

Supplementary Material: Memory-Efficient Incremental Learning Through Feature Adaptation

Ahmet Iscen¹, Jeffrey Zhang², Svetlana Lazebnik², and Cordelia Schmid¹

¹ Google Research

² University of Illinois at Urbana-Champaign

A Algorithm

An overview of our framework is described in Algorithm ??.

Algorithm 1 Memory-efficient incremental learning

```

1: procedure ALGORITHM(Training examples  $\mathcal{X}$ , labels  $\mathcal{Y}$ )
2:   Given  $T$  tasks
3:   *** Train first task ***
4:    $\mathcal{X}^1, \mathcal{Y}^1 \in \mathcal{X}, \mathcal{Y}$  ▷ Data examples for first task
5:    $\theta, W = \text{OPTIMIZE}(L_{CE}(f_{\theta, W}(\mathcal{X}^1), \mathcal{Y}^1))$  ▷ (4)
6:    $h_{\theta}^1 = h_{\theta}$  ▷ Freeze feature extractor
7:    $\mathbf{M}^1 = h_{\theta}^1(\mathcal{X}^1)$  ▷ Store feature descriptors of images
8:    $\mathbf{M}^1 = \text{HERDING}(\mathbf{M}^1)$  ▷ Reduce number of stored features
9:   for  $t \in [2, \dots, T]$  do
10:    *** Train incremental tasks ***
11:     $\mathcal{X}^t, \mathcal{Y}^t \in \mathcal{X}, \mathcal{Y}$  ▷ Data examples for current task
12:     $\theta, W = \text{OPTIMIZE}(L(f_{\theta, W}(\mathcal{X}^t), \mathcal{Y}^t))$  ▷ (6)
13:     $h_{\theta}^t = h_{\theta}$ 
14:     $\phi_{\psi} = \text{FEATURE ADAPTATION}(h_{\theta}^t, h_{\theta}^{t-1}, \mathcal{X}^t, \mathcal{Y}^t)$ 
15:     $\mathbf{M}^t = h_{\theta}^t(\mathcal{X}^t)$  ▷ Store new feature descriptors
16:     $\mathbf{M}^t = \text{HERDING}(\mathbf{M}^t)$ 
17:     $\mathbf{M}^t = \mathbf{M}^t \cup \phi_{\psi}(\mathbf{M}^{t-1})$  ▷ Adapt stored features
18:     $\bar{W} = \text{TRAIN CLASSIFIER}(\mathbf{M}^t, \mathcal{Y}^1, \dots, \mathcal{Y}^t)$  ▷ (Sec. 4.3)
1: procedure FEATUREADAPTATION( $h_{\theta}^{old}, h_{\theta}^{new}, \mathcal{X}, \mathcal{Y}$ )
2:   *** Returns transformation function ***
3:    $\bar{\mathbf{V}} = h_{\theta}^{old}(\mathcal{X})$  ▷ Feature descriptors of old extractor
4:    $\mathbf{V} = h_{\theta}^{new}(\mathcal{X})$  ▷ Feature descriptors of new extractor
5:    $\psi = \text{OPTIMIZE}(L_{FA}(\bar{\mathbf{V}}, \mathbf{V}, \mathcal{Y}))$  ▷ (7)
6:   return  $\phi_{\psi}$ 

```

B Feature Adaptation Quality

We evaluate the quality of our feature adaptation process by measuring the average similarity between the adapted features and their *ground-truth value*. We compute the feature adaptation quality as explained in Section 5.4. However, we compute two distinct measurements this time. ω_{t-1} measures the average feature adaptation quality of features extracted in the previous task (*i.e.* $y \in \mathcal{C}^{t-1}$). This measurement does not

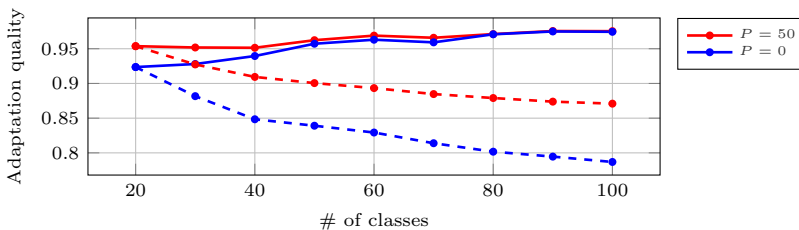


Fig. 1. Feature adaptation quality on CIFAR-100 for $M = 10$. P refers to the number of images preserved in the memory. Solid and dashed lines correspond to vectors from previous (ω_{t-1}) and first (ω_1) task respectively.

track the quality over time, but shows feature adaptation quality between two tasks. ω_1 measures the feature adaptation quality of features originally extracted in the first task (*i.e.* $y \in \mathcal{C}^1$), showing how much the adapted features can diverge from their optimal representation due to accumulated error. Adaptation quality is computed for all $L = 500$ feature descriptors per class.

Figure ?? shows the adaptation quality ω_{t-1} and ω_1 on CIFAR-100 with $M = 10$. We report the quality measures for when P number images are also preserved in the memory. We observe that $P = 0$ achieves ω_{t-1} greater than 0.9 in all tasks. It increases as more classes are seen, most likely due to the network becoming more stable. After 10 tasks, ω_1 is still close to 0.8, indicating our feature adaptation is still relatively successful after training 9 subsequent tasks with no preserved images. Adaptation quality improves as P increases, showing that preserving images also helps with learning a better feature adaptation.

C Balanced Feature Classifier Training

Wu *et al.* [53] illustrated that training a classifier on fewer examples for previous classes introduces a bias towards new classes. To verify the robustness of our method under this setting, we investigate the effect of training our feature classifier on class-balanced and class-imbalanced training sets.

In our main experiments, we train our feature classifier $g_{\bar{V}}$ with a balanced number of examples per class. We repeat our ImageNet-100 experiment in Table 1 without balancing the classifier training samples. In the unbalanced setting, the old classes contain 250 features per class in memory, whereas the new classes each contain ~ 1300 feature vectors. In the balanced setting, all classes contain 250 feature vectors. On ImageNet-100, we achieve 0.893 accuracy with class-imbalanced training, compared to 0.913 accuracy with class-balanced training (reported in Table 1). This shows that even though more training data is utilized in the class-imbalanced setting, the imbalanced class bias leads to a drop in overall performance.

Our method addresses this problem by building a large balanced feature set for training. Our feature adaptation method not only reduces the memory footprint compared to [53] (see Table 1), but also allows substantially more stored data points from old classes (250 features per class compared to 20 images per class for [53]). This may explain some improvement in our results over previous methods. Lastly, the significant increase in number of stored examples provides flexibility to remove examples to keep classes balanced.