TOWARDS CONTINUAL DOMAIN ADAPTION OF VISION-LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large-scale vision-language models have achieved remarkable performance on various downstream tasks. Nevertheless, how to efficiently adapt vision-language models to new data distributions without re-training, *i.e.*, domain incremental learning (DIL) of vision-language models, is still under-explored. Existing DIL methods for single modality are either not applicable to multi-modal settings or need exemplar buffers to store previous samples to avoid catastrophic forgetting, which is not memory-efficient. To address these limitations, we propose an exemplar-free paradigm to improve DIL of vision-language models based on prompt-tuning. We theoretically analyze and decompose the problem into two optimization objectives. Guided by the theoretical insights, we propose a novel framework named Multimodal Continual Domain Adaptation (MCDA), which incorporates two strategies: Multimodal Domain Alignment (MDA) and Maximum Softmax Gating (MSG). MDA enhances cross-domain performance by aligning visual and language representation spaces, while MSG improves the accuracy of domain identification by gating through Softmax probability. Extensive experimental results demonstrate that our method outperforms current state-of-the-art approaches.

026 027 028

029

024

025

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Large vision-language models have achieved remarkable performance on various downstream tasks (Zhang et al., 2022; Adel et al., 2019; Agarwal et al., 2022; Arani et al., 2021). After large-scale pretraining, these powerful models can make zero-shot predictions without requiring any task-specific training examples. This advantage of large vision-language models is promising for the development of more general-purpose models without further tuning. However, the zero-shot performance on certain tasks suffers from a lack of sufficient relevant image-text pairs in the pre-training corpus. For example, images containing domain shift (Aygün et al., 2022; Bhat et al., 2023; Cao et al., 2021) from one domain (*e.g.*sketch) to another (*e.g.*cartoon) with certain textural descriptions are difficult to collect, even though it's crucial for models' open-domain learning adaptability.

Unfortunately, expanding the knowledge of the vision-language model through re-training from 040 scratch would incur prohibitively high computational costs. One effective way to alleviate the issue 041 is to continually fine-tune or prompt-tune the vision-language model on various domains of data, 042 which is known as domain incremental learning (DIL) (Ardywibowo et al., 2022; Benzing, 2022; 043 Boschini et al., 2022b;a;c). In the DIL setting, the set of classes remains constant, but the domains 044 (data distributions) involved commonly vary a lot in sequence while domain indices are not provided at inference time. The continually learned model can handle any image-text input and can be further used for downstream tasks. Nevertheless, domain incremental learning remains challenging for 046 vision-language models. We find that when learning is performed sequentially on multiple domains, 047 pre-trained vision-language models tend to forget most of the knowledge related to previously learned 048 tasks (a task refers to learning on one domain). The phenomenon is commonly known as catastrophic forgetting. 050

Despite considerable efforts to apply DIL in single-modality settings, these methods either do not apply to multi-modal settings or require exemplar buffers to store previous samples to prevent catastrophic forgetting, which is not memory-efficient. Different from traditional methods that modify all or a subset of the network parameters or store examples in a buffer, a new paradigm arises for

054 continual learning by optimizing a limited number of learnable prompts. As a pioneer working under 055 such a paradigm, S-prompt (Wang et al., 2022a) treats the learning of the prompts independently, 056 which leads to the best performance per domain. It replaces the use of expensive buffers by optimizing per-domain prompts. During training time, they calculate centroids for each domain by applying 058 k-means on the training image features, which are generated with the fixed pre-trained transformer without any prompts. During inference, KNN is used to identify the nearest centroid to the test image and then add the associated domain prompt to the image tokens for classification. Despite the 060 empirical performance gains observed by S-prompting, there is still a gap between empirical success 061 and theoretical analysis. 062

In this paper, we theoretically analyze the paradigm of S-prompting and argue that the design of its component is sub-optimal. The theoretical analysis also shows that its performance is largely limited by wrong-domain prediction and out-of-domain prediction. Motivated by the limitation, we propose a novel framework called Multimodal Continual Domain Adaptation (MCDA) with two strategies:
 Multimodal Domain Alignment (MDA) and Maximum Softmax Gating (MSG) to tackle the issue.
 MSG transforms the problem of domain selection into out-of-domain detection. MDA alleviates the problem of catastrophic forgetting of the vision-language model by forcing the alignment matrix to be similar to that of the previous domain.

Extensive experimental results shows that our approach outperforms state-of-the-art methods. Specifically, MCDA outperforms existing advanced methods of L2P and S-liPrompts with the highest average accuracy score of 89.17% and the lowest average forgetting score of -0.17%. This suggests that MCDA is more adept at learning new information without significantly forgetting previously learned knowledge.

In summary, our contributions are as follows:

- From a theoretical perspective, we analyze the boundedness of a continual learning process of vision-language models, and a clear theorem is presented in Theorem 1.
- From a framework perspective, we propose the MCDA as a novel framework using both vision and language adaptation for enhancing continual learning of vision-language models.
- From an experimental perspective, our proposed MCDA archives new state-of-the-art performances on the CDDB-Hard, DomainNet and CORe50 datasets.

2 RELATED WORK

076

077

078

079

081

082

084 085

087

Continual learning refers to learning scenarios that require models to adapt to a sequence of tasks with varying data distributions. One of the major challenges of continual learning is known as catastrophic forgetting, where models tend to forget most of the knowledge they previously 091 learned after adapting to new data. To tackle catastrophic forgetting, numerous methods have been 092 proposed. Regularization-based methods (Kirkpatrick et al., 2017; Zenke et al., 2017; Aljundi et al., 2018; Chaudhry et al., 2018; Zenke et al., 2017) alleviate catastrophic forgetting by adding explicit regularization terms to balance the old and new tasks. Replay-based methods(Vitter, 1985; Chaudhry 094 et al., 2019; Riemer et al., 2018; Borsos et al., 2020; Caccia et al., 2020) try to approximate and 095 recover previous data distributions. Optimization-based methods(Lopez-Paz & Ranzato, 2017; Zeng 096 et al., 2019; Guo et al., 2022; Kong et al., 2022; Liu & Liu, 2021) focus on designing specific optimization procedures and programs. Architecture-based methods(Xue et al., 2022; Serra et al., 098 2018; Golkar et al., 2019; Jung et al., 2020; Gurbuz & Dovrolis, 2022) focus on constructing 099 task-specific parameters. 100

Domain-incremental learning is one of the most commonly seen scenarios of continual learning. In
 the setting of DIL, each task (domain) has the same data label space but different distributions. Task
 (domain) identities are not available at inference time. DIL is involved in many real-world problems
 such as autonomous driving, where the vehicle meets varying weather conditions in the wild (Mirza et al., 2022).

Prompt tuning methods are different from traditional methods that modify all or a subset of the network parameters or store examples in a buffer, L2P (Wang et al., 2022d) begins a new paradigm for continual learning by optimizing a limited number of learnable prompts. After that, several work



Figure 1: Training Pipeline of MCDA. During the training stage, Multimodal Domain Alignment (MDA) is used to align the image and text representation space between two domains.

(Lester et al., 2021; Li & Liang, 2021; Zhou et al., 2021; 2022; Bahng et al., 2022; Wang et al., 2022c; Douillard et al., 2022a) follow the paradigm and achieve great success. 130

131 Continual learning for vision–Language models is still under-explored. In (Srinivasan et al., 2022), the focus is on robust fine-tuning of VL. In (Wang et al., 2022b), the change between image and 132 text representation space during the pretraining stage of VL is explored. Recently, ZSCL (Zheng 133 et al., 2023) is proposed to solve the problem of zero-shot degradation during the fine-tuning process 134 of VL. However, little research has been done to explore the potential of prompt-tuning to solve 135 DIL. S-Prompts (Wang et al., 2022a) is the pioneering work to apply prompt-tuning to DIL. It trains 136 different prompts for each domain and dynamically selects the appropriate set during testing using 137 a fixed key/value dictionary. Recently, Hide-Prompt (Wang et al., 2023) decomposes the problem 138 of continual learning into several parts that can be optimized, and derives theoretical analysis of 139 performance. Despite its great success, it focuses on class incremental learning scenarios and could 140 not be directly applied to DIL.

141 142 143

144

145

125

126 127 128

129

3 METHOD

3.1 PROBLEM FORMULATION

146 Denote as $S = \{D_s\}_{s=1}^N$ as the sequence of datasets presented to the model in our incremental 147 learning scenario. Denote each dataset as $\mathcal{D}_s = \{\mathbf{x}_i^s, \mathbf{y}_i^s\}_{i=1}^{|\mathcal{D}_s|}$, where \mathbf{x}_i represents an image, and 148 $\mathbf{y}_i \in \{0,1\}^K$ is its corresponding one-hot label for K target classes. By convention, in the setting of 149 DIL, we are only allowed to one domain \mathcal{D}_s at a time. Each time a new domain \mathcal{D}_s arrives, the goal 150 of DIL is to improve the model's performance on \mathcal{D}_s and alleviate the catastrophic forgetting for past 151 domains $\mathcal{D}_{s-1}, \mathcal{D}_{s-2}, \ldots, \mathcal{D}_1$. 152

153 154

3.2 THEORETICAL ANALYSIS

For domain-incremental learning (DIL), let $\mathcal{X}_t = \bigcup_j \mathcal{X}_{t,j}$ and $\mathcal{Y}_t = \{\mathcal{Y}_{t,j}\}$, where $j \in \{1, \dots, |\mathcal{Y}_t|\}$ 156 denotes the j-th class in task t. Now assume we have a ground event denoted as $\mathcal{D} = \{\mathcal{D}_1, \ldots, \mathcal{D}_t\}$ 157 and a pre-trained model f_{θ} . For any sample $x \in \bigcup_{k=1}^{t} \mathcal{X}_k$, a general goal of the DIL problem is to 158 learn $P(x \in \mathcal{X}_{*,j} \mid \mathcal{D}, \theta)$, where $\mathcal{X}_{*,j}$ represents the *j*-th class domain in any task. Of note, $\mathcal{Y}_t = \mathcal{Y}_{t'}$, 159 $\forall t \neq t' \text{ for DIL.}$ 160

Denote domain identification (DI), within-domain prediction (WDP) and out-of-domain-prediction 161 (ODP) as $P(\boldsymbol{x} \in \mathcal{X}_i \mid \mathcal{D}, \theta), P(\boldsymbol{x} \in \mathcal{X}_{i,j} \mid \boldsymbol{x} \in \mathcal{X}_i, \mathcal{D}, \theta)$ and $P(\boldsymbol{x} \in \mathcal{X}_{k,j} \mid \boldsymbol{x} \in \mathcal{X}_k, k \neq i, \mathcal{D}, \theta)$ 162 respectively. Based on Bayes' theorem, we have 163

$$\begin{split} P\left(\boldsymbol{x} \in \mathcal{X}_{*,j} \mid \mathcal{D}, \theta\right) &= P\left(\boldsymbol{x} \in \mathcal{X}_{i,j} \mid \boldsymbol{x} \in \mathcal{X}_{i}, \mathcal{D}, \theta\right) P\left(\boldsymbol{x} \in \mathcal{X}_{i} \mid \mathcal{D}, \theta\right) \\ &+ \sum_{k \neq i} P\left(\boldsymbol{x} \in \mathcal{X}_{k,j} \mid \boldsymbol{x} \in \mathcal{X}_{k}, \mathcal{D}, \theta\right) P\left(\boldsymbol{x} \in \mathcal{X}_{k} \mid \mathcal{D}, \theta\right), \end{split}$$

166 167 168

169 170 where $\{*, j\}$ represents the j-th class in each domain. Thus, the problem of DIL can be written as

$$P(\boldsymbol{x} \in \mathcal{X}_{*,j} \mid \mathcal{D}, \theta) = P_{DI} \cdot P_{WDP} + (1 - P_{DI}) \cdot P_{ODP}.$$
(2)

(1)

From the formula above, it is clear to see that DIL is constrained by within-domain-prediction, 171 domain identification, and out-of-domain prediction. In other words, the performance of DIL could 172 be improved by enhancing these three objectives. 173

174 Within-domain-prediction is related to domain generalization of vision-language models, and there 175 have already been many works trying to improve it. One thing to mention is that our decomposition is similar to HiDe-prompt (Wang et al., 2023) which decomposes the problem of class incremental 176 learning into three objectives. However, in this paper, we focus on domain incremental learning and 177 the decomposition of Hide-prompt does not include ODP. Inspired by HiDe-prompt (Wang et al., 178 2023), we define 179

$$H_{\text{WDP}}(\boldsymbol{x}) = \mathcal{H}(\mathbf{1}_{\bar{j}}, \{P(\boldsymbol{x} \in \mathcal{X}_{\bar{i},j} | \boldsymbol{x} \in \mathcal{X}_{\bar{i}}, \mathcal{D}, \theta)\}_{j}),$$
(3)

 $H_{\mathrm{DI}}(\boldsymbol{x}) = \mathcal{H}(\mathbf{1}_{\overline{i}}, \{P(\boldsymbol{x} \in \mathcal{X}_i | \mathcal{D}, \theta)\}_i),$ (4)

$$H_{\text{ODP}}(\boldsymbol{x}) = \mathcal{H}(\boldsymbol{1}_{\bar{c}}, P(\boldsymbol{x} \in \mathcal{X}_{k,j} \mid \boldsymbol{x} \in \mathcal{X}_k, k \neq i, \mathcal{D}, \theta)),$$
(5)

where H_{WDP} , H_{DI} , and H_{ODP} are the cross-entropy values of WDP, DI, and ODP, respectively. The operation $\mathcal{H}(p,q) \triangleq -\mathbb{E}_n[\log q]$ stands for one-hot encoding function. We now present Theorem 1.

185 187

184

188 189

190

191

193

194

Theorem 1 If $\mathbb{E}_{\boldsymbol{x}}[H_{\text{WDP}}(\boldsymbol{x})] \leq \epsilon$, $\mathbb{E}_{\boldsymbol{x}}[H_{\text{DI}}(\boldsymbol{x})] \leq \gamma$, and $\mathbb{E}_{\boldsymbol{x}}[H_{\text{ODP}}(\boldsymbol{x})] \leq \eta$, we have the loss error $\mathcal{L} \in [0, \delta + \epsilon + \log(1 + e^{\epsilon - \delta}(e^{\gamma} - 1))]$

A detailed proof of the theorem is in Appendix. Theorem 1 shows that optimizing WDP, DI, and ODP can help improve the performance of DIL. In this paper, we focus on improving the last two 192 objectives, namely out-of-domain prediction (ODP) and domain identification, and propose two strategies: Multimodal Domain Alignment (MDA) and Maximum Softmax Gating (MSG). The overall training pipeline is shown in Figure 1.

195 196 197

3.3 MULTIMODAL DOMAIN ALIGNMENT

At training time, an independent set of prompts is trained for each domain and is frozen when training 199 on subsequent domain data. At inference time, the sample needs to identify which domain it comes 200 and select the corresponding set of prompts for that domain. If the sample successfully identify the 201 domain it belongs to, then forgetting will not happen. However, if the sample fail to identify the 202 correct domain, then it will select a wrong set of prompts. Therefore, the sample will be tested on 203 a model trained on domain data that are different from its original domain. In other words, if the 204 domain identification process goes wrong, then the model will be tested on out-of-distribution (OOD) data. Intuitively, the performance is expected to drop when a model is tested with OOD data. To 205 validate the statement, we use models prompt-tuned with different domain sources to test on data 206 from different domains and the results are in Figure 2 (a). 207

208 From the figure, we can see that no matter which source domain the model is trained on, it will suffer 209 performance drop when tested on OOD data. We hypothesis that the performance drop is caused by 210 the misalignment between the text representation space and image representation space. Formally, define f(x) as the image feature generated by the image encoder f, and $\{g(t_j)\}_j^C$ is a set of weight 211 212 vectors produced by the text encoder g(.) with each $g(t_i)$ representing the text feature of j-class. 213 The prediction probability of vision-language model is computed as 214

215

$$p(y_j \mid x) = \frac{\exp\left(\langle f(x), g(t_j) \rangle\right)}{\sum_{k=1}^C \exp\left(\langle f(x), g(t_k) \rangle\right)},\tag{6}$$

(1.0()) (.)))



Figure 2: (a) Performance of models prompt-tuned with different domain sources to test on data from different domains; (b) Average cosine similarity of 100 positive image-text pairs from five different domains

233 234

230

231

232

where $\langle .,. \rangle$ is the cosine similarity, and $\langle x, y \rangle = \frac{x \cdot y}{|x||y|}$. 235

236 From the formula, we see that the prediction process relies on the mutual interaction between image 237 representation space and text representation space. For models prompt-tuned on different domain 238 sources, since both the image and text are prompt-tuned with separate the sets of prompts, it is possible 239 that their image representation space and text representation space are changed synchronously.

240 To validate the hypothesis, we select 100 positive image-text pairs from five different domains, 241 feed them into models trained on different domain data and get their corresponding image and text 242 representations. Then we calculate their average cosine similarity and the result is shown in Figure 2 243 (b). We can see from the figure when using encoders trained on different domain data to generate 244 representations for these image-text pairs, the average cosine similarity indeed drops. Therefore, the 245 performance drop is caused by the misalignment between the text representation space and image 246 representation space.

247 To tackle the issue, we proposed a strategy called Multimodal Domain Alignment (MDA), which can 248 effectively alleviate the misalignment between image and text representation space during training. 249 Specifically, at training time, suppose that we are training on domain data D_t . Given the input 250 image-text batch, we first calculate the contrastive matrix M_t using the current set of domain-specific 251 prompts. Then, we calculate the contrastive matrix O_t generated by the original vision-language 252 model without any prompts. Then we force the alignment between these two contrastive matrices 253 using KL Divergence.

$$L_{KL}^{t}\left(M_{t},O_{t}\right) = -\sum O_{t} \ln\left(\frac{M_{t}}{O_{t}}\right).$$
(7)

256 The overall training loss is:

$$L = L_{contras} + \alpha L_{KL}^t, \tag{8}$$

262

254 255

257

3.4 MAXIMUM SOFTMAX GATING

263 At inference time, the test sample needs to identify which domain it comes from and then select the 264 corresponding set of domain-specific prompts. If the sample fails to identify the correct domains 265 it belongs to, then as discussed above, a performance drop will happen. Therefore, improving 266 the domain identification accuracy is another key factor in improving DIL performance under the 267 proposed paradigm. Previous work, such as S-prompt, takes a naive approach by storing prototype centroids at training time and calculating the similarity of testing sample and prototype centroids to 268 select domain at inference time. We argue that this is sub-optimal and that the problem of domain 269 identification can be transformed into the problem of OOD detection.

where $L_{contras}$ is the contrastive loss used in regular CLIP training and α is a constant.



Figure 3: Testing Pipeline of MCDA. During inference stage, Maximum Softmax Gating (MSG) is used to do OOD test of given test image on each domain. Then KNN selection will be performed on domains that survive the OOD tests.

Maximum Softmax Score (MSS) (Ming et al., 2022) is first proposed as a metric for zero-shot OOD detection. We find that MSS is also effective for prompt-tuned vision-language models, thus benefiting the procedure of domain identification. For a given domain with label set $\mathcal{Y}_{in} = \{y_1, y_2, \ldots, y_K\}$, the text prototypes are set as $\mathcal{T}(t_i)$, $i \in \{1, 2, \ldots, K\}$, where t_i is the text prompt. Then, image and text features are generated by models with a corresponding set of domain-specific prompts. For any test input image \mathbf{x}' , we can calculate the label-wise matching score based on the cosine similarity between the image feature $\mathcal{I}(\mathbf{x}')$ and the text vector $\mathcal{T}(t_i) : s_i(\mathbf{x}') = \frac{\mathcal{I}(\mathbf{x}') \cdot \mathcal{T}(t_i)}{\|\mathcal{I}(\mathbf{x}')\| \cdot \|\mathcal{T}(t_i)\|}$. Formally, the Maximum Softmax Score (MSS) is defined as:

where τ is the temperature. For in-domain data, it will be matched to one of the text features with a high score. The OOD detection function can be formulated as:

 $S_{\text{MSS}}\left(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}\right) = \max_{i} rac{e^{s_{i}\left(\mathbf{x}'
ight)/ au}}{\sum_{j=1}^{K} e^{s_{j}\left(\mathbf{x}'
ight)/ au}},$

$$G(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) = \begin{cases} 1 & S_{\text{MSS}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) \ge \lambda \\ 0 & S_{\text{MSS}}(\mathbf{x}'; \mathcal{Y}_{\text{in}}, \mathcal{T}, \mathcal{I}) < \lambda \end{cases},$$
(10)

(9)

where by convention 1 represents the positive class (ID) and 0 indicates OOD. λ is chosen so that a high fraction of ID data (e.g., 95%) is above the threshold. At inference time, given test data, we conduct an OOD test on all the domains using the OOD detection above. Ideally, only one domain will accept and the other will reject. If that is the case, then the domain that accepts is the one that the test data belongs to. However, we find that in fact, several domains accept, and others reject. Therefore, we propose Maximum Softmax Gating (MSG) to filter out all the domains that fail the OOD test. For domains that survive the OOD test, we will select among them using the same metric as the S-prompt. The overall procedure is shown in Figure 3.

- 315 4 EXPERIMENT

4.1 DATASETS

We perform experiments on three standard DIL benchmark datasets: CDDB (Li et al., 2023), CORe50 (Lomonaco & Maltoni, 2017), and DomainNet (Peng et al., 2019).

CDDB is a dataset for continual deepfake detection, which designs easy, long, and hard tracks.
 Particularly, we choose the most challenging track (*i.e.*, the Hard track) that requests learning on 5 se quential deepfake detection domains, which are GauGAN, BigGAN, WildDeepfake, WhichFaceReal, and SAN respectively.

327 328	Method	Prompts	Buffer size	AA (†)	$AF(\downarrow)$
329	LRCIL (Pellegrini et al., 2020)	×		76.39	-4.39
330	iCaRL (Marra et al., 2019)	×	100ex/class	79.76	-8.73
330	LUCIR (Hou et al., 2019)	×		82.53	-5.34
332	LRCIL (Pellegrini et al., 2020)	×		74.01	-8.62
222	iCaRL (Marra et al., 2019)	×	50ex/class	73.98	-14.50
333	LUCIR (Hou et al., 2019)	×		80.77	-7.85
334	DyTox (Douillard et al., 2022b)	\checkmark		86.21	-1.55
335	EWC (Kirkpatrick et al., 2017)	×		50.59	-42.62
336	LwF (Li & Hoiem, 2017)	×		60.94	-13.53
337	DyTox (Douillard et al., 2022b)	\checkmark	No buffer	51.27	-45.85
338	L2P (Wang et al., 2022d)	\checkmark	No buffer	61.28	-9.23
339	S-liPrompts (Wang et al., 2022a)	\checkmark	No buffer	88.65	-0.69
340	MCDA (Ours)	\checkmark	No buffer	89.27	-0.07

324 Table 1: Results on CDDB-Hard. Evaluation of existing state-of-the-art DIL methods in the standard 325 DIL setting. The best results are highlighted in **bold**.

CORe50 is a widely used dataset for continual object recognition that has 50 categories from 11 distinct domains. The continual learning setting uses 8 domains for incremental training and the rest domains as the test set.

346 DomainNet is a dataset for domain adaptation and domain incremental learning, which has 345 347 categories and roughly 600,000 images. The images in this dataset are split into 6 domains. The 348 DIL setup on DomainNet is the same as that of S-liprompt. We report the average forward detection 349 accuracy and the average forgetting degree on CDDB-Hard. For CORe50 and DomainNet, we report 350 the average forward classification accuracy. 351

352 353

360 361

362

326 327

341 342 343

344

345

4.2 COMPARISON METHODS

354 Following S-prompt (Wang et al., 2022a), we benchmark our proposed methods against state-of-355 the-art DIL methods. These include non-prompting methods: EWC (Kirkpatrick et al., 2017), LwF 356 (Li & Hoiem, 2017), ER (Chaudhry et al., 2019), GDumb (Prabhu et al., 2020), BiC (Wu et al., 357 2019), DER++ (Buzzega et al., 2020) and Co2L (Cha et al., 2021), prompting-based methods: L2P 358 (Wang et al., 2022d), DyTox (Douillard et al., 2022b) and S-liPrompts (Wang et al., 2022a) and a 359 self-supervised learning method: CaSSLe (Fini et al., 2022).

4.3 RESULTS

363 We evaluate our proposed approach in the standard DIL scenario. Table 1 shows the performance of our method on the challenging CDDB-Hard dataset: Referring to Table 1, we can see that MCDA 364 outperforms all previous state-of-the-art methods. Compared with methods without using prompts, MCDA has clear superiority with an average relative improvement of 17%. Besides, these prompt-366 free methods normally require exemplar buffers, whereas our method does not. Our method is also 367 superior to prompt-based DyTox by a large margin. 368

MCDA achieves better performance than DyTox without the use of a buffer, which is more memory-369 efficient. Compared with recent prompt-based approaches such as L2P and S-liPrompts, our method 370 is still superior. This is largely due to our design of MDA and MSG, which effectively alleviates the 371 forgetting issue of continual domain incremental training of vision-language models and improves 372 the domain selection accuracy. 373

374 Table 2 (a) shows the performance of our method on the DomainNet dataset. Compared with 375 the CDDB-hard dataset, DomainNet contains more heterogeneous domains and thus is more challenging to handle domain shift. Due to memory efficiency concerns, we see exemplar-376 free methods as our real competitors. CaSSLe is an exemplar method that can be used jointly 377 with other self-supervised learning methods and MCDA has nearly 20% performance gain over

378	Table 3: Results of ablating MDA on CDDB-Hard
379	for exemplar-free deepfake DIL. For a fair compari-
380	son, MSG is not applied. $\alpha = 0$: MDA is not used.

Method	Average Acc (†)	Forgetting (†)
MCDA ($\alpha = 0$)	88.65	-0.69
MCDA ($\alpha = 0.1$)	88.71	-0.63
MCDA ($\alpha = 0.5$)	88.95	-0.39
MCDA ($\alpha = 1$)	88.93	-0.41
MCDA ($\alpha = 5$)	79.15	-10.19

Table 4: Performance under different domain selection strategies on CDDB-hard. We show that MSG is effective even when it is applied to random domain selection.

_	Method	KNN/Random	MSG	$AA(\uparrow)$	
	MCDA	KNN	X	72.43	
	MCDA	Random	X	68.57	
	MCDA	KNN	\checkmark	89.17	
_	MCDA	Random	\checkmark	79.58	
_					-

Table 5: Average domain identification accuracy on CDDB-Hard. We show the domain identification accuracy improvement after applying MSG.

	GauGAN	BigGAN	WildDeepfake	WhichFaceReal	SAN	Average Acc (†)
KNN w/o MSG	0.67	0.78	0.94	0.95	0.56	0.63
KNN w/ MSG	0.83	0.89	0.97	0.97	0.79	0.88

CaSSLe. We hypothesize that although CaSSLe achieves memory-efficient domain-incremental learning by using prompts, its performance is restricted by the nature of self-supervised learning.

Compared with L2P, MCDA also 400 achieves a superior performance 401 gain of around 28%. This 402 is because L2P was initially 403 proposed for class-incremental 404 learning. Although the paradigm 405 of L2P can be applied to domain-406 incremental learning, it cannot 407 handle the domain shift well among different tasks. Con-408 trarily, S-liPrompts, and MCDA 409 both handle the domain shift 410

Table 2: Results on DomainNet. The results are reported as the Accuracy (Acc) metric, where the best values are highlighted in bold.

Method	Prompt	Buffer size	AA (†)
DyTox (Douillard et al., 2022b)	\checkmark	50ex/class	62.94
LwF 'te'pli2017learning	\checkmark		49.2
CaSSLe (Fini et al., 2022) w/ SimCLR (Chen et al., 2020)	X		44.2
CaSSLe w/ BYOL (Grill et al., 2020)	X		49.7
CaSSLe w/ Barlow Twins (Zbontar et al., 2021)	X	No buffer	48.9
CaSSLe w/ SupCon (Khosla et al., 2020)	X		50.9
L2P (Wang et al., 2022d)	\checkmark	No buffer	40.1
S-liPrompts (Wang et al., 2022a)	\checkmark	No buffer	67.7
MCDA (Ours)	\checkmark	No buffer	70.3

problem by maintaining a set of separately trained prompts for each domain, which effectively 411 solves the problem of domain shift. MCDA also achieves better performance than S-liPrompts, which 412 is largely due to the design of MDA which regularizes the alignment between visual and language 413 space, and MSG which improves the accuracy of domain selection. 414

Table 7 (b) shows the performance of our method on the CORe50 dataset. 415 Different from the previous two datasets, 416

CORe50 contains unseen tested domains 417 that do not appear in the incremental train-418 ing stage. This requires the method to be 419 capable of generalizing to unseen domains 420 well. Previous methods, no matter using 421 prompts or buffers or not, fail to take OOD 422 generalization into their design concern. 423 However, MDA in our methods naturally allows for better OOD generalization abil-424 ity, which is the reason why it outperforms 425 all previous methods on OOD testing do-426 mains. 427

4 78
- T & U

430

3

388

389

396

397

398 399

381

429 4.4 ABLATION STUDY

Table 7: Results on CORe50. The results are reported as the Accuracy (Acc) metric, where the best values are highlighted in **bold**.

Method	Prompt	Buffer size	$AA(\uparrow)$
GDumb (Prabhu et al., 2020)	X		74.92
DER++ (Buzzega et al., 2020)	X	50ex/class	79.70
DyTox (Douillard et al., 2022b)	\checkmark		79.21
L2P (Wang et al., 2022d)	\checkmark		81.07
EWC (Kirkpatrick et al., 2017)	\checkmark		74.82
LwF (Li & Hoiem, 2017)	\checkmark		75.45
L2P (Wang et al., 2022d)	\checkmark	No buffer	78.33
S-liPrompts (Wang et al., 2022a)	\checkmark	No buffer	89.06
MCDA (Ours)	\checkmark	No buffer	92.37

Effect of Multimodal Domain Alignment (MDA). To validate the effectiveness of our methods, 431 we first conduct an ablation study on the two proposed strategies in MCDA. Table 3 shows the result

	S 1	S2	S 3	OOD1	OOD2	OOD3	AA
S1 (ours)	99.78	87.72	47.63	51.09	63.46	72.38	66.41
S2 (ours)	99.92	99.21	52.31	73.56	74.38	81.50	73.26
S3 (ours)	99.90	98.75	82.16	73.29	66.83	78.15	75.31
S3 (S-liPrompts)	97.83	73.21	69.21	62.43	61.87	71.42	71.53

Table 6: Average accuracy on OOD domains on CDDB-hard. We show that MCDA enjoys superior
 domain generalization ability compared with previous methods.

of changing the value of hyperparameter α in the training loss on the CDDB-hard dataset. From the results, we can see that when applying MDA and setting α to be in a proper range, our method outperforms the baseline method. When α is set to 0, this means that we do not apply MDA in the training process. α stands for the balance between learning the distribution of the current domain and forcing the model to align the model with a previous domain. From the results, we can see that if α is too large, then performance will drop. We hypothesize that this is because if the weight assigned to the alignment matrix is too large, then it will affect the original adaptation performance.

449 Effect of Maximum Softmax Gating (MSG). MSG is proposed to improve the domain iden-450 tification accuracy. Table 5 shows the accuracy of domain selection when applying MSG on the 451 CDDB-hard dataset. From the results, we can see that after applying MSG, the average domain identification accuracy is improved from 63% to 88%. We find that without applying MSG, the model 452 tends to misidentify the test samples between the GauGAN domain and the BigGAN domain. As 453 shown in Table 5, the domain identification accuracy of these two domains both do not reach 80%. 454 However, after applying MSG, the model will first do an OOD test of each domain before selecting it. 455 Therefore, part of the testing samples that would be misclassified will be filtered out during the OOD 456 test. Consequently, the accuracy of domain selection will be improved. 457

Exploration of different domain identification strategies. MCDA adopts KNN as the domain selection strategy. At the inference stage, the extracted features of the testing image are used to query the domain using the KNN algorithm. Experiments have shown that MSG can help improve the domain selection accuracy when using KNN as the identification strategy. To explore whether MSG is still effective when using other identification strategies, we conduct experiments when domains are randomly selected during inference. Table 4 shows the between using KNN and random selection. We can see that even under the setting of random selection, MSG can still improve the performance.

465

442

443

444

445

446

447

448

Exploration of generalization capability of the model. Following Wang et al. (2022a), we apply 466 the trained MCDA prompts on \$1-\$3 to out-of-domain OOD1-OOD3 which are 3 unseen domains in 467 the CDDB-hard dataset. We choose S1-S3 to be: GauGAN, BigGAN, WildDeepfake; and choose 468 OOD1-OOD3 to be: FaceForensic++, Glow, StarGAN. From the results in Table 6, we can see that 469 our method can achieve good performance on OOD domains, with an average accuracy gain of 4% 470 compared with Sli-Prompts (which achieves the second-best generalization performance among 471 all methods). This indicates that our method is capable of handling cases when there are unseen 472 domains in the testing stage. 473

474 475

476 477

5 DISCUSSION AND CONCLUSION

In this paper, we theoretically analyze and decompose the problem of DIL into two optimization objectives. Guided by the theoretical insights, we propose two strategies: Multimodal Domain Alignment (MDA) and Maximum Softmax Gating (MSG). MDA improves the model's cross-domain performance by forcing the alignment between visual and language representation spaces. MSG improves the accuracy of domain identification by gating through softmax probability. Experiments show that our method outperforms existing state-of-the-art methods.

The limitation of the proposed MCDA is mainly in two aspects. On the one hand, MSG could still be
 improved to enhance the domain selection accuracy. On the other hand, our method is restricted for
 the setting of DIL. Making MCDA adaptable for all CL setting remains to be explored.

486 REFERENCES

493

510

516

523

- Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw). In
 International Conference on Learning Representations, 2019.
- Aishwarya Agarwal, Biplab Banerjee, Fabio Cuzzolin, and Subhasis Chaudhuri. Semantics-driven generative replay for few-shot class incremental learning. In *Proceedings of the ACM International Conference on Multimedia*, pp. 5246–5254, 2022.
- Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars.
 Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference* on Computer Vision, pp. 139–154, 2018.
- Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual
 learning method based on complementary learning system. In *International Conference on Learning Representations*, 2021.
- Randy Ardywibowo, Zepeng Huo, Zhangyang Wang, Bobak J Mortazavi, Shuai Huang, and Xiaoning Qian. Varigrow: Variational architecture growing for task-agnostic continual learning based on bayesian novelty. In *International Conference on Machine Learning*, pp. 865–877. PMLR, 2022.
- Eser Aygün, Ankit Anand, Laurent Orseau, Xavier Glorot, Stephen M Mcaleer, Vlad Firoiu, Lei M
 Zhang, Doina Precup, and Shibl Mourad. Proving theorems using incremental learning and
 hindsight experience replay. In *International Conference on Machine Learning*, pp. 1198–1210.
 PMLR, 2022.
- ⁵⁰⁸ Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modify ⁵⁰⁹ ing pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 2022.
- Frederik Benzing. Unifying importance based regularisation methods for continual learning. In International Conference on Artificial Intelligence and Statistics, pp. 2372–2396. PMLR, 2022.
- Prashant Shivaram Bhat, Bahram Zonooz, and Elahe Arani. Task-aware information routing from common representation space in lifelong learning. In *International Conference on Learning Representations*, 2023.
- Zalán Borsos, Mojmir Mutny, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *Advances in Neural Information Processing Systems*, 33:14879–14890, 2020.
- Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class incremental continual learning into the extended der-verse. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022a.
- Matteo Boschini, Lorenzo Bonicelli, Angelo Porrello, Giovanni Bellitto, Matteo Pennisi, Simone Palazzo, Concetto Spampinato, and Simone Calderara. Transfer without forgetting. *arXiv preprint arXiv:2206.00388*, 2022b.
- Matteo Boschini, Pietro Buzzega, Lorenzo Bonicelli, Angelo Porrello, and Simone Calderara. Continual semi-supervised learning through contrastive interpolation consistency. *Pattern Recognition Letters*, 162:9–14, 2022c.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020.
- Lucas Caccia, Eugene Belilovsky, Massimo Caccia, and Joelle Pineau. Online learned continual compression with adaptive quantization modules. In *International Conference on Machine Learning*, pp. 1240–1250. PMLR, 2020.
- Yue Cao, Hao-Ran Wei, Boxing Chen, and Xiaojun Wan. Continual learning for neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3964–3974, 2021.
 - Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co21: Contrastive continual learning. In ICCV, 2021.

540 541 542	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</i> , pp. 532–547, 2018.
543 544 545 546	Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc'Aurelio Ranzato. On tiny episodic memories in continual learning. <i>arXiv preprint arXiv:1902.10486</i> , 2019.
547 548 549	Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In <i>ICML</i> , 2020.
550 551	Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In <i>CVPR</i> , 2022a.
552 553 554 555	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, 2022b.
556 557 558	Enrico Fini, Victor G Turrisi Da Costa, Xavier Alameda-Pineda, Elisa Ricci, Karteek Alahari, and Julien Mairal. Self-supervised models are continual learners. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9621–9630, 2022.
559 560 561	Siavash Golkar, Michael Kagan, and Kyunghyun Cho. Continual learning via neural pruning. <i>arXiv</i> preprint arXiv:1903.04476, 2019.
562 563 564	Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. In <i>NeurIPS</i> , 2020.
565 566 567	Yiduo Guo, Wenpeng Hu, Dongyan Zhao, and Bing Liu. Adaptive orthogonal projection for batch and online continual learning. 2, 2022.
568 569 570	Mustafa B Gurbuz and Constantine Dovrolis. Nispa: Neuro-inspired stability-plasticity adaptation for continual learning in sparse networks. In <i>International Conference on Machine Learning</i> , pp. 8157–8174. PMLR, 2022.
571 572 573	Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In <i>CVPR</i> , 2019.
574 575 576 577	Sangwon Jung, Hongjoon Ahn, Sungmin Cha, and Taesup Moon. Continual learning with node- importance based adaptive group sparse regularization. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 33:3647–3658, 2020.
578 579	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In <i>NeurIPS</i> , 2020.
580 581 582 583 584	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. In <i>Proceedings of the national academy of sciences</i> , 2017.
585 586 587	Yajing Kong, Liu Liu, Zhen Wang, and Dacheng Tao. Balancing stability and plasticity through advanced null space in continual learning. In <i>European Conference on Computer Vision</i> , pp. 219–236. Springer, 2022.
588 589 590	Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. <i>arXiv preprint arXiv:2104.08691</i> , 2021.
591 592 593	Chuqiao Li, Zhiwu Huang, Danda Pani Paudel, Yabin Wang, Mohamad Shahbazi, Xiaopeng Hong, and Luc Van Gool. A continual deepfake detection benchmark: Dataset, methods, and essentials. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 1339–1349, 2023.

594 595	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. <i>arXiv</i> preprint arXiv:2101.00190, 2021.
590 597 598	Zhizhong Li and Derek Hoiem. Learning without forgetting. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 40(12):2935–2947, 2017.
599 600 601	Hao Liu and Huaping Liu. Continual learning with recursive gradient optimization. In International Conference on Learning Representations, 2021.
602 603	Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. <i>CoRR</i> , abs/1705.03550, 2017. URL http://arxiv.org/abs/1705.03550.
604 605 606	David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. Advances in Neural Information Processing Systems, 30, 2017.
607 608	Francesco Marra, Cristiano Saltori, Giulia Boato, and Luisa Verdoliva. Incremental learning for the detection and classification of gan-generated images. In <i>WIFS</i> , 2019.
609 610 611	Yifei Ming, Ziyang Cai, Jiuxiang Gu, Yiyou Sun, Wei Li, and Yixuan Li. Delving into out-of- distribution detection with vision-language representations, 2022.
612 613	M. Jehanzeb Mirza, Marc Masana, Horst Possegger, and Horst Bischof. An efficient domain- incremental learning approach to drive in all weather conditions, 2022.
615 616	Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for real-time continual learning. In <i>IROS</i> , 2020.
617 618 619 620	Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 1406–1415, 2019.
621 622 623	Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In <i>ECCV</i> , 2020.
624 625 626	Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. In <i>ICLR</i> , 2018.
627 628 629	Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In <i>ICML</i> , 2018.
630 631 632	Tejas Srinivasan, Ting-Yun Chang, Leticia Leonor Pinto Alva, Georgios Chochlakis, Mohammad Rostami, and Jesse Thomason. Climb: A continual learning benchmark for vision-and-language tasks. <i>arXiv preprint arXiv:2206.09059</i> , 2022.
633 634 635	Jeffrey S Vitter. Random sampling with a reservoir. ACM Transactions on Mathematical Software (TOMS), 11(1):37–57, 1985.
636 637	Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical decomposition of prompt-based continual learning: Rethinking obscured sub-optimality, 2023.
638 639 640 641 642	Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained transformers: An occam's razor for domain incremental learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), <i>Advances in Neural Information Processing Systems</i> , 2022a. URL https://openreview.net/forum?id=ZVe_WeMold.
643 644 645	Zhen Wang, Liu Liu, Yiqun Duan, Yajing Kong, and Dacheng Tao. Continual learning with lifelong vision transformer. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 171–181, 2022b.
040	Zifeng Wang, Zizhao Zhang, Chen-Vu Lee, Han Zhang, Ruovi Sun, Xiaogi Ren, Guolong Su, Vincent

648 649 650	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 <i>Proceedings</i> of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 139–149, 2022d.
652 653	Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In <i>CVPR</i> , 2019.
654 655 656	Mengqi Xue, Haofei Zhang, Jie Song, and Mingli Song. Meta-attention for vit-backed continual learning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 150–159, 2022.
657 658 659	Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In <i>ICML</i> , 2021.
660 661	Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning of context-dependent processing in neural networks. <i>Nature Machine Intelligence</i> , 1(8):364–372, 2019.
663 664	Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In <i>International Conference on Machine Learning</i> , pp. 3987–3995. PMLR, 2017.
665 666 667	Yanzhe Zhang, Xuezhi Wang, and Diyi Yang. Continual sequence generation with adaptive compositional modules. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 3653–3667, 2022.
668 669 670	Zangwei Zheng, Mingyuan Ma, Kai Wang, Ziheng Qin, Xiangyu Yue, and Yang You. Preventing zero-shot transfer degradation in continual learning of vision-language models, 2023.
671 672	Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision- language models. <i>arXiv preprint arXiv:2109.01134</i> , 2021.
673 674 675 676	Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In <i>CVPR</i> , 2022.
677	
678	
679	
680	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
694	
605	
696	
697	
698	
699	
700	
701	

702 Appendix

A **PROOF OF THEOREM 1**

$$H_{DIL}(x) = H\left(y, \{\mathbf{P} (x \in \mathbf{X}_{k,j} \mid D)\}_{k,j}\right)$$

= $-\sum_{k,j} y_{k,j} \log \mathbf{P} (x \in \mathbf{X}_{k,j} \mid D)$
= $-\log \mathbf{P} (x \in \mathbf{X}_{k_0,j_0} \mid D),$ (11)

where $H_{DIL}(x)$ represents the entropy of the distribution over the sets $\mathbf{X}_{k,j}$, given the data D and the probabilities of x being in these sets. The summation simplifies to a single term based on the highest likelihood for a specific set \mathbf{X}_{k_0,j_0} .

$$H_{WDP}(x) = H\left(\tilde{y}, \{\mathbf{P}\left(x \in \mathbf{X}_{k_0, j} \mid x \in \mathbf{X}_{k_0}, D\right)\}_j\right)$$
$$= -\sum_j y_{k_0, j} \log \mathbf{P}\left(x \in \mathbf{X}_{k_0, j} \mid x \in \mathbf{X}_{k_0}, D\right)$$
(12)

$$= -\log \mathbf{P} \left(x \in \mathbf{A}_{k_0, j_0} \mid x \in \mathbf{A}_{k_0}, D \right),$$

Here, $H_{WDP}(x)$ describes the conditional entropy given that x is in \mathbf{X}_{k_0} , evaluating the likelihood across subsets $\mathbf{X}_{k_0,j}$.

$$H_{DI}(x) = H\left(\bar{y}, \{\mathbf{P} (x \in \mathbf{X}_k \mid D)\}_k\right)$$

= $-\sum_k \bar{y}_k \log \mathbf{P} (x \in \mathbf{X}_k \mid D)$
= $-\log \mathbf{P} (x \in \mathbf{X}_{k_0} \mid D),$ (13)

In this case, $H_{DI}(x)$ is the entropy over the sets \mathbf{X}_k , which gives the likelihood of x belonging to any set k. This again reduces to the most probable set \mathbf{X}_{k_0} .

$$H_{ODP}(x) = H\left(\tilde{y}, \{\mathbf{P} (x \in \mathbf{X}_{k_i, j} \mid x \in \mathbf{X}_{k_i}, D, k_i \neq k_0)\}_j\right)$$

$$= -\sum_j y_{k_i, j} \log \mathbf{P} (x \in \mathbf{X}_{k_i, j} \mid x \in \mathbf{X}_{k_i}, D)$$

$$= -\log \mathbf{P} (x \in \mathbf{X}_{k_i, j} \mid x \in \mathbf{X}_{k_i}, D)$$
(14)

$$= -\log \mathbf{I} \ \left(x \in \mathbf{A}_{k_i, j_0} \mid x \in \mathbf{A}_{k_i}, D \right),$$

This describes the entropy $H_{ODP}(x)$ for the case where x belongs to a different set k_i (with $k_i \neq k_0$), and it evaluates the probability within the subsets $\mathbf{X}_{k_i,j}$.

$$\begin{array}{ll} 745 \\ 746 \\ 747 \\ 748 \\ 749 \\ 750 \\ 751 \end{array} \qquad H_{DIL}(x) = H\left(y, \{\mathbf{P}\left(x \in \mathbf{X}_{k,j} \mid D\right)\}_{k,j}\right) \\ = -\sum_{k,j} y_{k,j} \log \mathbf{P}\left(x \in \mathbf{X}_{k,j} \mid D\right) \\ \leq -\log\left(P_{WDP}P_{DI} + P_{ODP}(1 - P_{DI})\right) \\ = -\log\left(e^{-\epsilon}e^{-\delta} + e^{-\gamma}(1 - e^{-\delta})\right) \end{aligned}$$
(15)

$$= \delta + \epsilon + \log\left(1 + e^{\epsilon - \delta}(e^{\gamma} - 1)\right).$$

Finally, this inequality shows the bound on $H_{DIL}(x)$ by combining the weighted probabilities of xbelonging to either \mathbf{X}_{k_0,j_0} (with probability $P_{WDP}P_{DI}$) or \mathbf{X}_{k_i,j_0} (with probability $P_{ODP}(1-P_{DI})$). The resulting expression involves the exponents ϵ , δ , and γ , providing a closed-form solution.

B IMPLEMENTATION DETAILS

760

761

762

763

We implement the proposed MCDA framework in PyTorch with NVIDIA RTX 4090 GPU. The image encoder is implemented with the architecture of ViT-B/16 and the text encoder is the same as the text encoder in CLIP. The embedding dimension is 768 for both image and text encoder. We adopt SGD optimizer with a momentum of 0.9, an initial learning rate of 0.1, a batch size of 128, and a cosine scheduler. The learning epoches for CDDB dataset is 50, and for DomainNet and CORe5 are 10.

764 765 766

767 768 769

770

C ABLATION ON MORE DATASETS

C.1 DOMAINNET

771Table 8 presents the results of ablating the Multimodal Domain Alignment (MDA) strategy on the
DomainNet dataset for exemplar-free deepfake domain incremental learning (DIL). To ensure a fair
comparison, the Maximum Softmax Gating (MSG) strategy was not applied in these experiments.774The table shows the average accuracy and forgetting metric under different values of the α parameter,
which controls the strength of the MDA alignment.

776 The results reveal that introducing MDA with $\alpha = 0.5$ achieves the best overall performance, with an 777 average accuracy of 69.5% and the lowest forgetting value of -0.39. This indicates that a moderate 778 level of MDA alignment significantly improves cross-domain performance while effectively reducing 779 forgetting. Notably, the performance degrades when α is set to extremes: for $\alpha = 0$, which means 780 no MDA alignment is applied, the average accuracy drops to 67.1%, and the forgetting increases 781 to -0.69. Similarly, setting $\alpha = 5$ leads to a substantial decline in both accuracy (63.9%) and forgetting (-10.19), suggesting that overemphasizing MDA can negatively impact the model's ability 782 to generalize across domains. 783

⁷⁸⁴ In summary, the results demonstrate that MDA is crucial for improving the performance of MCDA ⁷⁸⁵ in domain incremental learning. The optimal balance is achieved at $\alpha = 0.5$, where the trade-off ⁷⁸⁶ between accuracy and forgetting is the most favorable.

Table 9 compares the performance of different domain selection strategies on DomainNet, with and without the application of MSG. The two domain selection methods under evaluation are K-NN and random domain selection.

The results show that the KNN domain selection strategy consistently outperforms random domain selection. When no MSG is applied, the KNN-based MCDA achieves an accuracy of 69.8%, while the random domain selection results in a significantly lower accuracy of 65.4%. This suggests that KNN is more effective in identifying domains that can benefit from domain incremental learning.

When MSG is applied, the performance of both strategies improves. Specifically, random domain
selection with MSG achieves an accuracy of 66.2%, while KNN with MSG achieves 67.1%. These
results indicate that MSG is effective in improving domain selection, even when the selection process
is randomized.

799

800Table 8: Results of ablating MDA on DomainNet for801exemplar-free deepfake DIL. For a fair comparison,802MSG is not applied. $\alpha = 0$: MDA is not used.

803			
804	Method	Average Acc (†)	Forgetting (†)
805	MCDA ($\alpha = 0$)	67.1	-0.69
806	MCDA ($\alpha = 0.1$)	68.4	-0.63
807	MCDA ($\alpha = 0.5$)	69.5	-0.39
808	MCDA ($\alpha = 1$)	66.1	-0.41
809	MCDA ($\alpha = 5$)	63.9	-10.19

Table 9: Performance under different domain selection strategies on DomainNet. We show that MSG is effective even when it is applied to random domain selection.

Method	KNN/Random	MSG	AA (†)
MCDA	KNN	X	69.8
MCDA	Random	X	65.4
MCDA	KNN	\checkmark	67.1
MCDA	Random	\checkmark	66.2

_

810 C.2 CORE50

811

812Table 10 shows the results of ablation studies on the CORe50 dataset, focusing on the impact of813the Multimodal Domain Alignment (MDA) strategy in exemplar-free deepfake domain incremental814learning (DIL). For these experiments, Maximum Softmax Gating (MSG) was not applied to ensure a815controlled comparison. The table reports average accuracy and forgetting metrics for different values816of the α parameter, which controls the strength of the MDA alignment.

817 The results indicate that introducing MDA with $\alpha = 0.5$ achieves the best performance, with an 818 average accuracy of 91.33% and the lowest forgetting value of -0.39. This shows that moderate MDA 819 alignment significantly enhances cross-domain performance while reducing forgetting. Interestingly, increasing α to 1 does not lead to further improvement, with a slight drop in average accuracy to 820 91.32% and forgetting remaining at -0.41. Furthermore, setting α to an extreme value of 5 leads to a 821 drastic performance degradation, with an accuracy of 88.54% and forgetting of -10.19, indicating that 822 excessive MDA alignment hampers the model's ability to generalize across domains. In summary, 823 the ablation study demonstrates that MDA is critical for improving the performance of MCDA in 824 domain incremental learning, with the optimal alignment strength being $\alpha = 0.5$, striking the best 825 balance between accuracy and forgetting. 826

Table 11 compares the performance of different domain selection strategies on the CORe50 dataset,
both with and without the application of MSG. The two domain selection methods evaluated are
K-Nearest Neighbors (KNN) and random selection.

830 When MSG is not applied, the KNN-based domain selection yields the highest accuracy of 91.65%, 831 while random selection results in a lower accuracy of 89.76%. This illustrates the effectiveness of KNN in selecting relevant domains for domain incremental learning. However, when MSG is 832 applied, both strategies see a slight performance decrease, with KNN achieving 89.72% and random 833 selection 88.79%. These results indicate that while MSG is generally beneficial, its interaction with 834 different domain selection methods may vary, and in this case, KNN alone proves more effective than 835 combining it with MSG. The result shows that selecting appropriate domain strategies is crucial, and 836 KNN consistently outperforms random selection. The results also suggest that the additional use of 837 MSG should be considered carefully based on the specific domain selection method being employed. 838

839 840

841

842

851

854

855

856

857 858

859

860

861

Table 10: Results of ablating MDA on CORe50 for exemplar-free deepfake DIL. For a fair comparison, MSG is not applied. $\alpha = 0$: MDA is not used.

Table 11: Performance under different domain selection strategies on CORe50. We show that MSG is effective even when it is applied to random domain selection.

Method	Average Acc (†)	Forgetting (†)	-				
MCDA ($\alpha = 0$)	89.02	-0.69	_	Method	KNN/Random	MSG	$AA(\uparrow)$
MCDA ($\alpha = 0.1$)	89.97	-0.63		MCDA	KNN	X	91.65
MCDA ($\alpha = 0.5$)	91.33	-0.39		MCDA	Random	X	89.76
MCDA ($\alpha = 1$)	91.32	-0.41		MCDA	KNN	\checkmark	89.72
MCDA ($\alpha = 5$)	88.54	-10.19	_	MCDA	Random	\checkmark	88.79

References

- Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw). In
 International Conference on Learning Representations, 2019.
 - Aishwarya Agarwal, Biplab Banerjee, Fabio Cuzzolin, and Subhasis Chaudhuri. Semantics-driven generative replay for few-shot class incremental learning. In *Proceedings of the ACM International Conference on Multimedia*, pp. 5246–5254, 2022.

Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference* on Computer Vision, pp. 139–154, 2018.

Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual
 learning method based on complementary learning system. In *International Conference on Learning Representations*, 2021.

 Randy Ardywibowo, Zepeng Huo, Zhangyang Wang, Bobak J Mortazavi, Shuai Huang, and Xiaoning Qian. Varigrow: Variational architecture growing for task-agnostic continual learning based on bayesian novelty. In *International Conference on Machine Learning*, pp. 865–877. PMLR, 2022.

- Eser Aygün, Ankit Anand, Laurent Orseau, Xavier Glorot, Stephen M Mcaleer, Vlad Firoiu, Lei M
 Zhang, Doina Precup, and Shibl Mourad. Proving theorems using incremental learning and
 hindsight experience replay. In *International Conference on Machine Learning*, pp. 1198–1210.
 PMLR, 2022.
- Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 2022.
- Frederik Benzing. Unifying importance based regularisation methods for continual learning. In International Conference on Artificial Intelligence and Statistics, pp. 2372–2396. PMLR, 2022.
- Prashant Shivaram Bhat, Bahram Zonooz, and Elahe Arani. Task-aware information routing from
 common representation space in lifelong learning. In *International Conference on Learning Representations*, 2023.
- Zalán Borsos, Mojmir Mutny, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *Advances in Neural Information Processing Systems*, 33:14879–14890, 2020.
- Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class incremental continual learning into the extended der-verse. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 2022a.
- Matteo Boschini, Lorenzo Bonicelli, Angelo Porrello, Giovanni Bellitto, Matteo Pennisi, Simone
 Palazzo, Concetto Spampinato, and Simone Calderara. Transfer without forgetting. *arXiv preprint arXiv:2206.00388*, 2022b.
- Matteo Boschini, Pietro Buzzega, Lorenzo Bonicelli, Angelo Porrello, and Simone Calderara. Continual semi-supervised learning through contrastive interpolation consistency. *Pattern Recognition Letters*, 162:9–14, 2022c.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark
 experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020.
- Lucas Caccia, Eugene Belilovsky, Massimo Caccia, and Joelle Pineau. Online learned continual com pression with adaptive quantization modules. In *International Conference on Machine Learning*,
 pp. 1240–1250. PMLR, 2020.
- Yue Cao, Hao-Ran Wei, Boxing Chen, and Xiaojun Wan. Continual learning for neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3964–3974, 2021.
- 903 Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *ICCV*, 2021.
- Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian
 walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision*, pp. 532–547, 2018.
- Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K
 Dokania, Philip HS Torr, and Marc'Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
 contrastive learning of visual representations. In *ICML*, 2020.
- Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *CVPR*, 2022a.
- Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers
 for continual learning with dynamic token expansion. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 9285–9295, 2022b.

918	Enrico Fini, Victor G Turrisi Da Costa, Xavier Alameda-Pineda, Elisa Ricci, Karteek Alahari, and	
919	Iulien Mairal Self-supervised models are continual learners. In <i>Proceedings of the IEEE/CVF</i>	
920	Conference on Computer Vision and Pattern Recognition pp 9621–9630 2022	
921		
922	Siavash Golkar, Michael Kagan, and Kyunghyun Cho. Continual learning via neural pruning. arXiv	
923	preprint arXiv:1903.04476, 2019.	
02/		
924	Jean-Bastien Grill, Florian Strub, Florent Altche, Corentin Tallec, Pierre Richemond, Elena	
925	Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaonan Guo, Mohammad Gheshlaghi Azar,	
926	et al. Bootstrap your own latent-a new approach to self-supervised learning. In <i>NeurIPS</i> , 2020.	
927	Yiduo Guo, Wenneng Hu, Dongyan Zhao, and Bing Liu, Adaptive orthogonal projection for batch	
928	and online continual learning. 2, 2022.	
929		
930	Mustafa B Gurbuz and Constantine Dovrolis. Nispa: Neuro-inspired stability-plasticity adaptation	
931	for continual learning in sparse networks. In International Conference on Machine Learning, pp.	
932	8157–8174. PMLR, 2022.	
933	Caibui Hay Visum Day Char Charge Ley 7:14 Ways and Dahus Lin Learning a wife datasifer	
934	incrementally via rabalancing. In CVDP, 2010	
935	incrementariy via rebalancing. In CVPK, 2019.	
936	Sangwon Jung, Hongjoon Ahn, Sungmin Cha, and Taesup Moon. Continual learning with node-	
937	importance based adaptive group sparse regularization. Advances in Neural Information Processing	
020	Systems, 33:3647–3658, 2020.	
930		
939	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron	
940	Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In <i>NeurIPS</i> , 2020.	
941	James Kirknatrick Dezven Desceny, Neil Pabinowitz, Joel Veness, Guilloume Designding, Andrei A	
942	Busu Kieran Milan John Quan Tiago Ramalho Agnieszka Grabska Barwinska et al. Overcoming	
943	constraint will all young and the second sec	
944	2017	
945	2017.	
946	Yajing Kong, Liu Liu, Zhen Wang, and Dacheng Tao. Balancing stability and plasticity through	
947	advanced null space in continual learning. In European Conference on Computer Vision, pp.	
948	219–236. Springer, 2022.	
949		
950	Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt	
051	tuning. arXiv preprint arXiv:2104.08691, 2021.	
951	Chuqiao Li Zhiwu Huang Danda Pani Paudel Vahin Wang Mohamad Shahhazi Visoneng Hong	
952	and Luc Van Gool. A continual deepfake detection benchmark. Dataset methods and essentials	
953	In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp.	
954	1339–1349, 2023.	
955		
956	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. arXiv	
957	preprint arXiv:2101.00190, 2021.	
958	Thishong Li and Derak Hojem Learning without forgetting IEEE transactions on rattern analysis	
959	and machine intelligence, 40(12):2025, 2047, 2017	
960	<i>una machine intentigence</i> , 40(12).2955–2947, 2017.	
961	Hao Liu and Huaping Liu. Continual learning with recursive gradient optimization. In International	
962	Conference on Learning Representations, 2021.	
963		
964	Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object	
965	recognition. Cokk, abs/1/05.03550, 2017. UKL http://arxiv.org/abs/1705.03550.	
330	David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning	
067	Advances in Neural Information Processing Systems 30 2017	
301	The access of the information 1 rocessing systems, 50, 2017.	
900	Francesco Marra, Cristiano Saltori, Giulia Boato, and Luisa Verdoliva. Incremental learning for the	
909	detection and classification of gan-generated images. In WIFS, 2019.	
970	Vitai Ming Ziyong Cai Liuviang Cu Viyou Sun Wai Li and Viyuan Li Dalving internet of	
971	distribution detection with vision-language representations, 2022.	

980

994

- M. Jehanzeb Mirza, Marc Masana, Horst Possegger, and Horst Bischof. An efficient domainincremental learning approach to drive in all weather conditions, 2022.
- Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for
 real-time continual learning. In *IROS*, 2020.
- 977 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching
 978 for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on* 979 *computer vision*, pp. 1406–1415, 2019.
- Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, 2020.
- Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro.
 Learning to learn without forgetting by maximizing transfer and minimizing interference. In *ICLR*, 2018.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *ICML*, 2018.
- Tejas Srinivasan, Ting-Yun Chang, Leticia Leonor Pinto Alva, Georgios Chochlakis, Mohammad
 Rostami, and Jesse Thomason. Climb: A continual learning benchmark for vision-and-language
 tasks. *arXiv preprint arXiv:2206.09059*, 2022.
- Jeffrey S Vitter. Random sampling with a reservoir. ACM Transactions on Mathematical Software (TOMS), 11(1):37–57, 1985.
- Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical decomposition of prompt-based continual learning: Rethinking obscured sub-optimality, 2023.
- Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained transformers:
 An occam's razor for domain incremental learning. In Alice H. Oh, Alekh Agarwal, Danielle
 Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022a.
 URL https://openreview.net/forum?id=ZVe_WeMold.
- Zhen Wang, Liu Liu, Yiqun Duan, Yajing Kong, and Dacheng Tao. Continual learning with lifelong vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 171–181, 2022b.
- 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*, 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 139–149, 2022d.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, 2019.
- Mengqi Xue, Haofei Zhang, Jie Song, and Mingli Song. Meta-attention for vit-backed continual
 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 pp. 150–159, 2022.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, 2021.
- Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning of context-dependent
 processing in neural networks. *Nature Machine Intelligence*, 1(8):364–372, 2019.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pp. 3987–3995. PMLR, 2017.
- Yanzhe Zhang, Xuezhi Wang, and Diyi Yang. Continual sequence generation with adaptive compositional modules. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3653–3667, 2022.

Zangwei Zheng, Mingyuan Ma, Kai Wang, Ziheng Qin, Xiangyu Yue, and Yang You. Preventing zero-shot transfer degradation in continual learning of vision-language models, 2023.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision language models. *arXiv preprint arXiv:2109.01134*, 2021.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, 2022.

1033 1034 1035

1036

A APPENDIX

You may include other additional sections here.

1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1070
10//

1078 1079