Algorithm 1: DualPrompt at training time

**Input:** Pre-trained transformer-based backbone f, final classification layer  $f_{\phi}$ , number of tasks T, training set  $\{(\boldsymbol{x}_{i,t}, y_{i,t})\}_{i=1}^{n_t}\}_{t=1}^T$ , G-Prompt  $\boldsymbol{g}$ , E-Prompt  $\mathbf{E} = \{ \boldsymbol{e}_t \}_{t=1}^T, \, \text{task keys } \mathbf{K} = \{ \boldsymbol{k}_t \}_{t=1}^T, \, \texttt{start}_{\boldsymbol{g}}, \texttt{end}_{\boldsymbol{g}}, \texttt{start}_{\boldsymbol{e}}, \texttt{end}_{\boldsymbol{e}}, \, \texttt{prompting} \}$ function  $f_{\text{prompt}}$ , number of training epochs of the *t*-th task  $M_t$ Initialize:  $\phi$ , g, E, K for  $t = 1, \cdots, T$  do Select the task-specific E-Prompt  $e_t$  and corresponding task key  $k_t$ Generate the prompted architecture  $f_{g,e_t}$ : attach g and  $e_t$  to start<sub>g</sub>-th to  $end_g$ -th and  $start_e$ -th to  $end_e$ -th MSA layers respectively, with  $f_{prompt}$ . for  $e = 1, \cdots, M_t$  do Draw a mini-batch  $B = \{(\boldsymbol{x}_{i,t}, y_{i,t})\}_{i=1}^{l}$ for  $(\boldsymbol{x}, y)$  in B do Calculate the prompted feature by  $f_{g,e_t}(x)$ Calculate the per sample loss  $\mathcal{L}_x$  via equation 6 end Update  $\phi$ , g,  $\mathbf{E}$ ,  $\mathbf{K}$  by backpropagation end  $\mathbf{end}$ 

Algorithm 2: DualPrompt at test time

Given components: Pre-trained transformer-based backbone f, trained classification layer  $f_{\phi}$ , G-Prompt g, E-Prompt  $\mathbf{E} = \{e_t\}_{t=1}^T$ , task keys  $\mathbf{K} = \{k_t\}_{t=1}^T$ , startg, endg, starte, ende, prompting function  $f_{\text{prompt}}$ Input: test example  $\boldsymbol{x}$ Generate query feature  $q(\boldsymbol{x})$ Matching for the index of E-Prompt via  $t_{\boldsymbol{x}} = \operatorname{argmin}_t \gamma(q(\boldsymbol{x}), \boldsymbol{k}_t)$ Select the task-specific E-Prompt  $e_{t_{\boldsymbol{x}}}$ Generate the prompted architecture  $f_{g,e_{t_{\boldsymbol{x}}}}$ : attach g and  $e_{t_{\boldsymbol{x}}}$  to startg-th to endg-th and starte-th to ende-th MSA layers respectively, with  $f_{\text{prompt}}$ . Prediction:  $f_{g,e_{t_{\boldsymbol{x}}}}(\boldsymbol{x})$ 

### A Algorithms for DualPrompt

The training and test time Algorithms for DualPrompt are illustrated in Algorithm 1 and 2, respectively.

### **B** Experimental details

For our method, DualPrompt, we use a constant learning rate of 0.005 usng Adam [18] optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ , and a batch size of 128 for all benchmarks. For methods with the ViT architecture, all input images are resized to  $224 \times 224$  and normalized to [0, 1]. Otherwise, we follow their

#### 20 Z. Wang, Z. Zhang et al.

original implementation in their paper. We set the balancing factor  $\lambda = 1$ in equation 6. We train Split CIFAR-100 and 5-datasets for 5 epochs per task and Split ImageNet-R for 50 epochs to ensure models converge properly for each task, thus the issue of forgetting is disentangled from possible underfitting [3]. We further sample 20% of the training set of Split ImageNet-R as a validation set for searching the optimal  $\mathtt{start}_g, \mathtt{end}_g, \mathtt{start}_e, \mathtt{end}_e$ , and empirically set  $\mathtt{start}_g = 1, \mathtt{end}_g = 2, \mathtt{start}_e = 3, \mathtt{end}_e = 5$  for all the setting, since we discover they perform consistently well for all datasets. Following the suggestion of prompt length by [59], we set  $L_g = 5$  and  $L_e = 20$  for all datasets as well, and we further verify the correctness of this choice in Appendix G. Note that for fair comparison, we set  $M = 30, L_p = 20, N = 5$  for L2P, which leads to similar amount of parameters as DualPrompt.

To ensure fair comparison, every aforementioned methods start from the same ImageNet pre-trained ViT-B/16 [9], following the setting in [59]. We carefully reimplement these method and use hyper-parameters by referring to their original source code. Moreover, we make the pre-trained model fully trainable for all methods (except L2P and DualPrompt), as we empirically observe they could not learn as good with a frozen backbone due to limited learning capacity.

### C Evaluation metrics

Let  $S_{t,\tau}$  be the evaluation score, *e.g.*, classification accuracy on the  $\tau$ -th task after training on the *t*-th task. After the model finishes training on the *t*-th task, we compute the Average Accuracy  $(A_t)$  and Forgetting  $(F_t)$  as follows:

$$A_{t} = \frac{1}{t} \sum_{\tau=1}^{t} S_{t,\tau}$$
$$F_{t} = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \max_{\tau' \in \{1, \cdots, t-1\}} \left( S_{\tau',\tau} - S_{t,\tau} \right)$$

Note that Average Accuracy is the overall evaluation metric for continual learning, which includes two aspects: greater learning capacity and less catastrophic forgetting, while Forgetting only serves as a measure of catastrophic forgetting.

#### D Details of comparing methods

- Regularization-based methods. EWC [19] and LwF [27] are representative regularization-based methods that are widely compared.
- Rehearsal-based methods. ER [7,13], GDumb [46], BiC [61], DER++ [3] and Co<sup>2</sup>L [4]. As earlier methods, ER and GDumb achieve very strong performance not only in their own work, but in later literature [34,3] as well. BiC is also a strong method in the class-incremental learning setting. DER++ and Co<sup>2</sup>L are the latest SOTA methods. We chose a medium and a large buffer size for these rehearsal-based methods, based on recommendation of prior work [3,40,4,59].



Fig. 5: Representative examples from Split ImageNet-R.

- Prompt-based method. L2P [59]is the current state-of-the-art promptbased method, we configure DualPrompt to have the similar amount of additional parameters as L2P for fair comparison.
- Architecture-based methods. All these methods are based on ResNet-18, as recommended in the original work. We either directly taken reported results in their original work or strictly follow the original implementation and hyper-parameter settings of these methods to reproduce the results. Some other well-known methods [17,11,45] are not compared here because they either have been outperformed by our compared methods or only verified on simpler task-incremental setting.

# E Split ImageNet-R: large intra-class diversity

Figure 5 shows some representative examples of three different classes from Split ImageNet-R. We can observe that although the images of the same row share the same label, they actually differ a lot. This observation accords well with the result in Figure 1 and Table 1 that rehearsal-based methods require a large buffer size to perform well due to large intra-class diversity in Split ImageNet-R.

## F Searching for multi-layered prompts

Based on the strategy in 4.2, we first search for the multi-layered E-Prompt on the validation set of Split ImageNet-R. Due to the large search space, we made a simplified assumption that the MSA layers to attach prompts should be contiguous. Moreover, since  $\mathtt{start}_e = \mathtt{end}_e = 5$  is the best option in the single layer case, it is natural to include the 5-th layer when searching for multi-layered E-Prompt. Nevertheless, we also include several cases when the 5-th layer is not included for completeness.

21

$i_g  j_g $	Avg. Acc	$\overline{i_e \ j}$	$i_e   \text{Avg. Acc}$
5 5	65.55	2 2	2 67.70
$3 \ 4$	65.76	1 2	2 68.46
4 5	66.59	1 3	3 67.79
$5 \ 6$	66.12	1 5	5 67.73
$3 \ 5$	67.12	6 8	8 65.10
4 6	66.41	1 1	2 63.13
$5 \ 7$	64.53		•
1 12	67.09		

Table 6: Searching results for multi-

layered G-Prompt.

Table 5: Searching results for multilayered E-Prompt.

We discover that  $\mathtt{start}_e = 3$ ,  $\mathtt{end}_e = 5$  yields the best performance in terms of average accuracy. Note that when we attach E-Prompt to every MSA layer ( $\mathtt{start}_e = 1, \mathtt{end}_e = 12$ ), it actually leads to comparable accuracy. However, we still choose  $\mathtt{start}_e = 3$ ,  $\mathtt{end}_e = 5$  since it has less additional parameters.

We then fix  $\mathtt{start}_e = 3$ ,  $\mathtt{end}_e = 5$ , and search for multi-layered G-Prompt, given that we have  $\mathtt{start}_g = \mathtt{end}_g = 2$  yields the best performance as a single-layered E-Prompt. We conduct similar searching process, with a preference of including the 2nd layer. The searching results are shown in Table 6.

We discover that  $\mathtt{start}_g = 1, \mathtt{end}_g = 2$  leads to best average accuracy. Moreover, simply share all MSA layers results in overall negative effect on the accuracy.

Although our search strategy is not exhaustive, we find the combination of  $\mathtt{start}_{g} = 1, \mathtt{end}_{g} = 2, \mathtt{start}_{e} = 3, \mathtt{end}_{e} = 5$  works quite well for all benchmark datasets.

## G Searching for prompt length

Following the suggestion by [59], we use set the base length of prompts as 5, and perform grid search on the lengths of G-Prompt and E-Prompt from  $\{5, 10, 20, 40\} \times \{5, 10, 20, 40\}$ , based on the optimal positions obtained in the previous step.

As shown in Figure 6, we indeed verify  $L_g = 5, L_e = 20$  is the optimal choice on the validation set of Split ImageNet-R. We also empirically observe this configuration works well on other datasets, thus we use this combination of prompt lengths for other datasets in experiments showed in the main text.

#### H Additional results on 5-datasets

For completeness, we also demonstrate the effectiveness of our method on 5datasets [11], which is a collection of five diverse image classification datasets, CIFAR-10 [21], MNIST [23], Fashion-MNIST [62], SVHN [42], and notMNIST [2].



Fig. 6: Grid search on prompt length.

Method	Buffer size	<b>5-da</b> Avg. Acc (†)	tasets Forgetting $(\downarrow)$
$\begin{array}{c} \text{ER} \\ \text{BiC} \\ \text{DER++} \\ \text{Co}^{2}\text{L} \end{array}$	250	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 15.69{\pm}0.89\\ 21.15{\pm}1.00\\ 14.38{\pm}0.35\\ 17.52{\pm}1.35\end{array}$
$\begin{array}{c} \text{ER} \\ \text{BiC} \\ \text{DER++} \\ \text{Co}^{2}\text{L} \end{array}$	500	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 12.85 {\pm} 0.62 \\ 10.27 {\pm} 1.32 \\ 10.46 {\pm} 1.02 \\ 12.28 {\pm} 1.44 \end{array}$
FT-seq EWC LwF L2P <b>DualPrompt</b>	0	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 94.63 {\pm} 0.68 \\ 34.94 {\pm} 0.07 \\ 38.01 {\pm} 0.28 \\ 4.64 {\pm} 0.52 \\ \textbf{2.21 {\pm} 0.69} \end{array}$
Upper-bound	-	$93.93 {\pm} 0.18$	-

Table 7: Additional results on 5-datasets [11].

Despite the simplicity of each task in 5-datasets, the benchmark mimics the real-world setting where task diversity is large, thus contributing to a more comprehensive evaluation of CL methods. Since each task in 5-dataset is relatively easier than that of Split CIFAR-100 and ImageNet-R, rehearsal-based methods generally require smaller buffer size to perform well. However, a buffer size of 500 is already considered large for 5-datasets [11,40]. Although DualPrompt still consistently outperforms competing methods, we argue that in real world continual learning scenarios, it is rare that tasks are too diverse: for example, we don't expect a digit classifier to continually learn to classify animals. Thus, we only show the performance on 5-datasets as a proof-of-concept that our method works well when task diversity is large.

23

#### 24 Z. Wang, Z. Zhang et al.

Table 8: Comparison between DualPrompt with query strategy introduced in Section 4.1, and with the perfect match (known test time task identity) on Split ImageNet-R.

Method	Matching Acc	$ $ Avg Acc ( $\uparrow$ ) $ $	Forgetting $(\downarrow)$
DualPrompt-Query	55.8	68.13	4.68
DualPrompt-Perfect Match	100	71.97	3.95

# I Relationship between query accuracy and performance

To demonstrate the relationship between query accuracy and performance, we compare DualPrompt with the query strategy introduced in 4.1 to DualPrompt with known test time task identity to select E-Prompt. The result is shown in Table 8. Interestingly, although the matching accuracy is not that great, DualPrompt is quite robust to it and still achieves an accuracy very close to perfect match. We contribute this robustness to mismatching to the design of our method. First, the task-invariant instruction captured in G-Prompt still remains useful for prediction even if E-Prompt is noisy. Second, the query strategy is based on the input feature, thus implicitly taking into account task similarity. Even when mismatching happens, our method tends to choose the E-Prompt from one of the most similar tasks. We also note that their is still forgetting even if we use ground truth task identity to select the corresponding task-specific E-Prompt. This part of forgetting results from the bias in the final softmax classification head [34,46,67], a common issue in class-incremental learning that could be mitigated in parallel.