 \pm 0.18}$	$\underline{58.47} \pm 0.87$	$\underline{7.80\pm0.22}$	
	Ours	$\textbf{47.47} \pm \textbf{0.34}$	$\textbf{55.45} \pm \textbf{1.75}$	$\textbf{7.77} \pm \textbf{0.44}$	54.90 ± 0.10	62.23 ± 1.15	$\textbf{7.59} \pm \textbf{0.39}$	
	L2P	33.26 ± 2.82	50.62 ± 2.37	15.74 ±2.28	49.42 ± 0.86	61.02 ± 1.08	10.29 ±0.46	
	DualPrompt	30.77 ± 2.59	46.28 ± 2.25	$23.21 \pm \!$	$46.91 \pm\! 0.80$	57.92 ± 1.17	13.25 ± 0.61	
iBOT-21K	CODA-Prompt	$36.52 \pm\! 1.04$	$53.43 \ {\pm}1.44$	$21.05 \ {\pm} 2.72$	60.28 ± 0.43	70.10 ± 0.18	11.89 ± 0.43	
	HiDe-Prompt	$\underline{60.24 \pm 1.52}$	$\underline{63.44}\pm\!3.70$	$\textbf{4.45} \pm \textbf{0.16}$	$\underline{68.23 \pm 0.50}$	$\underline{71.67 \pm 1.37}$	$\textbf{3.07} \pm \textbf{0.14}$	
	Ours	$\textbf{66.58} \pm \textbf{0.77}$	71.36 ± 1.59	$\underline{5.15 \pm 0.52}$	71.51 ± 0.30	$\textbf{75.15} \pm \textbf{0.80}$	$\underline{5.03\pm}0.13$	
	L2P	34.82 ± 2.33	51.06 ±1.87	17.55 ±1.18	52.40 ± 0.53	62.90 ± 1.20	10.33 ±0.99	
	DualPrompt	34.25 ± 1.79	49.27 ± 2.37	21.67 ± 2.70	53.10 ± 0.62	$64.75 \ {\pm} 0.84$	14.88 ± 0.54	
iBOT-1K	CODA-Prompt	39.24 ± 0.60	56.10 ± 1.65	18.56 ± 2.32	62.11 ± 1.09	72.48 ± 0.77	10.74 ± 0.86	
	HiDe-Prompt	$\underline{61.37 \pm 1.68}$	64.14 ± 3.57	$\textbf{4.59} \pm \textbf{0.54}$	70.67 ± 0.11	$\underline{73.91 \pm 1.16}$	$\textbf{3.95} \pm \textbf{0.25}$	
	Ours	65.54 ± 0.52	70.34 ± 1.55	$\underline{5.82 \pm 0.22}$	$\textbf{74.94} \pm \textbf{0.04}$	$\textbf{78.58} \pm \textbf{0.76}$	$\underline{4.87 \pm 0.52}$	
	L2P	36.52 ± 0.97	50.08 ±1.42	17.85 ± 3.34	51.68 ± 0.69	62.36 ± 0.63	11.02 ± 1.97	
	DualPrompt	36.05 ± 2.46	$52.15\ {\pm}2.81$	18.46 ± 2.68	52.97 ± 0.83	$64.21 \pm\!\! 1.25$	11.47 ± 1.32	
DINO-1K	CODA-Prompt	42.71 ± 1.42	$57.92 \ {\pm}0.82$	17.66 ± 0.93	$62.43 \ {\pm}0.74$	$72.13 \ {\pm} 0.51$	10.29 ± 0.24	
	HiDe-Prompt	$\underline{62.14} \pm 1.44$	$\underline{65.55 \pm 3.55}$	$\textbf{4.49} \pm \textbf{0.49}$	$\underline{71.93} \pm 0.08$	$\underline{75.57 \pm 0.91}$	$\textbf{3.82} \pm \textbf{0.22}$	
	Ours	66.12 ± 0.63	$\textbf{71.11} \pm \textbf{1.81}$	$\underline{5.88\pm}0.34$	$\textbf{75.39} \pm \textbf{0.11}$	$\textbf{78.95} \pm \textbf{0.79}$	$\underline{4.61 \pm 0.17}$	
	L2P	26.84 ± 1.07	43.69 ±2.25	9.08 ± 1.47	39.10 ± 0.39	54.22 ± 0.81	3.22 ±0.15	
	DualPrompt	27.68 ± 2.01	43.55 ± 2.03	9.58 ± 2.90	41.73 ± 0.80	56.68 ± 1.00	3.29 ± 0.38	
MoCo-1K	CODA-Prompt	35.75 ± 1.56	52.11 ± 1.90	20.12 ± 2.78	54.06 ± 0.21	65.25 ± 1.17	16.57 ±0.51	
	HiDe-Prompt	$\underline{53.05\pm}0.89$	58.74 ± 2.79	$\textbf{4.89} \pm \textbf{0.20}$	$\underline{66.09 \pm 0.36}$	$\underline{70.75\pm1.30}$	3.87 ± 0.11	
	Ours	57.52 ± 0.52	63.75 ± 1.80	7.12 ± 0.29	$\textbf{68.64} \pm \textbf{0.14}$	$\textbf{72.92} \pm \textbf{1.11}$	5.03 ± 0.45	

Table 3: Performance comparison of various methods on Split Aircrafts and Split Cars-196 Datasets,
 we present FAA, CAA, and FFM, each metric with mean and standard deviation over three different
 random seeds. The best outcome is marked in bold, with the second-best underlined. All experi mental results for compared methods were reproduced by ourselves.

6 EXPERIMENTS

6.1 EXPERIMENTAL SETUP

Datasets. For evaluating the CIL performance of different methods under OOD scenarios, we chose Split Aircrafts Maji et al. (2013) and Split Cars-196 Krause et al. (2013) datasets. Additionally, we also conducted experiments on CIFAR-100 Krizhevsky et al. (2009), ImageNet-R Hendrycks et al. (2021), and CUB-200 Wah et al. (2011); the first two are subclasses of ImageNet, whereas CUB-200 has 52 overlapping categories as ImageNet Ostapenko et al. (2022); Wen et al. (2022). Cars-196 was divided into 7 tasks, all other datasets were divided into 10 tasks each. Detailed descriptions of these datasets are available in the appendix.

Implementation Details. Following Wang et al. (2024a), we used ViT-Base model with a patch size of 16 (except for DINOv2 Oquab et al. (2023), which is 14), set the projection dimension to 2048 for computing SDL, and used d=95 to determine the threshold. The weight of SDL was set to 0.1, with all other training hyperparameters consistent with those in Wang et al. (2024a). Following the evaluation metrics in Wang et al. (2024a), we reported the final average accuracy (FAA) of all seen classes, the average accuracy over tasks refereed as cumulative average accuracy (CAA), and final forgetting measure (FFM) of all previous tasks. The models are either supervised pretrained Ridnik et al. (2021) or self-supervised pretrained Zhou et al. (2021); Caron et al. (2021); Chen et al. (2021) on ImageNet-21K/1K, aligning with the benchmarks for a fair comparison. In extended experiments, we also utilized DINOv2 Oquab et al. (2023) pretrained on LVD-142M.

379													
380	Mathad		Split	t ImageNe	nageNet-R			Split CUB-200					
381	Method	Sup-21K	iBOT-21K i	BOT-1K	DINO-1K	MoCo-1	K Sup	o-21K	iBOT-21	K iBOT	-1K	DINO-1K	MoCo-1K
501	L2P	59.61	61.23	64.16	61.01	55.01	73	3.92	52.57	57.5	55	53.79	52.94
382	DualPrompt	65.83	60.50	64.81	60.75	54.06	74	4.55	48.09	56.6	52	55.75	52.04
383	CODA-Prompt	59.93	66.72	69.00	63.60	61.50	72	2.25	54.73	59.9	96	60.01	45.79
303	HiDe-Prompt	73.26	74.73	76.82	73.90	67.58	83	3.26	70.73	77.2	24	76.75	71.82
384	Ours	72.20	75.07	76.29	73.58	68.55	83	3.59	71.92	77.9	94	76.53	72.59
385													
386		-				Split	CIFAR	-100			-		
387		-	Method	Sup-2	21K iBOT	-21K iE	OT-1K	DIN	IO-1K	MoCo-1K	_		
388			L2P	83.4	43 72	.09	75.59	79	9.79	77.32			
000			DualPrompt	87.9	98 73	.65	77.49	78	8.10	73.80			
389			CODA-Promp	ot 81.	30 77	.25	78.79	8	1.48	79.56			
390			HiDe-Prompt	92.9	96 <u>91</u>	.92	92.04	90	0.10	90.57			
			Ours	92.3	35 92	.49	92.38	9	0.86	90.30			

Table 4: FAA(\uparrow) comparison for various methods on three different datasets.

393

378

6.2 COMPARISON TO STATE-OF-THE-ART METHODS

394 **Results on OOD scenarios.** Table 3 displays the overall results of our method compared to four popular prompt-based CL methods as we mentioned above on two OOD datastets: the Split Aircrafts and Split Cars-196. For all approaches², self-supervised pretrained models (without *) demonstrate 396 markedly better performance than supervised pretrained models (with *), consistent with our ob-397 servations that self-supervised pretrained models have significant advantages under OOD scenarios. 398 Among all the methods, the proposed method demonstrated significant improvement over all four 399 compared methods across all models For supervised pretrained models, we achieved an FAA in-400 crease of 3.26% and a CAA increase of 2.70% on Split Aircrafts, with improvements of 5.15% and 401 3.76% respectively on Split Cars-196, while also achieving the lowest forgetting on both datasets. 402 Similar improvements were observed for self-supervised models, notably a 6.34% FAA and 7.92% 403 CAA boost for the iBOT model pretrained on ImageNet-21K. The improvements in CAA highlight 404 our method's consistent performance improvements across all stages of incremental learning, rather 405 than only after all tasks have been trained. Additionally, our method exhibits the lowest standard 406 deviation across almost all metrics, further indicating the stability of our approach. Compared to 407 HiDe-Prompt, our method exhibited slightly higher forgetting on self-supervised models.

408 **Results on non-OOD scenarios.** We also evaluated our proposed method on Split ImageNet-R, 409 CIFAR-100, and CUB-200, which are commonly assessed in previous methods. The overall results 410 are shown in Table 4. We achieved improvements of 0.55% and 0.16% over the SOTA method 411 across five different models on two ID datasets Split CUB-200 and Split CIFAR-100, respectively. 412 For Split ImageNet-R, we observed certain improvements on two self-supervised pretrained models iBOT-21K and MoCo-1K, while experiencing slight declines on other models. In summary, although 413 existing methods already perform excellently under ID scenarios, our method can still achieve cer-414 tain improvements. However, as we previously mentioned, our method shows even greater advan-415 tages under OOD scenarios. 416

417

419

6.3 FURTHER ANALYSIS 418

Impact of pretrain data scale. In analyzing the effect of the scale of pretrain data, results in 420 Table 5 present the final average accuracy (FAA) provided by the first-stage model and second-421 stage model across different data scales. We found that, in the first stage, supervised pretrained 422 model significantly underperforming all self-supervised models on Split CIFAR-100 and falling far 423 behind on Split Aircrafts, its strength lies in the second stage under ID scenarios (Split CIFAR-100), 424 where it achieves a 12.94% improvement over the first stage-the highest among all models. For 425 self-supervised pretrained models, except for MoCo-1K on Split Aircrafts, other models pretrained 426 on ImageNet-1K exhibit comparable performance in both stages. With increased scale of pretrain 427 data, the first-stage performance improves notably, while the improvement from the second stage 428 over the first stage lessens. For DINOv2 with the largest scale of pretrain data, the second stage's

⁴³⁰ ²We denote the supervised pretrained model as *, where 'Sup' stands for 'supervised' and '21K' refers to the ImageNet-21K dataset. The other four models are self-supervised, labeled in the format 'A-B', with A 431 representing the pretrain method and B the dataset.

	Pretrain	Pretrain	Split C	IFAR-100	Split Aircrafts	
	method	data	first-stage FAA	second-stage FAA	first-stage FAA	second-stage FAA
Supervised	-	ImageNet-21K	79.41	92.35	42.42	47.37
Self-supervised	MoCoV3 DINO iBOT	ImageNet-1K	81.33 82.28 83.39	90.30 90.86 92.38	53.23 61.12 60.72	58.18 67.36 66.49
-	iBOT DINOv2	ImageNet-21K LVD142M	85.18 88.67	92.49 92.51	63.37 73.39	68.08 75.46

Table 5: Performance comparison across models on Split CIFAR-100 and Split Aircrafts datasets.

Table 6: The performance on Split CIFAR-100 and Aircrafts of different methods employing diverse combinations of the first-stage model and second-stage model, compared to our proposed method.

Dataset	First-stage model	Second-stage model	L2P	Dual	HiDe	Ours
	Sup-21K		83.43	87.98	92.96	92.35
Split CIFAR-100	DINO-1K	Sup-21K	81.83	87.77	93.01	92.53
-	DINOv2-LVD142M	-	82.15	88.09	<u>93.88</u>	94.18
	Sup-21K		24.07	24.75	45.06	47.37
Split Aircrafts	DINO-1K	Sup-21K	25.26	28.30	<u>53.29</u>	61.06
	DINOv2-LVD142M		25.56	29.43	<u>58.18</u>	70.01

452 453

444

432

improvement over the first on Split CIFAR-100 is only 3.84%, and the relative improvement on Split
 Aircrafts is merely 2.07% (for comparison, although not shown, it is only 0.33% on HiDe-Prompt).

456 **Combination of different models.** We explored different combinations of first-stage and second-457 stage models on Split CIFAR-100 and Split Aircrafts datasets, with results shown in Table 6. Ex-458 isting methods that often use identical architectures and weights for both models are outlined in 459 the first row. For Split Aircrafts, by replacing the first-stage model with self-supervised pretrained 460 models (DINO-1K and DINOv2-LVD142M), we observed significant improvements both in explicit task identity prediction methods (HiDe-Prompt and ours) and in implicit task identity usage 461 methods like L2P and DualPrompt, which rely on matching results between pretrained features and 462 learnable keys. This indicates that supervised pretrained models might lack sufficient capability 463 to encode more fine-grained features under OOD scenarios, highlighting the necessity of employ-464 ing self-supervised pretrained models. For our method, we achieved substantial improvements by 465 enhancing the accuracy of task identity and filtering first-stage predictions. On Split CIFAR-100, 466 previous methods that implicitly infer task identity experienced a slight decrease when the first-stage 467 model was replaced with a self-supervised model, except for a modest increase when using DINOv2 468 with DualPrompt. However, explicit methods continue to show significant enhancements, particu-469 larly when applying DINOv2 to our method. Combining the high first-stage accuracy of DINOv2 470 with the extremely high efficacy of the supervised pretrained model in the second stage, we achieved 471 an impressive accuracy of 94.18%. We show more results from different model combinations in Ta-472 ble 9, further analysis can be found in the appendix.

473
474
474
475
475
476
476
478
479
479
479
470
470
470
470
471
471
472
473
474
475
475
476
476
476
476
476
477
478
478
479
479
479
470
470
470
470
470
471
471
472
473
474
475
474
475
476
476
476
476
477
478
478
478
478
479
479
479
470
470
470
470
470
471
471
471
472
473
474
475
476
476
476
476
476
476
476
476
477
478
478
478
479
479
470
470
470
470
470
471
471
471
472
473
474
474
475
476
476
476
476
476
476
476
476
476
477
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478
478

477 478

479

 $S = \frac{1}{T} \sum_{i=1}^{c} \sum_{j=i+1}^{c} \cos(p_i, p_j)$ (10)

Where T is the number of possible pairs of classes, cos is the cosine similarity function, p_i is the mean vector of extracted features belong to class i. On one hand, S reflects the difficulty of classification on that dataset; the lower the similarity, the more orthogonal the features, and consequently, the easier the classification. On the other hand, this metric can also be used to indicate the extent to which the pretrained model leaks information about downstream data.

The heatmaps of the similarity matrices between classes are shown in Figure 3. The specific values for S are in Table 8, which can be found in the appendix. The results indicate that class similarity S



Figure 3: Class similarity heatmaps for different models on CIFAR-100 (top row) and Aircrafts (bottom row) datasets. Each heatmap represents the pairwise class similarity matrix, with the color scale indicating similarity values. The lower the feature similarity indicates more orthogonal features and lower classification difficulty. The higher suggests that the downstream dataset is more out-ofdistribution compared to the pretrained data.

Table 7: Ablation studies of our proposed method, we present $FAA(\uparrow)$ for comparison.

Mathad	Split Aircrafts					Split Cars-196					
Method	Sup-21K	iBOT-21K	iBOT-1K	DINO-1K	MoCo-1K	Sup-21K	iBOT-21K	iBOT-1K	DINO-1K	MoCo-1K	
HiDe-Prompt	45.06	63.19	64.72	64.99	54.76	49.47	68.86	70.80	71.86	66.56	
+SDL	46.09	64.39	64.51	66.82	56.41	51.47	70.07	73.85	74.43	67.74	
+threshold	47.37	68.08	66.49	67.36	58.18	54.72	72.08	74.88	75.49	68.83	

on CIFAR-100 is relatively low across all models, suggesting more orthogonal features and easier classification. This is particularly obvious in Sup-21K, which shows a similarity as low as 0.228. As all categories of CIFAR-100 are included in ImageNet Wen et al. (2022), the low similarity observed in supervised pretrained models on this dataset is expected. On the contrast, these models exhibit relatively higher class similarity on Aircrafts, despite the category "*aircraft*" being included in ImageNet. Nonetheless, accurately classifying specific aircraft types within this fine-grained dataset remains challenging for the models.

521

522

523

524

499

500

501

502

509 510

511

512

513

514

515

6.4 ABLATION STUDY

520 We conducted ablation studies to validate the effectiveness of each component of our proposed method, with the results shown in Table 7. Both modules consistently improved performance across different models. Specifically, on Split Aircrafts, SDL improved performance in all models except for a slight decrease in iBOT-1K, while our proposed threshold filtering strategy provided further enhancements. On Split Cars-196, both modules showed stable improvements across all models.

525 526

7 **CONCLUSIONS**

527 528

529 In this work, we revisited existed prompt-based CL methods through comprehensive analysis. 530 Through empirical analysis, we revealed several limitations of existing methods. To overcome these 531 limitations, we significantly boosted the performance of various pretrained models under OOD scenarios by introducing self-supervised learning objectives in the first stage and proposing a simple 532 threshold filtering strategy. Moreover, we explored the efficiency of self-supervised pretrained mod-533 els in providing task identities, thereby achieving further improvements and contributing to estab-534 lishing a more unified framework for these approaches. 535

536 **Limitations.** Although our proposed method shows significant improvement under OOD scenarios, the gains under ID scenarios are relatively less. Additionally, even though we have overcome the limitations of previous methods-allowing the first-stage model in our framework to be any archi-538 tecture (ViT, CNN, etc.), the second-stage model still relies on the ViT structure to integrate prompts and generate instructed features for prediction.

540 REFERENCES 541

542 543 544	Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 139–154, 2018.
545 546 547	Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 9650–9660, 2021.
548 549 550 551	Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 532–547, 2018.
552 553	Haoran Chen, Zuxuan Wu, Xintong Han, Menglin Jia, and Yu-Gang Jiang. Promptfusion: Decoupling stability and plasticity for continual learning. <i>arXiv preprint arXiv:2303.07223</i> , 2023.
554 555 556 557	Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 9640–9649, 2021.
558 559 560	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. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 44(7):3366–3385, 2021.
561 562 563	Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. <i>arXiv preprint arXiv:2010.11929</i> , 2020.
564 565 566	Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9285–9295, 2022.
567 568 569 570	Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer- peeled model: Minority collapse in imbalanced training. <i>Proceedings of the National Academy</i> <i>of Sciences</i> , 118(43):e2103091118, 2021.
571 572	Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learn- ing. <i>arXiv preprint arXiv:2112.15121</i> , 2021.
573 574 575 576	Qiankun Gao, Chen Zhao, Yifan Sun, Teng Xi, Gang Zhang, Bernard Ghanem, and Jian Zhang. A unified continual learning framework with general parameter-efficient tuning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 11483–11493, 2023.
577 578 579	Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empiri- cal investigation of catastrophic forgetting in gradient-based neural networks. <i>arXiv preprint</i> <i>arXiv:1312.6211</i> , 2013.
580 581 582 583	Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. <i>ICCV</i> , 2021.
585 586	Geoffrey Hinton. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> , 2015.
587 588 589 590	Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 831–839, 2019.
591 592 593	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. <i>Proceedings of the national academy of sciences</i> , 114(13):3521–3526, 2017.

594	Alexander Kolesnikov, Lucas Bever, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly,
595	and Neil Houlsby. Big transfer (bit): General visual representation learning. In Computer Vision-
596	ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part
597	V 16, pp. 491–507. Springer, 2020.

- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision work-shops*, pp. 554–561, 2013.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 2009.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2021.
- Kiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
 Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers),
 pp. 4582–4597, 2021.
- ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁴
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
- Kialei Liu, Yu-Song Hu, Xu-Sheng Cao, Andrew D Bagdanov, Ke Li, and Ming-Ming Cheng. Long-tailed class incremental learning. In *European Conference on Computer Vision*, pp. 495–512.
 Springer, 2022.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost
 Van De Weijer. Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533, 2022.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.

- Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton van den Hengel.
 Ranpac: Random projections and pre-trained models for continual learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Daniel Morales-Brotons, Thijs Vogels, and Hadrien Hendrikx. Exponential moving average of
 weights in deep learning: Dynamics and benefits. *Transactions on Machine Learning Research*,
 2024.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Oleksiy Ostapenko, Mihai Puscas, Tassilo Klein, Patrick Jahnichen, and Moin Nabi. Learning to
 remember: A synaptic plasticity driven framework for continual learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11321–11329, 2019.
- Oleksiy Ostapenko, Timothee Lesort, Pau Rodriguez, Md Rifat Arefin, Arthur Douillard, Irina Rish, and Laurent Charlin. Continual learning with foundation models: An empirical study of latent replay. In *Conference on lifelong learning agents*, pp. 60–91. PMLR, 2022.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40): 24652–24663, 2020.

657

659

660

661 662

685

- 648 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 649 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 650 models from natural language supervision. In International conference on machine learning, pp. 651 8748-8763. PMLR, 2021.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: 653 Incremental classifier and representation learning. In Proceedings of the IEEE conference on 654 Computer Vision and Pattern Recognition, pp. 2001–2010, 2017. 655
- 656 Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. arXiv preprint arXiv:2104.10972, 2021. 658
 - Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International conference on machine learning*, pp. 4548-4557. PMLR, 2018.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative 663 replay. Advances in neural information processing systems, 30, 2017. 664
- 665 James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, 666 Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual de-667 composed attention-based prompting for rehearsal-free continual learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11909–11919, 2023. 668
- 669 Zhe Tao, Lu Yu, Hantao Yao, Shucheng Huang, and Changsheng Xu. Class incremental learning 670 for light-weighted networks. *IEEE Transactions on Circuits and Systems for Video Technology*, 671 2024. 672
- 673 Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. arXiv preprint 674 arXiv:1904.07734, 2019. 675
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd 676 birds-200-2011 dataset. 2011. 677
- 678 Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical de-679 composition of prompt-based continual learning: Rethinking obscured sub-optimality. Advances 680 in Neural Information Processing Systems, 36, 2024a. 681
- 682 Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual 683 learning: theory, method and application. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2024b. 684
- Yabin Wang, Zhiheng Ma, Zhiwu Huang, Yaowei Wang, Zhou Su, and Xiaopeng Hong. Isolation 686 and impartial aggregation: A paradigm of incremental learning without interference. In Proceed-687 ings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 10209–10217, 2023. 688
- 689 Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, 690 Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for 691 rehearsal-free continual learning. In European Conference on Computer Vision, pp. 631-648. 692 Springer, 2022a.
- 693 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vin-694 cent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In Pro-695 ceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 139–149, 696 2022b. 697
- Shixian Wen, Amanda Sofie Rios, Kiran Lekkala, and Laurent Itti. What can we learn from mis-699 classified imagenet images? arXiv preprint arXiv:2201.08098, 2022.
- I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-701 supervised learning for image classification. arXiv preprint arXiv:1905.00546, 2019.

702 Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Ping Hu, Dong Wang, Huchuan Lu, and You He. Boosting 703 continual learning of vision-language models via mixture-of-experts adapters. In Proceedings of 704 the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 23219–23230, 2024. 705 Lu Yu, Vacit Oguz Yazici, Xialei Liu, Joost van de Weijer, Yongmei Cheng, and Arnau Ramisa. 706 Learning metrics from teachers: Compact networks for image embedding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2907–2916, 2019. 708 709 Lu Yu, Bartlomiej Twardowski, Xialei Liu, Luis Herranz, Kai Wang, Yongmei Cheng, Shangling 710 Jui, and Joost van de Weijer. Semantic drift compensation for class-incremental learning. In 711 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 6982– 6991, 2020. 712 713 Lu Yu, Xialei Liu, and Joost Van de Weijer. Self-training for class-incremental semantic segmenta-714 tion. IEEE Transactions on Neural Networks and Learning Systems, 34(11):9116–9127, 2022. 715 716 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In International conference on machine learning, pp. 3987–3995. PMLR, 2017. 717 718 Jiang-Tian Zhai, Xialei Liu, Andrew D Bagdanov, Ke Li, and Ming-Ming Cheng. Masked autoen-719 coders are efficient class incremental learners. In Proceedings of the IEEE/CVF International 720 Conference on Computer Vision, pp. 19104–19113, 2023a. 721 Jiang-Tian Zhai, Xialei Liu, Lu Yu, and Ming-Ming Cheng. Fine-grained knowledge selection and 722 restoration for non-exemplar class incremental learning. In Proceedings of the AAAI Conference 723 on Artificial Intelligence, volume 38, pp. 6971-6978, 2024. 724 725 Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. In-726 vestigating the catastrophic forgetting in multimodal large language models. arXiv preprint 727 arXiv:2309.10313, 2023b. 728 Gengwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learner 729 with classifier alignment for continual learning on a pre-trained model. In Proceedings of the 730 IEEE/CVF International Conference on Computer Vision, pp. 19148–19158, 2023. 731 732 Da-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards 733 memory-efficient class-incremental learning. arXiv preprint arXiv:2205.13218, 2022. 734 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. Image 735 bert pre-training with online tokenizer. In International Conference on Learning Representations, 736 2021. 737 Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation 738 739 and self-supervision for incremental learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5871–5880, 2021. 740 741 Kai Zhu, Wei Zhai, Yang Cao, Jiebo Luo, and Zheng-Jun Zha. Self-sustaining representation expan-742 sion for non-exemplar class-incremental learning. In Proceedings of the IEEE/CVF Conference 743 on Computer Vision and Pattern Recognition, pp. 9296–9305, 2022. 744 745 746 747 748 749 750 751 752 754 755

756 A APPENDIX

758 A.1 CL DATASETS 759

Aircrafts includes 100 different categories, each representing a specific type of aircraft. It contains
 a total of 6667 training samples and 3333 testing samples.

Cars-196 comprises 196 distinct categories, with each category corresponding to a specific type of
 car. The dataset consists of 8144 training images and 8041 test images.

CIFAR-100 comprises 100 categories, divided into 20 superclasses, each containing 5 fine-grained categories. Each category consists of 600 images of size 32x32, with 500 designated for the training set and 100 for the testing set.

ImageNet-R consists of 200 subclasses extracted from ImageNet, each containing challenging or stylistically recollected samples. The dataset includes a total of 30000 samples without a standard-ized split between training and testing sets. Typically, previous methods Smith et al. (2023); Wang et al. (2024a) assign 80% of the samples (24000) for training and the remaining 20% (6000) for testing.

CUB-200-2011 comprises 200 distinct categories of birds, featuring 5994 training samples and 5794 testing samples.

A.2 CLASS SIMILARITY OF DIFFERENT MODELS

CIFAR-100

Aircrafts

Table 8: Average class similarity for different models on testsets of CIFAR-100 and Aircrafts.

iBOT-1K

0.649

0.923

DINO-1K

0.655

0.923

MoCo-1K

0.775

0.946

778 779

775 776

777

781

781 782

Sup-21K

0.228

0.782

783 784 785

796

797

A.3 MORE EXPERIMENTS FOR DIFFERENT COMBINATIONS OF TWO-STAGE MODELS

iBOT-21K

0.702

0.905

We display more experimental results of different combinations of first-stage and second-stage mod-786 els in Table 9. Here, DINOv2-LVD142M is consistently used as the first-stage model, while various 787 self-supervised pretrained models are employed for the second stage. Our approach achieves sub-788 stantial improvements over SOTA methods under both ID and OOD scenarios. Specifically, on 789 Split CIFAR-100 (ID scenario), we achieve an average improvement of 0.94% over the other four 790 self-supervised models besides DINOv2-LVD142M, with all results outperforming the case where 791 DINOv2-LVD142M is used in both stages. This breaks the limitation discussed in our analysis 792 in Table 5 regarding the relatively limited second-stage improvement, highlighting the superiority 793 of our method. For Split Aircrafts (OOD scenario), we achieve significant gains in all situations, with an average improvement of 6.01% across five models. Particularly, the final performance on 794 iBOT-1k and iBOT-21k surpasses that of using DINOv2-LVD142M in both stages. 795

Table 9: The performance on Split CIFAR-100 and Aircrafts of different methods employing diverse combinations of the first-stage model and second-stage model, compared to our proposed method.

dataset	first-stage model	second-stage model	L2P	Dual	HiDe	Ours
		DINOv2-LVD142M	85.60	85.84	92.54	92.51
		MoCo-1K	75.28	78.00	<u>91.81</u>	93.06
Split CIFAR-100	DINOv2-LVD142M	DINO-1K	72.45	73.20	91.69	92.77
1		iBOT-1K	75.92	74.21	92.92	93.68
		iBOT-21K	78.68	76.76	<u>93.13</u>	93.80
	DINOv2-LVD142M	DINOv2-LVD142M	33.07	29.06	72.61	75.46
		MoCo-1K	29.40	32.46	66.23	73.70
Split Aircrafts		DINO-1K	36.84	38.67	69.29	75.95
		iBOT-1K	38.29	39.54	70.55	75.62
		iBOT-21K	37.21	35.44	<u>68.18</u>	76.16

A.4 PREDICTION CONFIDENCE HISTOGRAMS OF DIFFERENT MODELS

Figure 4 illustrates the prediction confidence histograms of five different models on the test sets of CIFAR-100 and Aircrafts after the first stage of training. Red indicates incorrect classifications, and blue indicates correct classifications. The x-axis represents prediction confidence, and the y-axis represents the number of samples. It is evident that on CIFAR-100, all models exhibit very high confidence for correctly classified samples, whereas misclassified samples have generally lower confidence. Under OOD scenario with Split Aircrafts, supervised pretrained models show a significantly smaller number of correctly predicted samples with high confidence compared to self-supervised pretrained models. For incorrectly predicted samples, self-supervised pretrained models not only have generally lower confidence but also significantly fewer such samples.

