

---

# Coarse-To-Fine Incremental Few-Shot Learning

---

Xiang Xiang, Yuwen Tan, Qian Wan, Jing Ma  
School of Artificial Intelligence and Automation  
Huazhong University of Science and Technology, China  
xex@hust.edu.cn

## Abstract

Different from fine-tuning models pre-trained on a large-scale dataset of preset classes, class-incremental learning (CIL) aims to recognize novel classes over time without forgetting pre-trained classes. However, a given model will be challenged by test images with finer-grained classes, e.g., a basenji is at most recognized as a dog. Such images form a new training set (i.e., support set) so that the incremental model is hoped to recognize a basenji (i.e., query) as a basenji next time. This paper formulates such a hybrid natural problem of coarse-to-fine few-shot (C2FS) recognition as a CIL problem named C2FSCIL, and proposes a simple, effective, and theoretically-sound strategy Knowe: to learn, normalize, and freeze a classifier’s weights from fine labels, once learning an embedding space contrastively from coarse labels. Besides, as CIL aims at a stability–plasticity balance, new overall performance metrics are proposed. In that sense, on CIFAR-100, BREEDS, and tieredImageNet, Knowe outperforms all recent relevant CIL/FSCIL methods that are tailored to the new problem setting for the first time.

## 1 Introduction

Product visual search is normally driven by a deep model pre-trained on a large-scale private image-set, while at inference it needs to recognize consumer images at a finer granularity. Such a model is expected to evolve on-the-fly [1] over time as being used, because fine-tuning (FT) it for specific novel classes induces an increasing number of separate models retrained, and thus is inefficient. This expectation is also generally valid for vision-driven autonomous systems or intelligent agents. For example, a self-driving car needs to gradually grow its perception capabilities as it runs on the road.

As shown in Fig. 1, we are interested in such a coarse-to-fine recognition problem that fits the class-incremental learning (CIL) setting. Moreover, fine classes appear asynchronously, which again fits CIL. It is also a few-shot learning problem, as there is no time to collect abundant samples per new class. We name such an incremental few-shot learning problem Coarse-to-Fine Few-Shot Class-Incremental Learning (C2FSCIL), and aim to propose a method that can evolve a generic model to both avoid catastrophic forgetting of source-blind coarse classes (see also Fig. 1) and prevent over-fitting the new few-shot fine-grained classes. However, **what exactly is the knowledge?** Incremental learning (IL) is aimed for the learning model to adapt to new data without forgetting its existing knowledge, which is also called catastrophic forgetting, a concept in connectionist networks [3, 4] – it occurs when the new weight vector is completely inappropriate as a solution for the originally learned pattern. In deep learning (DL), knowledge distillation (KD) is one of the most effective approaches to IL, while there lacks a consensus about what exactly the knowledge is in deep networks. Will it similarly be the weight vectors?

**Is a coarsely-learned embedding space generalizable?** We aim to achieve a superior performance at both the coarse and fine granularity. However, considering the diversity of fine labels, it is infeasible to train a comprehensive fine-grained model beforehand. Instead, can a model be trained, using coarsely-labeled samples, to classify finely-labeled samples with accuracy comparable to that of a

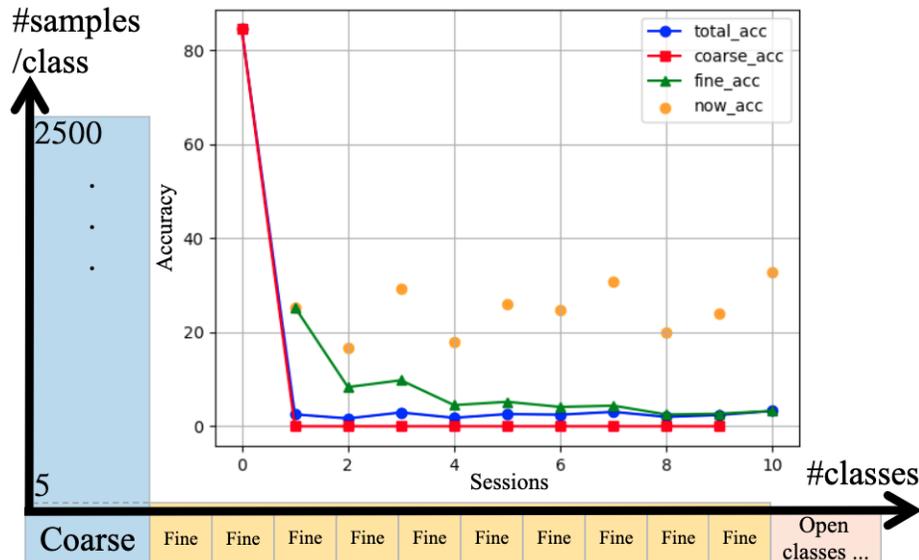


Figure 1: Catastrophic forgetting when FT-ing a coarsely-trained model on fine samples presently available w/o freezing any weight. We pre-set 10 sessions from CIFAR-100 [2]. There is a fine-class accuracy from the 1st session and yet no coarse-class accuracy as all samples are with fine labels.

model trained with fine labels [5]? Our hypothesis is yes; then, the next question is how to pre-train a generalizable base model? How to explore a finer embedding space from coarse labels? Namely, what type of knowledge is useful for fine classes and how can we learn and preserve them [6]?

**Can we balance old knowledge and current learning?** (*a.k.a.*, solving the stability–plasticity dilemma [7, 8]). We aim to remember cues of both the pre-trained base classes and fine classes in the previous few-shot sessions. Our hypothesis is yes and our preference is a linear classifier as it is flexible, data in-demanding, and efficient to train as well as simple for derivation. Furthermore, the next question is how a linear classifier can evolve the model effectively with a few shots and yet a balanced performance. As presumed, if the knowledge is weights, then freezing weights retains knowledge while updating weights is evolving the model.

To answer those questions, we will first fine-tune a coarse model to test our hypothesis. Motivated to make CIL as simple as fine-tuning, our contributions are four-fold.

1. We propose a new problem and empirical insights for incrementally learning coarse-to-fine with a few shots.
2. We propose to **learn, normalize, and freeze** weights, a simple process (Knowe, pronounced as 'now') that can effectively solve the problem once we have a base model contrastively learned from coarse labels.
3. We theoretically analyze why Knowe is a valid solver.
4. We propose a way to measure balanced performance.

## 2 Related Problems

### 2.1 Catastrophic Forgetting (CF)

To learn over time (*i.e.*, sequential learning), it is suggested in [9, 3] that neural networks can be limited by catastrophic forgetting (CF) just like Perceptron is unable to solve X-OR. Knowledge forgetting, or called catastrophic forgetting/interference is about a learner’s memory (*e.g.*, LSTM) and is a result of the stability–plasticity dilemma regarding how to design a model that is sensitive to, but not radically disrupted by, new input [9, 3]. Often, maintaining plasticity results in forgetting old classes while maintaining stability prevents the model from learning new classes, which may be caused by a single set of shared weights [3].

## 2.2 Weakly-Supervised Learning

Judging from the fine-class stage (Fig. 1 middle to right), if we combine a pre-training set and the support set as a holistic training set, then the few-shot fine-grained recognition using a model pre-trained on coarse samples are similar to the *weakly-supervised learning* and specifically *learning from coarse labels* [10, 5, 11, 12], *e.g.*, C2FS [10]. Ristel *et. al.* investigates how coarse labels can be used to recognize sub-categories using random forests [13] (say, NCM [14]).

## 2.3 Open-World Recognition

Judging from the coarse-class stage [15] (see the left side of Fig. 1), CIL [16] can be dated back to the support vector machine [17] and random forest [14, 13], where a new class can be added as a new node, and now seen as a progressive case of *continual/lifelong learning* [18, 1], where CF is a challenge as data are hidden. The topology structure is also favored in DL [19, 20]. *Few-shot learning* (FSL) measures models’ ability to quickly adapt to new tasks [21] and has a flavor of CIL considering novel classes in the support set, *e.g.*, DFSL [22], IFSL [23], FSCIL [19, 24, 25], and so on.

## 2.4 Uniqueness of Proposed Problem

Different from existing settings [24, 19, 25] that focus on remembering the pre-trained base classes only, our setting requires remembering the knowledge gained in both the base coarse and previous fine sessions. We add finer classes instead of new classes at the same granularity. Our setting requires a balance between coarse and fine performance unexplored by existing works, as shown in Fig. 2.

Method	Class hierarchy	Few-shot Learning	Incremental Learning
LwF [22]			✓
CEC[25]		✓	✓
ANCOR [10]	✓	✓	
IIRC[26]	✓		✓
C2FSCIL (Ours)	✓	✓	✓

Table 1: Comparison of settings with related works.

## 3 State of the Art (SOTA)

### 3.1 Incremental Learning (IL)

IL allows a model to be continually updated on new data without forgetting, instead of training a model once on all data. There are two settings: class-IL [27] and task-IL [18].

They share main approaches, such as regularization and rehearsal methods. Regularization methods prevent the drift of consolidated weights and optimize network parameters for the current task, *e.g.*, parameter control in EWC [4]. CIL is our focus and aims at learning a classifier that maintains a good performance on all classes seen in different sessions.

In addition, Li *et.al.* first introduces KD [28] to IL literature in LwF [29] by modifying a cross-entropy loss to retain the knowledge in the original model. Recent works focus on retaining old-class samples to compute the KD loss. For example, iCaRL [16] learns both features and strong classifiers by combining KD and feature learning, *e.g.*, NME.

### 3.2 Operating Weights for IL

The IL literature since 2017 has seen various weight operations (op. for short) in the sense of consolidation (*e.g.*, EWC [4]), aligning [30, 31], normalization [30, 32], standardization [33], regularization [4, 34], aggregation [35], calibration [36], rectification [37], transfer [38, 39], sharing [40], masking

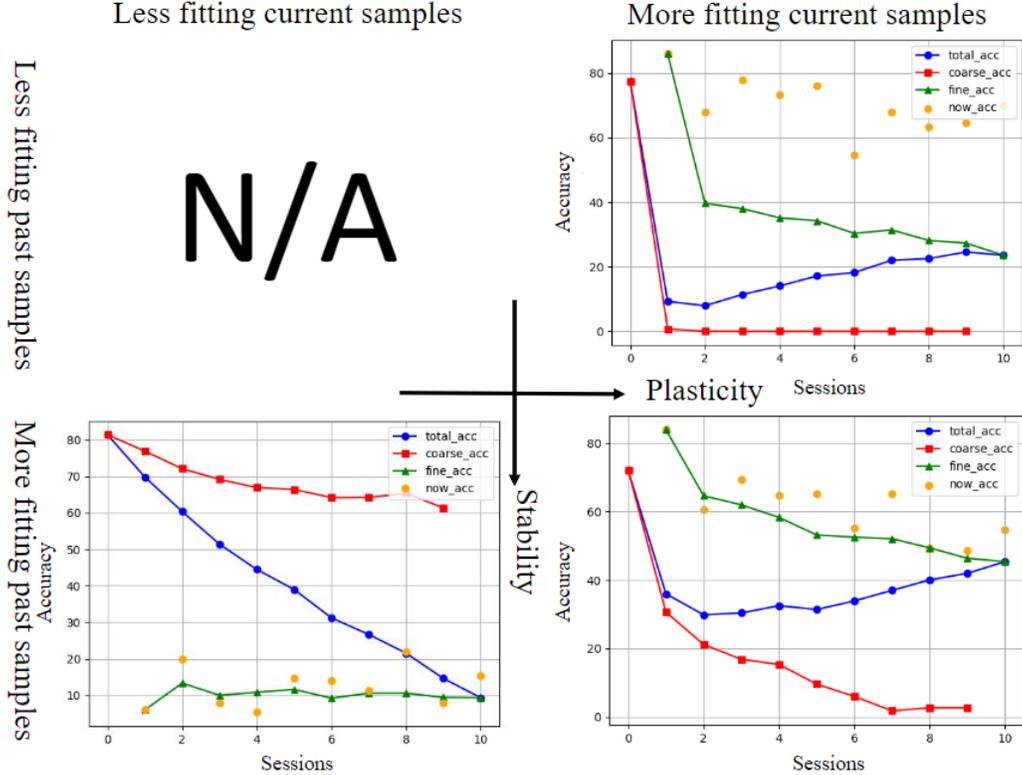


Figure 2: The stability-plasticity trade-off. Top-right is FT w/o IL; bottom-left represents most IFSL methods; bottom-right is our approach; top-left does not apply. (CIFAR-100)

[41], imprinting [42], picking [43], scaling [44], merging [45], pruning [46], quantization [47], weight importance [48], assignment [49], restricting weights to be positive [30], constraining weight changes [50], and so on.

### 3.3 Few-Shot Learning (FSL)

The prosperity of DL has pushed large-scale supervised learning, so far, to be the most popular learning paradigm. However, FSL is human-like learning [51] in the case of only a few samples [52, 53]. For example, Finn *et al.* proposed Model-Agnostic Meta Learning (MAML) to train a model that can quickly adapt to a new task using only a few samples and training iterations [54]. Prototypical Network learns a metric space in which classification can be performed by computing distances to prototype representations of each class [55]. Ren *et al.* proposes a meta-learning model, the Attention Attractor Network, which regularizes the learning of novel classes [23]. It is shown that decoupling the embedding learner and classifier is feasible [25]. Tian *et al.* demonstrates that using a good learned embedding model can be more effective than meta learning [21].

### 3.4 Incremental Few-Shot Learning (IFSL)

In the IFSL [23] or similarly FSCIL [19] setting, samples in the incremental session are relatively scarce, different from conventional CIL. While IFSL is based on meta learning, IFSL and DFSL [22] both utilize attentions. In FSCIL, a model named TOPIC is proposed, which contains a single neural gas (NG) network to learn feature-space topologies as knowledge, and adjust NG to preserve the stabilization and enhance the adaptation. In [24], Dong *et al.* propose an exemplar relation KD-IL framework to balance the tasks of old-knowledge preserving and new-knowledge adaptation as done in [56]. CEC [25] is proposed to separate classifier from the embedding learner, and use a graph attention network to propagate context cues between classifiers for adaptation. In [57], Hou *et al.* address the imbalance between old and new classes by cosine normalization [58, 22, 57].

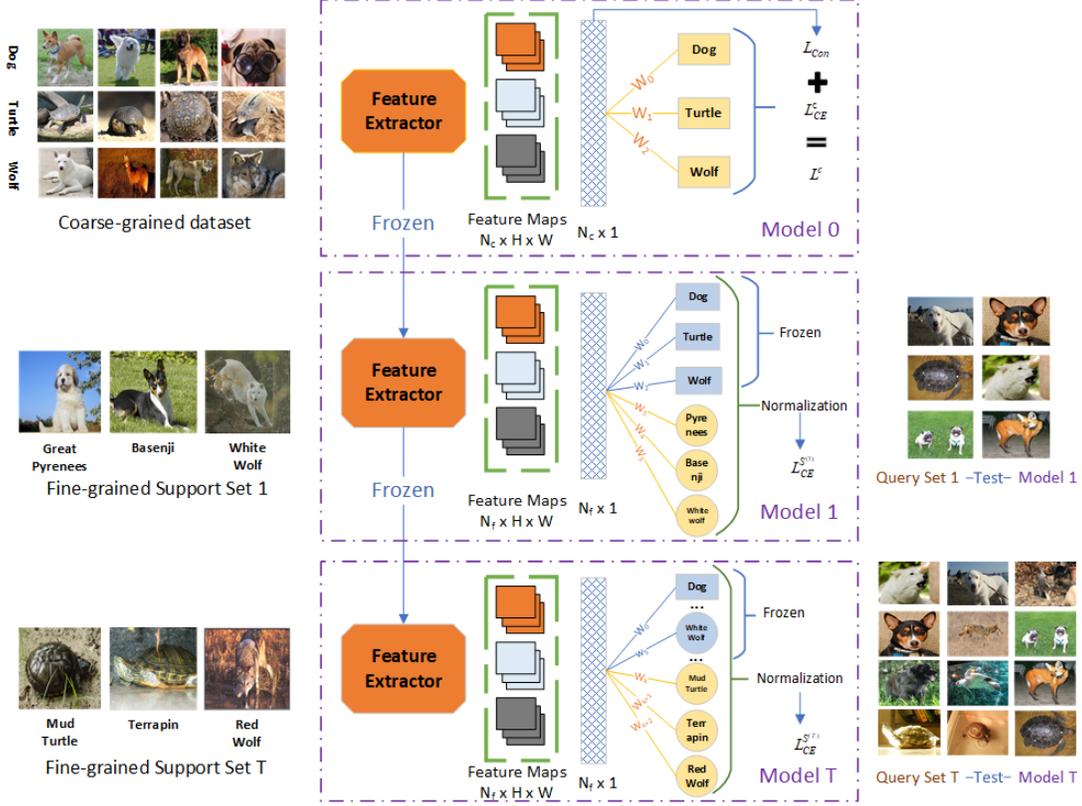


Figure 3: C2FSCIL and basic idea. In base session we train  $\Theta$  on  $\mathbb{D}$  to get  $\Theta^{(0)}$ . Per incremental session,  $\Theta^{(t)}$  is trained on  $C$ -way  $K$ -shot support set  $\mathbb{S}^{(t)}$  based on  $\Theta^{(t-1)}$ ,  $t \geq 1$  and then tested on any class seen in either  $\mathbb{D}$  or  $\mathbb{S}^{(1)}, \dots, \mathbb{S}^{(t)}$ .

### 3.5 Uniqueness of Proposed Approach

Different from state-of-the-art approaches to IFSL, we do not follow rehearsal methods, namely, our model learns without memorizing samples [59]. However, retaining samples is often practically infeasible, say, when learning on-the-fly [1]. Even if there is memory for storing previous samples, there often is a budget, buffer, or queue. Thus, we aim to examine the extreme case of knowledge forgetting, and thus design IFSL methods to the upper-bound extent. For example, although in [50] they do not use any base-class training samples and keep the weights of the base classifier frozen, they still use previous samples in their third phase.

## 4 A New Problem C2FSCIL and Our Insights

Given a model parameterized by  $\Theta$  and pre-trained on  $\mathbb{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $y_i \in \mathbb{Y} = \{\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_R\}$ , a set of  $R$  coarse labels  $\mathcal{Y}$ , we have a stream of  $C$ -way  $K$ -shot support sets  $\mathbb{S}^{(1)}, \mathbb{S}^{(2)}, \dots, \mathbb{S}^{(t)}, \dots, \mathbb{S}^{(T)}$  where  $\mathbb{S}^{(t)} = \{(\mathbf{x}_j^{(t)}, y_j^{(t)})\}_{j=1}^{C \cdot K}$  and  $y_j^{(t)} \in \mathbb{Z}^{(t)} = \{\mathcal{Z}_1^{(t)}, \dots, \mathcal{Z}_C^{(t)}\}$ , a set of  $C$  fine-grained labels  $\mathcal{Z}$ . Then, we adapt our model to  $\mathbb{S}^{(1)}, \mathbb{S}^{(2)}, \dots, \mathbb{S}^{(t)}$  over time and update the parameter set  $\Theta$  from  $\Theta^{(0)}$  all the way to  $\Theta^{(t)}$ , as shown levelwise in Fig. 3.

For testing, we also have a stream of  $(C \cdot t + R)$ -way  $H$ -shot query sets  $\mathbb{Q}^{(1)}, \mathbb{Q}^{(2)}, \dots, \mathbb{Q}^{(t)}, \dots, \mathbb{Q}^{(T)}$  where  $\mathbb{Q}^{(t)} = \{(\mathbf{x}_k^{(t)}, y_k^{(t)})\}_{k=1}^{(C \cdot t + R)H}$  and  $y_k^{(t)} \in \cup_{l=1}^t \mathbb{Z}^{(l)} \cup \mathbb{Y}$ , which is the generalized union of all label sets till the  $t$ -th session.

Notably,  $\mathbb{Z}^{(t_1)} \cap \mathbb{Z}^{(t_2)} = \emptyset, \forall t_1, t_2$ . We assume no sample can be retained (unlike rehearsal methods) and the CIL stage only includes (sub-classes of) base classes. At the  $t$ -th session, only the support set  $\mathbb{S}^{(t)}$  can be used for training.

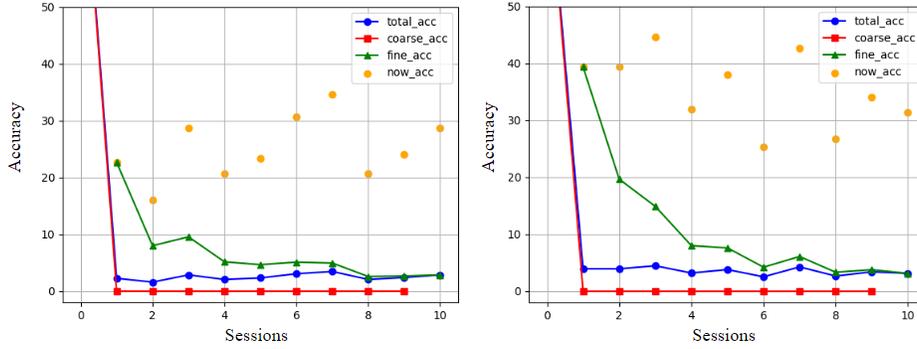


Figure 4: Ablation study of contrastive learning when fine-tuning ResNet12 w/o IL. Left: w/o; right: w/. (CIFAR-100)

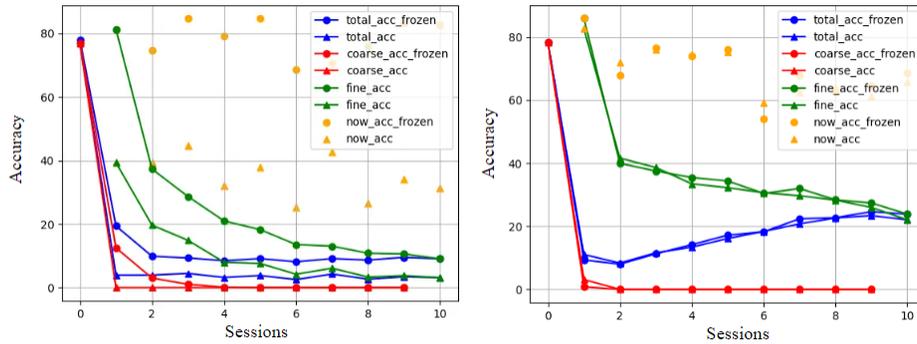


Figure 5: Ablation study of freezing embedding-weights for fine-tuning a contrastive model. Left: when not freezing classifier-weights. Right: when freezing them. (CIFAR-100)

#### 4.1 Embeddings need to be contrastively learned

As shown in Fig. 4, straightforward training on coarse labels does not help much the subsequent FSL on fine labels (now\_acc at  $\sim 25\%$ ), while contrastive learning self-supervised by the fine cues does help (now\_acc at  $\sim 35\%$ ). Thus, **coarsely-trained embedding can be generalizable**.

Fig. 5-left shows that freezing embedding-weight outperforms not freezing them. It implies the embedding space without any update is generalizable, and that, *if classifier-weights are not frozen, freezing embedding-weights helps*.

#### 4.2 Freezing weights helps, surely for classifiers

However, Fig. 5-right implies that, *if classifiers weights are frozen, then freezing embedding-weights does not help*.

Comparing left with right of Fig. 5, we find that freezing classifier-weights (right) outperforms not doing so (left), either freezing embedding-weights (circle) or not (triangle).

#### 4.3 Weights need to be normalized

As shown in Fig. 6, samples of classes seen in the 1st session are totally classified to classes seen in the 2nd session while only samples of the present classes can be correctly classified. We plot weight norms to find them grow and propose a conjecture implying a need of normalization. Please see our analysis in the Appendix.

**Conjecture 1** (FC weights grow over time). Let  $\|\mathbf{W}^{(t)}\|_F$  denotes the Frobenius norm of the weight matrix formed by all weight vectors in the FC layer for new classes in the  $t$ -th session. With training converged and norm outliers ignored, it holds that  $\|\mathbf{W}^{(t)}\|_F > \|\mathbf{W}^{(t-1)}\|_F, \forall t \in \{1, \dots, T\}$ .

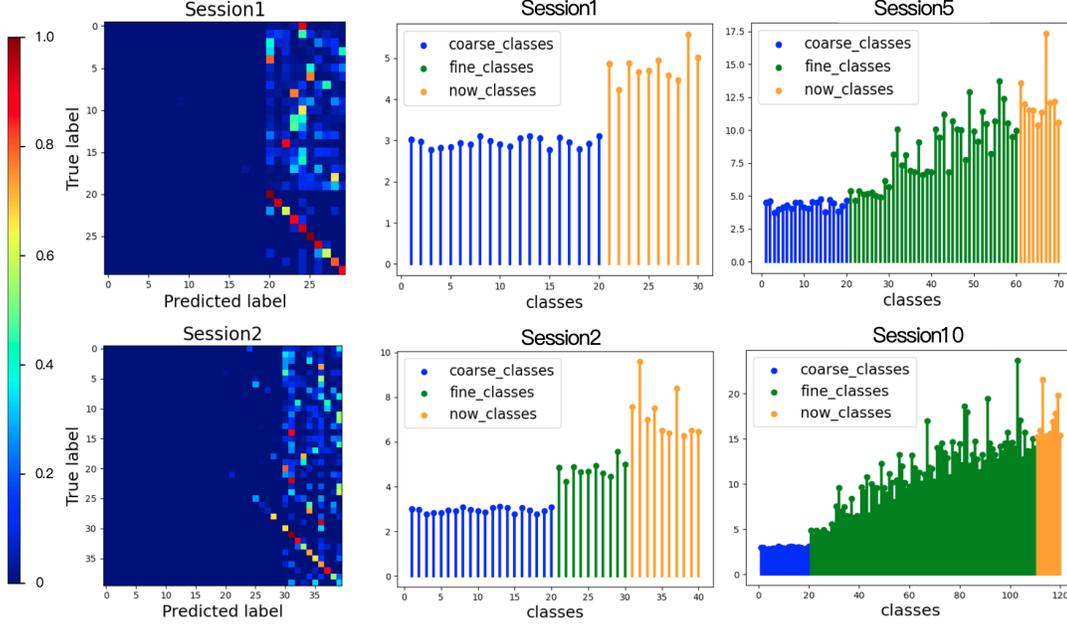


Figure 6: 10-way 5-shot confusion matrix (left) and visualization of the norm of raw weights (mid-right) in the last layer for old/new classes. As each session can only access labels of the present classes, a linear classifier will have a larger weight for the current classes’ neurons, inducing the queries of previous classes to be likely assigned into current classes’ region (left) in the embedding space. (CIFAR-100)

#### 4.4 The balance need to be measured

As already shown in the bottom-right sub-plot of Fig. 2, it is possible to avoid the collapse of coarse-class accuracy, and slow down the accuracy drop of all the previous classes, while still maintaining a high accuracy on present classes.

**Old knowledge and current learning can be balanced**, which can be achieved not only on CIFAR-100 but also more generally on BREEDS [60]. Fig. 7 shows the balanced performance on its various subsets. In order to better measure how good the balance is, we also need new overall metrics.

## 5 A New Approach: Know-weight (Knowe)

### 5.1 Learning Embedding-Weights Contrastively

Now, we elaborate on the Model 0 of Fig. 3 about how we train a generalizable base embedding space [21, 61].

We follow ANCOR [10] to use MoCo [62] as the backbone, and keep two network streams each of which contains a backbone with the last-layer FC replaced by a Multi-Layer Perceptron (MLP). The hidden layer of two streams’ MLP outputs intermediate  $\mathbf{q}$  and  $\mathbf{k}$ , respectively. Given coarse labels, the total loss is defined as  $\mathcal{L}^c = \mathcal{L}_{Con} + \mathcal{L}_{CE}^c$  where

$$\mathcal{L}_{Con} = - \sum_{n=1}^N \log \frac{\exp(\mathbf{q}_n^T \mathbf{k}_n^+ / \tau)}{\exp(\mathbf{q}_n^T \mathbf{k}_n^+ / \tau) + \sum_{m \neq n} \exp(\mathbf{q}_n^T \mathbf{k}_m^- / \tau)}, \quad (1)$$

and  $\mathcal{L}_{CE}^c$  is the standard cross-entropy loss that captures the inter-class cues. We also use angular normalization [10] to improve their synergy. Note that  $m, n$  index samples,  $\tau$  is a temperature parameter,  $\mathbf{k}_m^-$  denotes the intermediate output of the  $m$ -th sample, a negative sample, in the same class with the  $n$ -th sample, a positive sample, so as to capture intra-class cues (fine cues), and reduce unnecessary noises to the subsequent fine-grained classification [11].  $\mathcal{L}_{Con}$  will be small when  $\mathbf{q}_n$  is similar with  $\mathbf{k}_n^+$  and different from  $\mathbf{k}_m^-$ .

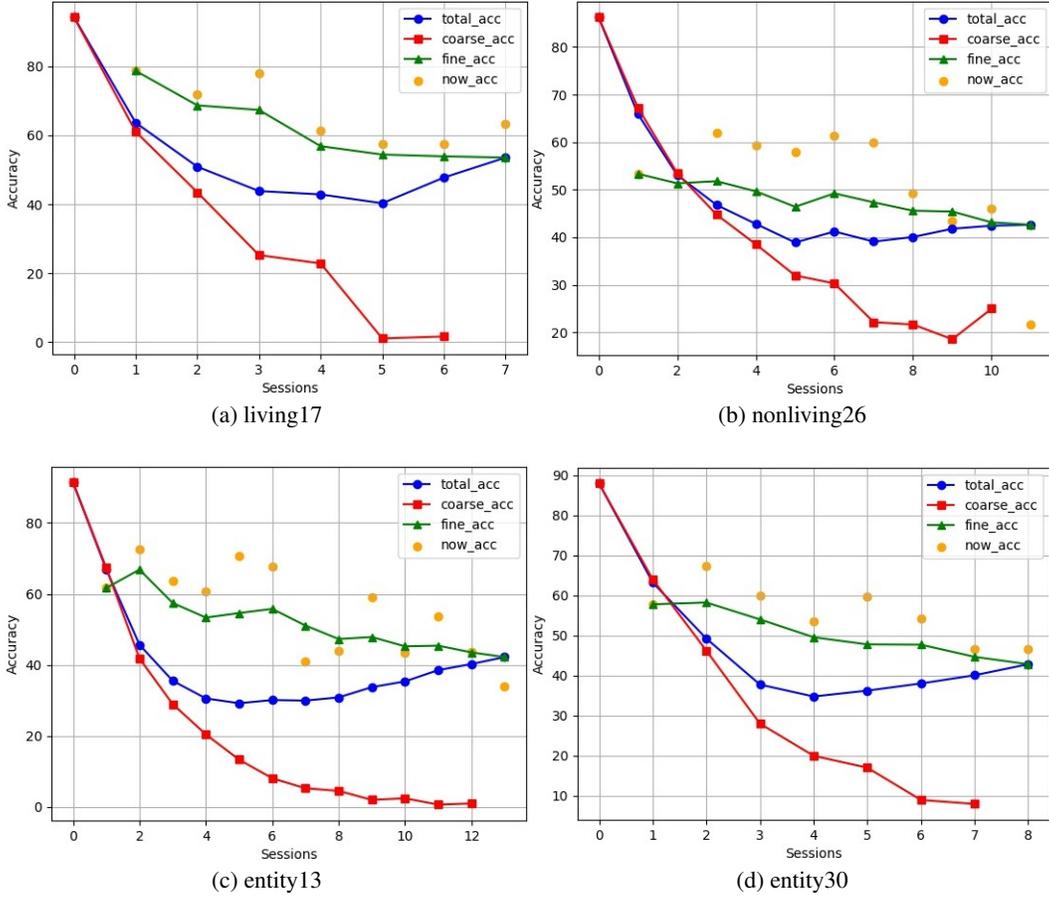


Figure 7: Reaching a balance on BREEDS. More in Sec. 6.

## 5.2 Normalizing Classifier-Weights

In the last layer, we set the bias term to 0. For a sample  $\mathbf{x}$ , once a neuron has its output logit  $o = \mathbf{w}^T \mathbf{f}(\mathbf{x})$  ready, then a Softmax activation function  $Sm.x(\cdot)$  is applied to convert  $o$  to a probability so that we can classify  $\mathbf{x}$ . ( $T$  is transpose)

However, such an inner-product linear classification often favors new classes [57]. Instead, we compute the logit using the normalized inner-product [58] (*a.k.a.*, cosine similarity, cosine normalization [22, 57]) as  $\tilde{o} = \tilde{\mathbf{w}}^T \tilde{\mathbf{f}}(\mathbf{x})$  where  $\mathcal{L}_2$ -normalized  $\tilde{\mathbf{f}}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) / \|\mathbf{f}(\mathbf{x})\|_2$  and  $\tilde{\mathbf{w}}_i = \mathbf{w}_i / \|\mathbf{w}_i\|_2$ , and then apply Softmax to the rescaled logit  $\tilde{o}$  as

$$p_i(\mathbf{x}) = Sm.x(\tilde{o}/\lambda) = \frac{\exp(\tilde{\mathbf{w}}_i^T \tilde{\mathbf{f}}(\mathbf{x})/\lambda)}{\sum_j \exp(\tilde{\mathbf{w}}_j^T \tilde{\mathbf{f}}(\mathbf{x})/\lambda)} \quad (2)$$

where  $i$  is the class index,  $\lambda$  is a temperature parameter that rescales the Softmax distribution, as  $\tilde{o}$  is ranged of  $[-1, 1]$ .

## 5.3 Freezing Memorized Classifier-Weights

As shown in Fig. 3, in the  $t$ -th incremental session, the task is similar to FSL where a support set  $\mathbb{S}^{(t)}$  is offered to train a model to be evaluated on a query set  $\mathbb{Q}^{(t)}$ . However, FSL only evaluates the classification accuracy of the classes appeared in the support set  $\mathbb{S}^{(t)}$ . In our setting, the query set  $\mathbb{Q}^{(t)}$  contains base classes, and all classes in previous support sets. As shown in Fig. 5, no matter freezing

embedding-weights helps or not, it does not hurt. We do so, hoping it to reduce model complexity to avoid over-fitting.

As past samples are not retained, we store the classifier-weights per session to implicitly retain the label information by augmenting a weight matrix  $\mathbf{W}$  where in the  $t$ -th session, we have  $\mathbf{W}_{[B:E]} = [\mathbf{w}_1^{(t)} | \mathbf{w}_2^{(t)} | \dots | \mathbf{w}_C^{(t)}]_{d \times C}$  with  $B = R + C \cdot (t - 1) + 1$ ,  $E = R + C \cdot t$  for  $t \geq 1$ , except  $\mathbf{W}_{[:R]} = [\mathbf{w}_1^{(0)} | \mathbf{w}_2^{(0)} | \dots | \mathbf{w}_R^{(0)}]_{d \times R}$  where  $d$  is the feature dimension.

In the  $t$ -th session, we minimize the following regularized cross-entropy loss on the support set  $\mathbb{S}^{(t)}$ :

$$\mathcal{L}_{CE}^{\mathbb{S}^{(t)}} = -\frac{1}{C \cdot K} \sum_{n=1}^{C \cdot K} \sum_{i=1}^{R+t \cdot C} \delta_{y_n^{(t)}=i} \log[p_i(\mathbf{x}_n^{(t)})] \quad (3)$$

where  $\delta_{y_n^{(t)}=i}$  is the indicator function and  $p_i(x_n^{(t)})$  is the output probability (*i.e.*, Softmax of logits) of the  $i$ -th class.

#### 5.4 Theoretical Guarantee for Stability-Plasticity

We extend definitions in [8] to set the base of our analysis. Please see also proofs in the Appendix.

**Definition 1** (Stability Decay). *For the same input sample, let  $\hat{\mathbf{o}}_i^{(t)}$  denote the output logits of the  $i$ -th neuron in the last layer in the  $t$ -th session. After the loss reaches the minimum, we define the degree of stability as  $\mathcal{D} = \sum_i (\frac{\hat{\mathbf{o}}_i^{(T)} - \hat{\mathbf{o}}_i^{(t)}}{\hat{\mathbf{o}}_i^{(t)}})^2$ .*

**Definition 2** (Relative Stability). *Given models  $\Theta_a$  and  $\Theta_b$ , if  $0 \leq \mathcal{D}_a < \mathcal{D}_b$ , then we say  $\Theta_a$  is more stable than  $\Theta_b$ .*

Assuming embedding-weights are frozen, then we have:

**Proposition 1** (Normalizing or freezing weights improves stability; doing both improves the most). *Given  $\Theta_a$ , if we only normalize weights of a linear FC classifier, we obtain  $\Theta_b$ ; if we only freeze them, we obtain  $\Theta_c$ ; if we do both, we obtain  $\Theta_d$ . Then,  $\mathcal{D}_d < \mathcal{D}_b < \mathcal{D}_a$  and  $\mathcal{D}_d < \mathcal{D}_c < \mathcal{D}_a$ .*

Our second claim is about normalization for plasticity.

**Proposition 2** (Weights normalized, plasticity remains). *To train our FC classifier, if we denote the loss as  $\mathcal{L}(\mathbf{w})$  where  $\mathbf{w}$  is normalized, the weight update at each step as  $\Delta \mathbf{w}$ , and the learning rate as  $\alpha$ , then we have  $\mathcal{L}(\mathbf{w} - \alpha \Delta \mathbf{w}) < \mathcal{L}(\mathbf{w})$ .*

Notably, freezing the weights does not affect plasticity.

#### 5.5 New Overall Performance Measures

In this section, we evaluate the model after each session with the query set  $\mathbb{Q}^{(t)}$ , and report the Top-1 accuracy. The base session only contains coarse labels, and thus is evaluated by the coarse-grained classification accuracy  $\mathcal{A}_c$ . We evaluate  $\mathcal{A}_c$ , the fine-grained accuracy  $\mathcal{A}_f$ , and the total accuracy  $\mathcal{A}_t$  per incremental session, except the last session when only fine labels are available and  $\mathcal{A}_c$  is not evaluated. We average  $\mathcal{A}_t$  to obtain an overall performance score as

$$\bar{\mathcal{A}} = \frac{1}{T+1} \sum_{i=0}^T \mathcal{A}_t^i. \quad (4)$$

Inspired by [44], we define the fine-class forgetting rate

$$\mathcal{F}_f^t = \frac{\mathcal{A}_f^{t-1} - \mathcal{A}_f^t}{\mathcal{A}_f^{t-1}}, \quad (5)$$

and the forgetting rate for the base coarse class as

$$\mathcal{F}_c^t = \frac{\mathcal{A}_c^0 - \mathcal{A}_c^t}{\mathcal{A}_c^0}. \quad (6)$$

With them, we can evaluate the model with an overall measure to represent the catastrophic forgetting rate as

$$\mathcal{F} = \frac{1}{T-1} \left( \sum_{t=2}^T \mathcal{F}_f^t * \frac{c_t}{N_f} + \sum_{t=1}^{T-1} \mathcal{F}_c^t * \left(1 - \frac{c_t}{N_f}\right) \right) \quad (7)$$

where  $T$  is the number of incremental sessions;  $c_t$  is the number of appeared fine classes until the  $t$ -th session, and  $N_f$  is fine-class total number;  $\mathcal{A}_c$  and  $\mathcal{A}_f$  are the accuracy of coarse and fine classes per session, respectively.

## 6 Experiments

### 6.1 Datasets and Results

**CIFAR-100** contains 60,000 32x32 images from 100 fine classes, each of which has 500 training images and 100 test images [2]. They can be grouped into 20 coarse classes, each of which includes 5 fine classes, e.g., *trees* contains *maple*, *oak*, *pine*, *palm*, and *willow*. The 100 fine classes are divided into 10 10-way 5-shot incremental sessions.

**BREEDS** is derived from ImageNet with class hierarchy re-calibrated by [60] and contains 4 subsets named living17, nonliving26, entity13, and entity30. They have 17, 26, 13, 30 coarse classes, 4, 4, 20, 8 fine classes per coarse class, 88K, 132K, 334K, 307K training images (224x224), 3.4K, 5.2K, 13K, 12K test images, respectively. See also Table 2.

**tieredImageNet** (tIN) is a subset of ImageNet and contains 608 classes [63] that are grouped into 34 high-level super-classes to ensure that the training classes are distinct enough from the test classes semantically. The train/val/test set have 20, 6, 8 coarse classes, 351, 97, 160 fine classes, 448K, 124K, 206K images (sized at 84x84), respectively.

Table 2 summarizes our performance. Fig. 7 shows our separated accuracy on BREEDS, and Fig. 8 visualizes confusion matrices to show the evolving of per-class accuracy.

Dataset	coarse#	fine#	total#	sessions	way/shot	queries	$\bar{\mathcal{A}}$	$\mathcal{F}$
CIFAR-100	20	100	120	10	10/5	15	38.50	0.42
living17	17	68	85	7	10/1	15	54.62	0.33
nonliving26	26	104	130	11	10/1	15	48.41	0.25
entity13	13	260	273	13	20/1	15	41.45	0.38
entity30	30	240	270	8	30/1	15	47.79	0.32
tieredImageNet	20	351	371	10	36/5	15	33.24	0.39

Table 2: Dataset setting and performance. # is class num.

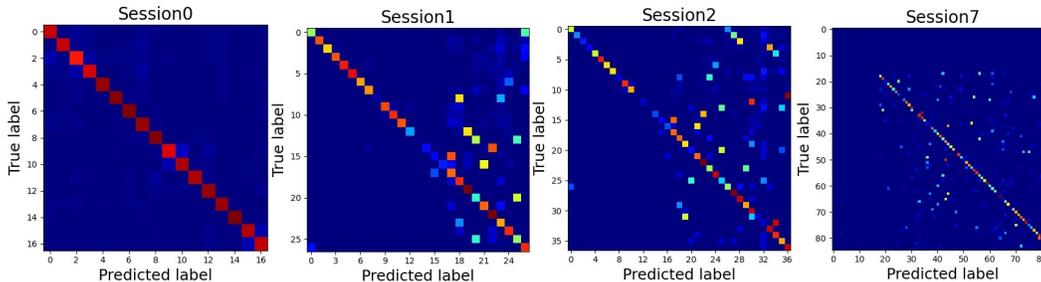


Figure 8: Confusion matrices of Knowe tested on living17.

### 6.2 Implementation Details

We use ResNet-50 on BREEDS, ‘-12’ on CIFAR-100 and ‘-12’ on tIN, train  $\Theta^{(0)}$  except FC using ANCOR, use SGD with a momentum 0.9, as well as set weight decay to  $5e-4$ , batch size to 256,  $\tau$  to 0.2, and  $\lambda$  to 0.5. The learning rate is 0.12 for  $\Theta^{(0)}$ , and is 0.1 for  $\Theta^{(1)}, \Theta^{(2)}$ , etc. for 200 epochs.

### 6.3 Ablation Study

**Impact of base contrastive learning.** Fig. 4 already illustrates its benefit for a simple model without any weight operation. As shown in Fig 9a, Knowe also obtains a better performance than not using MoCo in Knowe’s base, which verifies that the contrastively-learned base model helps fine-grained recognition. Starting from almost the same fine accuracy in the 2nd session, the gap between w/ MoCo and w/o MoCo increases, as the former stably outperforms the latter on current classes. It verifies that the former can learn more fine knowledge than the latter. Given that there are only a few fine-class samples, the extra fine-grained knowledge is likely from the contrastively-learned base model.

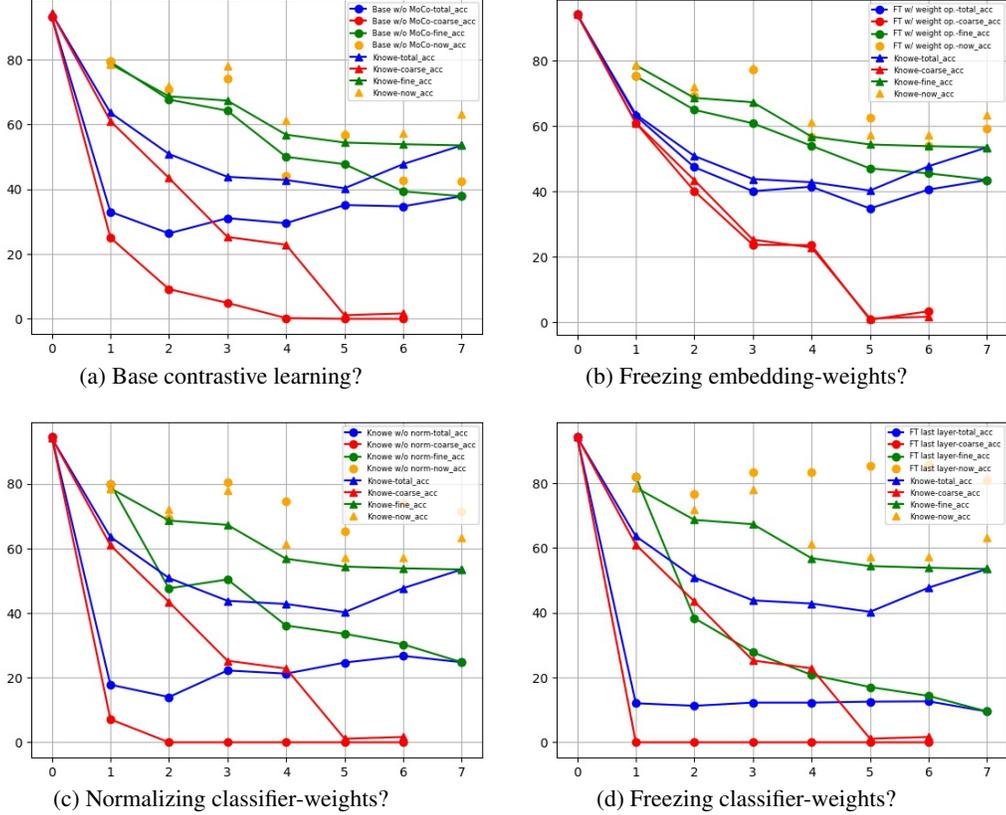


Figure 9: 4-factor ablation study on living17, separated acc.

**Impact of freezing embedding-learner weight (decoupling).** It has been verified in Sec.4.1 that, if neither freezing nor normalizing classifier-weights, freezing the embedding-weights helps. We have a **Conjecture 2**:  $\neg p \wedge \neg q \Rightarrow r$  where  $p$  is a premise that classifier-weights are normalized,  $q$  is another that classifier-weights are frozen, and  $r$  is a conclusion that freezing embedding-weights improves Knowe’s performance. However, Fig. 9b illustrates that, for Knowe, freezing embedding-weights induces a slightly better performance than not freezing them. It implies, if classifier-weights are normalized and frozen, then freezing the embedding-weights does not help ( $p \wedge q \Rightarrow \neg r$ ), which is shown by small changes of  $\bar{\mathcal{A}}$  and  $\bar{\mathcal{F}}$  in Table 3.

**Impact of normalizing classifier-weight.** Fig 6 has already shown that, with a linear classifier, the weight norms of new classes totally surpass the weight norms of previous classes, which causes that the linear classifier biases towards new classes (i.e., any sample of previous class can be classified as a new class). That implies a need of normalizing the classifier-weights. As shown in Fig 9c, when we freeze weights of previous classes and only tune the weights of new classes without normalization, the model performs stably worse than Knowe with normalization, which verifies that normalizing classifier-weights plays a positive role.

**Impact of freezing memorized classifier-weights.** As shown in Fig 9d, there is severe CF of both fine and coarse knowledge when not freezing the weights of previously-seen classes, which implies that little knowledge is retained. Although embedding-weights are frozen and classifier-weights are normalized, the coarse knowledge is totally forgotten. It implies that, *if classifier-weights are normalized and yet not frozen, freezing the embedding-weights does not help* ( $p \wedge \neg q \Rightarrow \neg r$ ). It can be explained that fine-tuning on a few samples normally induces little change to the embedding-weights and yet great change to classifier-weights. Moreover, the model without freezing classifier-weights performs much worse than Knowe that freezes previous weights. The gap of the fine accuracy increases over time and is larger than the gap of the present accuracy. It implies that they also differ in the performance of previous fine classes, which is the CF of learned fine knowledge.

**More about freezing embedding-weights.** In Sec.4.2, we know  $\neg p \wedge q \Rightarrow \neg r$ . Thus, we have a **Conjecture 3:**  $p \vee q \Leftrightarrow \neg r$ , meaning *if and only if classifier-weights are either normalized or frozen, then freezing embedding-weights does not help*. Please see also the Appendix for the analysis.

**Overall finding.** A decent *now\_acc* seems to be a condition for weight freezing and normalization to be effective.

Method	Contr. learn.	Decoupled	Normalization	Frozen	0	1	2	Total accuracy per session								$\bar{A} \uparrow$	$\mathcal{F} \downarrow$
								3	4	5	6	7	8				
(a) Base w/o MoCo					93.18	33.04	26.37	31.08	29.51	35.10	34.71	37.84	N/A	40.10	0.50		
(b) FT w/ weight op.	✓				94.21	63.14	47.45	40.10	41.47	34.80	40.59	43.53	N/A	50.66	0.35		
(c) Knowe w/o norm.	✓	✓			94.50	17.84	14.02	22.26	21.28	24.71	26.77	24.80	N/A	30.77	0.57		
(d) FT last layer	✓	✓			94.21	12.06	11.28	12.26	12.26	12.55	12.65	9.51	N/A	22.09	0.66		
LwF+[29]	✓				<b>94.50</b>	<i>61.47</i>	<i>44.61</i>	27.45	19.12	11.28	6.37	4.22	N/A	33.63	0.51		
ScalL[44]	✓				<b>94.50</b>	38.63	25.59	31.08	30.29	35.10	37.84	41.08	N/A	41.76	0.48		
Weight Align+ [30]	✓	✓			<b>94.50</b>	50.98	37.94	<i>38.43</i>	<i>37.06</i>	<i>35.20</i>	<i>39.80</i>	<i>43.24</i>	N/A	<i>47.14</i>	<i>0.40</i>		
Subsp. Reg.+ [64]	✓	✓			<b>94.50</b>	59.41	39.51	33.43	29.31	25.59	27.84	26.47	N/A	42.01	0.40		
Knowe (Ours)	✓	✓	✓		<i>94.21</i>	<b>63.63</b>	<b>50.88</b>	<b>43.82</b>	<b>42.84</b>	<b>40.29</b>	<b>47.75</b>	<b>53.53</b>	N/A	<b>54.62</b>	<b>0.33</b>		
ANCOR [10]	✓				94.50	11.86	11.18	12.35	11.77	12.55	10.78	9.02	N/A	21.75	0.66		
Jt. train. (upp. bd.)	✓	✓	✓		94.21	63.63	58.53	52.26	46.28	47.75	36.96	42.75	N/A	55.29	0.25		
LwF+[29]	✓				<b>89.48</b>	<b>65.03</b>	<i>48.69</i>	22.72	9.36	6.03	4.61	2.86	3.33	28.01	0.47		
ScalL[44]	✓				<b>89.48</b>	39.25	25.50	22.44	23.69	25.75	30.81	32.08	35.25	36.03	0.48		
Weight Align+ [30]	✓	✓			<b>89.48</b>	47.36	37.06	<i>31.72</i>	<i>30.56</i>	<i>32.28</i>	<i>34.11</i>	<i>36.39</i>	<i>37.06</i>	<i>41.78</i>	<i>0.42</i>		
Subsp. Reg.+ [64]	✓	✓			<b>89.48</b>	42.39	28.94	20.86	16.14	16.44	16.75	16.17	16.06	29.25	0.48		
Knowe (Ours)	✓	✓	✓		<i>87.90</i>	<i>63.22</i>	<b>49.22</b>	<b>37.75</b>	<b>34.78</b>	<b>36.25</b>	<b>38.03</b>	<b>40.08</b>	<b>42.83</b>	<b>47.79</b>	<b>0.32</b>		
ANCOR[10]	✓				89.48	8.67	8.28	9.50	6.83	8.75	9.53	8.19	8.69	17.55	0.61		
Jt. train. (upp. bd.)	✓	✓	✓		87.90	63.22	56.56	53.72	47.36	44.78	41.61	38.06	36.75	52.22	0.20		

Table 3: Ablation study of 4 factors and comparison with others on BREEDS living17 (top) and entity30 (bottom). Best seen on computer.

Method	Contr. learn.	Decoupled	Normalization	Frozen	0	1	2	3	4	5	6	7	8	9	10	11	12	13	$\bar{A} \uparrow$	$\mathcal{F} \downarrow$	
LwF+[29]	✓				<b>86.94</b>	65.57	<b>58.14</b>	<i>44.17</i>	22.76	14.36	9.68	6.92	5.90	5.19	5.32	3.40	N/A	N/A	27.36	0.38	
ScalL[44]	✓				<b>86.94</b>	36.09	24.10	21.47	23.27	23.65	27.95	31.80	34.23	36.09	<i>37.76</i>	<i>38.14</i>	N/A	N/A	35.12	0.43	
Weight Align+ [30]	✓	✓			<b>86.94</b>	61.41	46.03	40.00	35.77	<i>34.10</i>	<i>35.96</i>	<i>33.21</i>	<i>35.51</i>	<i>36.60</i>	<i>37.56</i>	<i>37.76</i>	N/A	N/A	<i>43.40</i>	<i>0.29</i>	
Subsp. Reg.+ [64]	✓	✓			<b>86.94</b>	63.59	52.56	42.95	35.96	31.41	28.01	26.15	23.27	19.68	19.36	20.19	N/A	N/A	37.51	<b>0.25</b>	
Knowe (Ours)	✓	✓	✓		<i>86.23</i>	<b>65.90</b>	<i>53.08</i>	<b>46.80</b>	<b>42.82</b>	<b>38.91</b>	<b>41.22</b>	<b>39.10</b>	<b>40.06</b>	<b>41.80</b>	<b>42.44</b>	<b>42.63</b>	N/A	N/A	<b>48.41</b>	<b>0.25</b>	
ANCOR [10]	✓				<b>86.94</b>	5.83	6.03	6.92	5.90	6.60	7.63	7.05	7.05	7.50	7.44	2.63	3.60	N/A	N/A	13.13	0.61
Jt. train. (upp. bd.)	✓	✓	✓		<b>86.23</b>	65.90	60.51	59.04	53.53	53.85	46.73	46.60	43.85	36.67	37.31	36.80	N/A	N/A	52.25	0.16	
LwF+[29]	✓				<b>92.03</b>	59.10	<i>43.64</i>	18.49	10.49	6.82	3.59	2.54	3.10	2.56	2.10	2.23	1.77	1.54	17.86	0.52	
ScalL[44]	✓				<b>92.03</b>	37.10	13.92	13.36	14.87	18.36	21.72	23.28	24.33	27.62	29.59	31.54	32.36	34.08	29.58	0.49	
Weight Align+ [30]	✓	✓			<b>92.03</b>	36.74	24.15	<i>20.51</i>	<i>22.31</i>	<i>24.82</i>	<i>26.41</i>	<i>26.85</i>	<i>27.26</i>	<i>31.49</i>	<i>32.26</i>	<i>35.28</i>	<i>36.72</i>	<i>37.69</i>	<i>33.89</i>	<i>0.46</i>	
Subsp. Reg.+ [64]	✓	✓			<b>92.03</b>	52.72	28.95	15.92	12.08	10.82	10.90	11.49	12.05	12.03	11.77	11.72	12.54	14.36	22.10	0.45	
Knowe (Ours)	✓	✓	✓		<i>91.35</i>	<b>66.90</b>	<b>45.69</b>	<b>35.54</b>	<b>30.56</b>	<b>29.21</b>	<b>30.10</b>	<b>29.95</b>	<b>30.85</b>	<b>33.74</b>	<b>35.36</b>	<b>38.54</b>	<b>40.26</b>	<b>42.21</b>	<b>41.45</b>	<b>0.38</b>	
ANCOR[10]	✓				92.03	5.36	5.67	5.49	5.18	6.51	5.82	4.80	5.39	6.28	5.36	5.13	5.26	5.62	11.71	0.57	
Jt. train. (upp. bd.)	✓	✓	✓		91.35	66.90	57.54	49.92	50.59	48.64	47.69	44.41	41.72	39.13	39.62	40.72	38.49	37.26	49.57	0.24	

Table 4: Comparison with others on BREEDS nonliving26 (top) and entity13 (bottom). **Bold** is the best, *slanted* is 2nd. Best seen on computer.

Method	Contr. learn.	Decoupled	Normalization	Frozen	0	1	2	3	4	5	6	7	8	9	10	$\bar{A} \uparrow$	$\mathcal{F} \downarrow$
LwF+[29]	✓				<b>78.39</b>	<b>41.87</b>	28.00	23.80	14.93	10.53	8.00	8.80	6.47	7.33	6.73	21.35	0.51
ScalL[44]	✓				<b>78.39</b>	14.47	14.13	18.07	21.00	25.20	26.20	31.87	32.60	36.53	38.20	30.61	0.52
Weight Align+ [30]	✓	✓			<b>78.39</b>	13.20	14.13	18.20	21.20	24.60	26.93	32.33	32.60	38.93	38.46	30.82	0.53
Subsp. Reg.+ [64]	✓	✓			<b>78.39</b>	<i>41.47</i>	<b>31.80</b>	<b>32.87</b>	26.73	25.73	25.27	26.73	24.27	25.73	24.00	<i>33.00</i>	<i>0.42</i>
Knowe (Ours)	✓	✓	✓		<i>72.07</i>	36.00	<i>28.13</i>	<i>30.27</i>	<b>32.20</b>	<b>31.20</b>	<b>30.93</b>	<b>36.33</b>	<b>39.27</b>	<b>43.20</b>	<b>43.93</b>	<b>38.50</b>	<b>0.42</b>
ANCOR[10]	✓				78.39	7.93	7.13	8.27	7.80	8.60	6.40	7.53	6.93	8.20	8.33	14.14	0.59
Jt. train. (upp. bd.)	✓	✓	✓		<i>72.07</i>	36.00	<i>37.07</i>	<i>40.27</i>	<i>40.13</i>	<i>41.33</i>	38.60	41.13	40.47	41.40	43.47	42.90	0.33
LwF+[29]	✓				<b>87.64</b>	<b>69.36</b>	13.88	4.22	4.05	4.03	3.02	2.74	1.44	1.05	1.06	17.50	0.55
ScalL[44]	✓				<b>87.64</b>	48.51	<b>33.12</b>	<b>26.15</b>	<b>22.66</b>	22.77	23.42	22.72	23.38	25.17	26.65	32.93	0.40
Weight Align+ [30]	✓	✓			<b>87.64</b>	25.13	18.63	18.37	20.08	22.20	24.22	<b>24.73</b>	26.71	29.00	<i>30.45</i>	29.74	0.48
Subsp. Reg.+ [64]	✓	✓			<b>87.64</b>	49.73	32.06	24.35	20.95	20.76	20.84	21.12	21.79	23.15	24.31	31.52	0.42
Knowe (Ours)	✓	✓	✓		<i>76.15</i>	48.24	30.60	25.60	22.34	<b>23.48</b>	<b>24.79</b>	<i>24.69</i>	<b>27.65</b>	<b>30.26</b>	<b>31.87</b>	<b>33.24</b>	<b>0.39</b>
ANCOR[10]	✓				87.64	7.10	6.69	6.55	6.36	6.57	6.42	6.55	6.55	6.40	5.17	13.82	0.61
Jt. train. (upp. bd.)	✓	✓	✓		76.15	48.24	39.89	34.09	32.21	30.85	28.81	29.86	28.57	28.74	29.06	36.95	0.32

Table 5: Comparison with others on CIFAR-100 (top table) and tieredImageNet (bottom table). '+' means improvement. Best seen on computer.

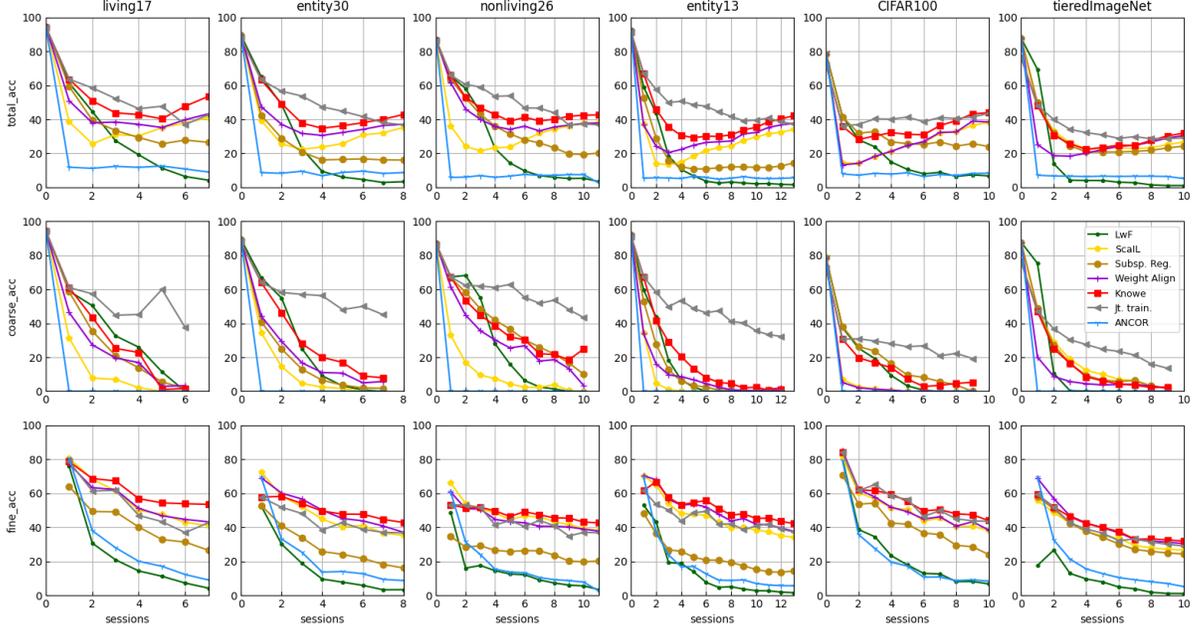


Figure 10: Separated accuracy comparison on all datasets. Top-down: total, corese, fine; red is Knowe, grey is joint training.

## 6.4 Performance Comparison and Analysis

Table 3,4,5 and Fig. 10 compare Knowe with SOTA FSCIL/IL methods including LwF [29], ScaIL [44], Weight Aligning [30] and Subspace Regularizers (Sub. Reg.) [64]. Joint training is non-IL and an  $acc$  upper bound in principle.

**Overall metrics:** average  $acc \bar{A}$  and forgetting rate  $\mathcal{F}$ . As presented in Table 3, 4, 5, Knowe has the smallest  $\mathcal{F}$  and the largest  $\bar{A}$  on all datasets. From both metrics, Weight Aligning ranks 2nd on BREEDS, Sub. Reg. ranks 2nd on CIFAR-100, and ScaIL ranks 2nd on tIN. There is a consistency of two metrics. LwF often has poor numbers, which implies that, with no samples retained, KD does not help.

**Total accuracy** per session decreases over time yet slower and slower for Knowe and SOTA methods. However, outstanding ones decrease first and then rise, because that the proportion of fine classes in the query set gets higher and their accuracy plays a leading role in the total accuracy. Knowe is the best, with a strong rising trend, which satisfies the aim of CIL the most and envisions Knowe continuing performing well when more sessions are added (Table 4). Sub. Reg. and Weight Align. often have 2nd-best numbers (both freeze weights); ScaIL and LwF occasionally do.

**Coarse class accuracy** decreases over time unavoidably (see Fig. 10), while Knowe and SOTA methods slow down the decay, with comparable rates. As IL methods, Weight Aligning, ScaIL, and LwF do not forget knowledge totally although they do not operate weights as done by Knowe. As a non-IL approach, ANCOR totally forgets old knowledge from the 1st session because it fine-tunes on the few fine shots without any extra operation to retain coarse knowledge. The joint training on all fine classes till the present is non-IL, and in principle should bound the fine-class performance. Interestingly, it also suffers less from coarse  $acc$  decay, the rate of which is much lower (Fig. 10). Differently, the cause can be imbalance between increasing fine classes and existing coarse classes. Knowe’s performance is very competitive and indeed bounded by joint training.

**Fine class’s total accuracy** normally decreases over time yet slower and slower for Knowe and SOTA methods (Fig. 10), and can be maintained in a similar range for most methods, among which Knowe often stays the highest, ScaIL and Weight Aligning are in the middle, Sub Reg. often stays in a low level, and LwF and ANCOR perform stably the worst. Knowe is the most balanced, while Sub. Reg. biases towards stability that is its drawback. Joint training does not necessarily bound the accuracy, possibly due to few shots.

## 7 Conclusion

In this paper, we present a challenging new problem, new metrics, insights, and an approach that solves it well in the sense of getting more balanced performance than the state-of-the-art approaches.

While it is not new to freeze or normalize weights, we are unaware of them previously being presented as a principled approach (to CIL) that is as simple as fine-tuning. It makes pre-trained big models more useful for finer-grained tasks.

For C2FSCIL with a linear classifier, **weights seem to be the knowledge**. However, how generic are our findings in practice? Can they be applied to general FSCIL? If yes, we are more comfortable with that answer, but then how does a class hierarchy make a difference? Future work will include examining those questions, non-linear classifiers, and so on.

## References

- [1] Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online continual learning in image classification: An empirical survey. *ArXiv preprint:2101.10423*, 2021.
- [2] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Unvieristy of Toronto: Technical Report*, 2009.
- [3] Robert M.French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3, 1999.
- [4] 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. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [5] Dimitris Fotakis, Alkis Kalavasis, Vasilis Kontonis, and Christos Tzamos. Efficient algorithms for learning from coarse labels. In *34th Annual Conference on Learning Theory*, 2021.
- [6] Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *ICCV*, 2021.
- [7] Martial Mermillod, Aurélie Bugajska, and Patrick Bonin. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers in Psychology*, 4, 2013.
- [8] Shipeng Wang, Xiaorong Li, Jian Sun, and Zongben Xu. Training networks in null space of feature covariance for continual learning. In *CVPR*, 2021.
- [9] M McCloskey and N Cohen. Catastrophic interference in connectionist networks: the sequential learning problem. *The Psychology of Learning and Motivation*, 24:109–164, 1989.
- [10] Guy Bukchin, Eli Schwartz, Kate Saenko, Ori Shahar, Rogerio Feris, Raja Giryes, and Leonid Karlinsky. Fine-grained angular contrastive learning with coarse labels. In *CVPR*, 2021.
- [11] Yuanhong Xu, Qi Qian, Hao Li, Rong Jin, and Juhua Hu. Weakly supervised representation learning with coarse labels. In *ICCV*, 2021.
- [12] Jinhai Yang, Hua Yang, and Lin Chen. Towards cross-granularity few-shot learning: Coarse-to-fine pseudo-labeling with visual-semantic meta-embedding. In *ACM Conference on Multimedia*, 2021.
- [13] Marko Ristin, Juergen Gall, Matthieu Guillaumin, and Luc Van Gool. From categories to subcategories: Large-scale image classification with partial class label refinement. In *CVPR*, 2015.
- [14] Marko Ristin, Matthieu Guillaumin, Juergen Gall, and Luc Van Gool. Incremental Learning of NCM Forests for Large-Scale Image Classification. In *CVPR*, 2014.
- [15] Abhijit Bendale and Terrance Boult. Towards open world recognition. In *CVPR*, 2015.
- [16] Sylvestre-Alvise Rebuff, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. iCaRL: Incremental Classifier and Representation Learning. In *CVPR*, 2017.
- [17] Ilja Kuzborskij, Francesco Orabona, and Barbara Caputo. From n to n+1: Multiclass transfer incremental learning. In *CVPR*, 2013.

- [18] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [19] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong. Few-shot class-incremental learning. In *CVPR*, 2020.
- [20] Xiaoyu Tao, Xinyuan Chang, Xiaopeng Hong, Xing Wei, and Yihong Gong. Topology-preserving class-incremental learning. In *ECCV*, 2020.
- [21] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B. Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? In *ECCV*, 2020.
- [22] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *CVPR*, 2018.
- [23] Mengye Ren, Renjie Liao, Ethan Fetaya, and Richard S. Zemel. Incremental few-shot learning with attention attractor networks. In *NeurIPS*, 2019.
- [24] Songlin Dong, Xiaopeng Hong, Xiaoyu Tao, Xinyuan Chang, and Xing Wei. Few-shot class-incremental learning via relation knowledge distillation. In *AAAI*, 2021.
- [25] Chi Zhang, Nan Song, Guosheng Lin, Yun Zheng, Pan Pan, and Yinghui Xu. Few-shot incremental learning with continually evolved classifiers. In *CVPR*, 2021.
- [26] Mohamed Abdelsalam, Mojtaba Faramarzi, Shagun Sodhani, and Sarath Chandar. IIRC: Incremental Implicitly-Refined Classification. In *CVPR*, 2021.
- [27] 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. *ArXiv preprint:2010.15277*, 2020.
- [28] Hinton Geoffrey, Vinyals Oriol, and Dean Jeff. Distilling the knowledge in a neural network. In *NeurIPS*, 2015.
- [29] Zhizhong Li and Derek Hoiem. Learning without forgetting. In *ECCV*, 2016.
- [30] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shutao Xia. Maintaining discrimination and fairness in class incremental learning. In *CVPR*, 2020.
- [31] Chen He, Ruiping Wang, and Xilin Chen. A tale of two cils: The connections between class incremental learning and class imbalanced learning, and beyond. In *CVPR Workshops*, 2021.
- [32] Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation and self-supervision for incremental learning. In *CVPR*, 2021.
- [33] Eden Belouadah, Adrian Popescu, and Ioannis Kanellos. Initial classifier weights replay for memoryless class incremental learning. *ArXiv preprint:2008.13710*, 2020.
- [34] Pingbo Pan, Siddharth Swaroop, Alexander Immer, Runa Eschenhagen, Richard E. Turner, and Mohammad Emtiyaz Khan. Continual deep learning by functional regularisation of memorable past. In *NeurIPS*, 2020.
- [35] Yaoyao Liu, Bernt Schiele, and Qianru Sun. Adaptive aggregation networks for class-incremental learning. In *CVPR*, 2021.
- [36] Pravendra Singh, Vinay Kumar Verma, Pratik Mazumder, Lawrence Carin, and Piyush Rai. Calibrating cnns for lifelong learning. In *NeurIPS*, 2020.
- [37] Pravendra Singh, Pratik Mazumder, Piyush Rai, and Vinay P. Namboodiri. Rectification-based knowledge retention for continual learning. In *CVPR*, 2021.
- [38] Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming catastrophic forgetting by incremental moment matching. In *NIPS*, 2017.
- [39] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and Qianru Sun. Mnemonics training: Multi-class incremental learning without forgetting. In *CVPR*, 2020.
- [40] 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*, 2019.

- [41] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *ECCV*, 2018.
- [42] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *CVPR*, 2018.
- [43] Steven C. Y. Hung, Cheng-Hao Tu, Cheng-En Wu, Chien-Hung Chen, Yi-Ming Chan, and Chu-Song Chen. Compacting, picking and growing for unforgetting continual learning. In *NeurIPS*, 2019.
- [44] Eden Belouadah and Adrian Popescu. ScaLL: Classifier weights scaling for class incremental learning. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020.
- [45] Janghyeon Lee, Hyeong Gwon Hong, Donggyu Joo, and Junmo Kim. Continual learning with extended kronecker-factored approximate curvature. In *CVPR*, 2020.
- [46] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *CVPR*, 2018.
- [47] Yujun Shi, Li Yuan, Yunpeng Chen, and Jiashi Feng. Continual learning via bit-level information preserving. In *CVPR*, 2021.
- [48] Sangwon Jung, Hongjoon Ahn, Sungmin Cha, and Taesup Moon. Continual learning with node-importance based adaptive group sparse regularization. In *NeurIPS*, 2020.
- [49] Xinting Hu, Kaihua Tang, Chunyan Miao, Xian-Sheng Hua, and Hanwang Zhang. Distilling causal effect of data in class-incremental learning. In *CVPR*, 2021.
- [50] Anna Kukleva, Hilde Kuehne, and Bernt Schiele. Generalized and incremental few-shot learning by explicit learning and calibration without forgetting. In *ICCV*, 2021.
- [51] Yaqing Wang and Quanming Yao. Few-shot learning: A survey. *ArXiv preprint:1904.05046*, 2019.
- [52] Jun Shu, Zongben Xu, and Deyu Meng. Small sample learning in big data era. *arXiv preprint arXiv:1808.04572*, 2018.
- [53] Nihar Bendre, Hugo Terashima Marín, and Peyman Najafirad. Learning from few samples: A survey. *ArXiv preprint:2007.15484*, 2020.
- [54] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017.
- [55] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017.
- [56] Guile Wu, Shaogang Gong, and Pan Li. Striking a balance between stability and plasticity for class-incremental learning. In *ICCV*, 2021.
- [57] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, 2019.
- [58] Feng Wang, Xiang Xiang, Jian Cheng, and Alan L. Yuille. Normface: L2 hypersphere embedding for face verification. In *ACM Conference on Multimedia*, 2017.
- [59] Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyang Wu, and Rama Chellappa. Learning without memorizing. In *CVPR*, 2019.
- [60] Shibani Santurkar, Dimitris Tsipras, and Aleksander Madry. Breeds: Benchmarks for subpopulation shift. *ArXiv preprint:2008.04859*, 2020.
- [61] Chen Liu, Yanwei Fu, Chengming Xu, Siqian Yang, Jilin Li, Chengjie Wang, and Li Zhang. Learning a few-shot embedding model with contrastive learning. In *AAAI*, 2021.
- [62] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- [63] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classification. In *ICLR*, 2018.
- [64] Afra Feyza Akyürek, Ekin Akyürek, Derry Wijaya, and Jacob Andreas. Subspace regularizers for few-shot class incremental learning. *Arxiv preprint:2110.07059*, 2021.

# APPENDIX

## A Introduction

In this analysis, we decouple the embedding learner and classifier, a linear FC layer, freeze weights of the embedding learner, and use the conventional Softmax cross-entropy loss. Different from convention FC layer, we freeze weights of neurons corresponding to previously-seen classes.

## B Prior Art

We set the base of our analysis with two definitions [8]. As we only analyze the last layer, we take off layer index  $l$ .

**Definition A** (Stability). *When the model  $\Theta$  is being trained in the  $t$ -th session,  $\Delta \mathbf{w}_{t,s}$  in each session should lie in the null space of the uncentered feature covariance matrix  $\bar{\mathcal{X}}_{t-1} = [\mathbf{X}_{1,1}^T, \dots, \mathbf{X}_{t-1,t-1}^T]^T$ , namely, if  $\bar{\mathcal{X}}_{t-1} \Delta \mathbf{w}_{t,s} = 0$  holds, then  $\Theta$  is stable at the  $t$ -th session's  $s$ -th step.*

Note that  $\mathbf{w}$  is the classification-layer's weight vector,  $\Delta \mathbf{w}$  is the change of  $\mathbf{w}$ ,  $t$  indexes the session, and  $s$  indexes the training step.  $\mathbf{X}_{p,p}$  where  $p < t$  in  $\bar{\mathcal{X}}_{t-1}$  is the input features of classification-layer on  $p$ -th session using classification-layer's weight trained on  $p$ -th session. We call it the absolute stability where the equality condition is strict.

**Definition B** (Plasticity). *Assume that the network  $\Theta$  is being trained in the  $t$ -th session, and  $\mathbf{g}_{t,s} = \{g_{t,s}^1, \dots, g_{t,s}^L\}$  denotes the parameter update generated by Gradient Descent for training  $\Theta$  at step  $s$ . If  $\langle \Delta \mathbf{w}_{t,s}, \mathbf{g}_{t,s} \rangle > 0$  holds, then  $\Theta$  preserves plasticity at the  $t$ -th session's  $s$ -th step.*

Notably, if the inequality condition holds, then the  $\Theta$ 's loss deceases, which is the essence, and thus  $\Theta$  is learning.

## C Our Extension of Stability

**Definition 1** (Stability Decay). *For the same input sample, let  $\tilde{\mathbf{o}}_i^{(t)}$  denote the output logits of the  $i$ -th neuron in the last layer in the  $t$ -th session. After the loss reaches the minimum, we define the degree of stability as*

$$\mathcal{D} = \sum_i \left( \frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}} \right)^2. \quad (8)$$

**Definition 2** (Relative Stability). *Given models  $\Theta_a$  and  $\Theta_b$ , if  $0 \leq \mathcal{D}_a < \mathcal{D}_b$ , then we say  $\Theta_a$  is more stable than  $\Theta_b$ .*

## D Our Proof of Proposition 1

**Proposition 1** (Normalizing or freezing weights improves stability; doing both improves the most). *Given  $\Theta_a$ , if we only normalize weights of a linear FC classifier, we obtain  $\Theta_b$ ; if we only freeze them, we obtain  $\Theta_c$ ; if we do both, we obtain  $\Theta_d$ . Then,  $\mathcal{D}_d < \mathcal{D}_b < \mathcal{D}_a$  and  $\mathcal{D}_d < \mathcal{D}_c < \mathcal{D}_a$ .*

**Proof.** (1) Stability Degree of model  $\Theta_a$ .

It is assumed that the training for all sessions will reach minimum loss. For the training sample  $m$  in 0-th session, the probability that  $m$  belongs to superclass is one, i.e.,  $p_{t,c_{super}}^m = 1$  and  $p_{t,i}^m = 0 (i \neq c_{super})$ . According to  $p_i^m = \frac{\exp(o_i^m)}{\sum_{j=1}^J \exp(o_j^m)}$ , the following conditions are satisfied,

$$\tilde{\mathbf{o}}_{c_{super}}^{(t)} = a (a \in \mathbb{R}), \tilde{\mathbf{o}}_i^{(t)} (i \neq c_{super}) = -\infty. \quad (9)$$

After training of  $T$ -th session has reached the minimum loss,  $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = b (b \in \mathbb{R})$ ,  $\tilde{\mathbf{o}}_i^{(T)} (i \neq c_{sub}) = -\infty$ , then,

$$\begin{aligned} \mathcal{D}_a &= \sum_i \left( \frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}} \right)^2 \\ &= \left( \frac{-\infty - a}{a} \right)^2 + \left( \frac{b - (-\infty)}{-\infty} \right)^2 = \infty. \end{aligned} \quad (10)$$

(2) Stability Degree of model  $\Theta_b$ .

Under the same conditions above, the following conditions are satisfied according to  $p_i^m = \frac{\exp(\cos \theta_i^m)}{\sum_{j=1}^J \exp(\cos \theta_j^m)}$ ,

$$\tilde{\mathbf{o}}_{c_{super}}^{(t)} = 1, \tilde{\mathbf{o}}_i^{(t)} (i \neq c_{super}) = -1. \quad (11)$$

After training of  $T$ -th session has reached minimum loss,  $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = 1$ ,  $\tilde{\mathbf{o}}_i^{(T)} (i \neq c_{sub}) = -1$ , then the following holds:

$$\begin{aligned} \mathcal{D}_b &= \sum_i \left( \frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}} \right)^2 \\ &= \left( \frac{-1 - 1}{1} \right)^2 + \left( \frac{1 - (-1)}{-1} \right)^2 = 8. \end{aligned} \quad (12)$$

(3) Stability Degree of model  $\Theta_c$ .

Compared with  $\Theta_a$ , model  $\Theta_c$  freezes weights of neurons corresponding to previously-seen classes. After training of  $T$ -th session has reached its minimum loss,  $\tilde{\mathbf{o}}_{c_{super}}^{(T)} = a$ ,  $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = \infty^+$ ,  $\tilde{\mathbf{o}}_i^{(T)} (i \neq c_{super} \vee i \neq c_{sub}) = -\infty$ , where  $\infty^+ > \infty$  in order to offset the influence of  $\tilde{\mathbf{o}}_{c_{super}}^{(T)}$ , then,

$$\begin{aligned} \mathcal{D}_c &= \sum_i \left( \frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}} \right)^2 \\ &= \left( \frac{\infty^+ - (-\infty)}{-\infty} \right)^2, \\ &9 > \mathcal{D}_c > 4. \end{aligned} \quad (13)$$

(4) Stability Degree of model  $\Theta_d$ .

Compared with  $\Theta_b$ , model  $\Theta_d$  freezes weights of neurons corresponding to previously-seen classes. After training of  $T$ -th session has reached its minimum loss,  $\tilde{\mathbf{o}}_{c_{super}}^{(T)} = 1$ ,  $\tilde{\mathbf{o}}_{c_{sub}}^{(T)} = 1$ ,  $\tilde{\mathbf{o}}_i^{(T)} (i \neq c_{super} \vee i \neq c_{sub}) = -1$ , then,

$$\begin{aligned} \mathcal{D}_d &= \sum_i \left( \frac{\tilde{\mathbf{o}}_i^{(T)} - \tilde{\mathbf{o}}_i^{(t)}}{\tilde{\mathbf{o}}_i^{(t)}} \right)^2 \\ &= \left( \frac{1 - (-1)}{-1} \right)^2 = 4. \end{aligned} \quad (14)$$

Comparing the stability degree of different models, we have  $\mathcal{D}_{max} = \mathcal{D}_a$ ,  $\mathcal{D}_{min} = \mathcal{D}_d$  and  $\Theta_d$  is the most stable.

## E Our Proof of Proposition 2

**Proposition 2** (Weights normalized, plasticity remains). *To train our FC classifier, if we denote the loss as  $\mathcal{L}(\mathbf{w})$  where  $\mathbf{w}$  is normalized, the weight update at each step as  $\Delta\mathbf{w}$ , and the learning rate as  $\alpha$ , then we have*

$$\mathcal{L}(\mathbf{w} - \alpha\Delta\mathbf{w}) < \mathcal{L}(\mathbf{w}). \quad (15)$$

**Proof.** For a sample  $m$  whose feature vector is  $\mathbf{x}$ , the output of the  $i$ -th neuron in linear FC layer is denoted as

$$\mathbf{o}_i = \sigma(\mathbf{x} \cdot \mathbf{w}^i) = \cos \theta_i = \frac{\mathbf{x} \cdot \mathbf{w}^i}{\|\mathbf{x}\| \|\mathbf{w}^i\|}. \quad (16)$$

The probability of sample  $m$  belonging to  $i$ -th class is

$$p_i = \frac{\exp(\mathbf{o}_i)}{\sum_{j=1}^I \exp(\mathbf{o}_j)} \quad (17)$$

And the loss of training is denoted as

$$\mathcal{L}(\mathbf{w}) = - \sum_i y_i \log(p_i) \quad (18)$$

where  $y_i$  denotes the label of sample  $m$ . Denote the weights update of the  $i$ -th neuron in linear FC layer as  $\Delta\mathbf{w}^i$ , then

$$\Delta\mathbf{w}^i = \begin{cases} (p_i - 1) \left( \frac{\mathbf{x}}{\|\mathbf{x}\| \|\mathbf{w}^i\|} - \frac{\mathbf{w}^i (\mathbf{x} \cdot \mathbf{w}^i)}{\|\mathbf{x}\| \|\mathbf{w}^i\|^3} \right), & i = c \\ p_i \left( \frac{\mathbf{x}}{\|\mathbf{x}\| \|\mathbf{w}^i\|} - \frac{\mathbf{w}^i (\mathbf{x} \cdot \mathbf{w}^i)}{\|\mathbf{x}\| \|\mathbf{w}^i\|^3} \right), & i \neq c \end{cases} \quad (19)$$

According to  $\hat{\mathbf{w}} = \mathbf{w} - \alpha\Delta\mathbf{w}$ , we have

$$\hat{\mathbf{w}}^i = \begin{cases} \mathbf{w}^i + \alpha(1 - p_i) \frac{1}{\|\mathbf{w}^i\|} \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} - \frac{\mathbf{w}^i}{\|\mathbf{w}^i\|} \cos \theta_i \right), & i = c \\ \mathbf{w}^i - \alpha p_i \frac{1}{\|\mathbf{w}^i\|} \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} - \frac{\mathbf{w}^i}{\|\mathbf{w}^i\|} \cos \theta_i \right), & i \neq c \end{cases} \quad (20)$$

By denoting  $h(\alpha) \triangleq \mathcal{L}(\mathbf{w} - \alpha\Delta\mathbf{w})$ , according to Taylor's theorem, we have

$$\mathcal{L}(\mathbf{w} - \alpha\Delta\mathbf{w}) = \mathcal{L}(\mathbf{w}) - \alpha \langle \Delta\mathbf{w}, \mathbf{g} \rangle + o(\alpha) \quad (21)$$

where  $\frac{|o(\alpha)|}{\alpha} \rightarrow 0$  when  $\alpha \rightarrow 0$ . Therefore, there exists  $\bar{\alpha} > 0$  such that

$$|o(\alpha)| < \alpha |\langle \Delta\mathbf{w}, \mathbf{g} \rangle|, \forall \alpha \in (0, \bar{\alpha}) \quad (22)$$

With  $\mathbf{g} = \frac{\partial \mathcal{L}(\hat{\mathbf{w}})}{\partial \hat{\mathbf{w}}}$ , we have  $\langle \Delta\mathbf{w}, \mathbf{g} \rangle = \sum_i \Delta\mathbf{w}^i \Delta\hat{\mathbf{w}}^i > 0$ , and thus  $\mathcal{L}(\mathbf{w} - \alpha\Delta\mathbf{w}) < \mathcal{L}(\mathbf{w})$  for all  $\alpha \in (0, \bar{\alpha})$ . Therefore, weights update  $\Delta\mathbf{w}$  is the descent direction.

## F Our Analysis of Conjecture 1

Considering a convention linear FC layer without weight normalization nor weight frozen of previously-seen classes. Let  $\mathbf{w}^i$  denotes a weight vector where  $i$  indexes the classes. When the sample's ground-truth label is  $c$ , we have

$$\Delta \mathbf{w}^i = \begin{cases} (p_i - 1)\mathbf{x}, & i = c \\ p_i \mathbf{x}, & i \neq c \end{cases} \quad (23)$$

where  $\mathbf{x}$  is the feature vector of a training sample.

**Conjecture 1** (FC weights grow over time). *Let  $\|\mathbf{W}^{(t)}\|_F$  denotes the Frobenius norm of the weight matrix formed by all weight vectors in the FC layer for new classes in the  $t$ -th session. With training converged and norm outliers ignored, it holds that  $\|\mathbf{W}^{(t)}\|_F > \|\mathbf{W}^{(t-1)}\|_F, \forall t \in \{1, \dots, T\}$ .*

**Analysis.** For a convention linear FC layer, the output of neural network directly determines the probability of which class the sample belongs to. So we use  $\Delta \mathbf{o}_i$  to represent the reward ( $\Delta \mathbf{o}_i > 0$ ) or penalty ( $\Delta \mathbf{o}_i < 0$ ) for different neurons after sample  $\mathbf{x}$  with label  $c$  is trained, where  $\mathbf{o}_i = \mathbf{x} \cdot \mathbf{w}^i$  is the output of the  $i$ -th neuron and  $\alpha > 0$  is the learning rate. Then, we have

$$\Delta \mathbf{o}_i = \begin{cases} \alpha(1 - p_i)|\mathbf{x}|^2 \geq 0, & i = c \\ -\alpha p_i |\mathbf{x}|^2 \leq 0, & i \neq c \end{cases} \quad (24)$$

For a sample  $m$  with superclass label  $c_{super}$  and subclass label  $c_{sub}$ , when we train sample  $m$  only with label  $c_{super}$  and reach a relatively good state in 0-th session, we will get  $p_{c_{super}}^m \rightarrow 1$  and  $p_i^m (i \neq c_{super}) \rightarrow 0$ . When we train sample  $m$  only with label  $c_{sub}$  in other sessions and reach a relatively good state, the penalty for superclass of sample  $m$  will be much larger than other classes, meanwhile the reward for subclass of sample  $m$  will be much larger too. Therefore, if  $i$  belongs to previously-seen classes,  $i \neq c$  will hold most of the time during training. Thus, previously-seen classes will keep being penalized during the gradient descent. As a result, the weights of previously-seen classes are prone to be smaller than those for the newly added classes. And because we train new classes in stages and reach a relatively good state (say, the training loss converges to small value) for all sessions, the FC weights will piecewisely grow over time. Therefore, the model is consequently biased towards new classes.

## G Our Analysis of Conjecture 2 and 3

Since the Conjecture 2 and Conjecture 3 are drawn from empirical observations, the following inductions will be conditioned on that the observations are always true. As a result, we present our analysis, rather than calling it a proof.

**Conjecture 3** (Sufficient & necessary condition of no impact of freezing embedding-weights).

$p \vee q \Leftrightarrow \neg r$  where

$p$ : classifier-weights are normalized,

$q$ : classifier-weights are frozen,

$r$ : freezing embedding-weights improves the performance.

As the name hints, Conjecture 2&3 is an integration of Conjecture 2 and Conjecture 3. Since  $\neg r \rightarrow p \vee q$  is the contrapositive proposition of  $\neg p \wedge \neg q \rightarrow r$ , they have the same truth value. Since  $\neg p \wedge \neg q \Rightarrow r$ , we have  $\neg r \Rightarrow p \vee q$ . Furthermore, given  $p \vee q \Rightarrow \neg r$ , we have  $p \vee q \Leftrightarrow \neg r$ , which means that  $p \vee q$  is sufficient (if) and necessary (only if) for  $\neg r$ . Namely, *iff classifier-weights are either normalized or frozen, then freezing embedding-weights does not help*. In the following, we analyze Conjecture 2 and 3, respectively.

**Conjecture 2** (the 'only if' part).  $\neg p \wedge \neg q \Rightarrow r$

**Analysis.** Although Conjecture 2 is a direct formulation of the corresponding observation, we will analyze it in a general sense. We have four propositions that are all true according to our empirical observations:

- ①  $\neg p \wedge \neg q \rightarrow r$
- ②  $p \wedge q \rightarrow \neg r$

- ③  $p \wedge \neg q \rightarrow \neg r$   
 ④  $\neg p \wedge q \rightarrow \neg r$ .

They share a similar composition pattern, and thus we can summarize them as Table 6.

$P$	$Q$	$P \wedge Q$	$R$	$P \wedge Q \rightarrow R$
$\neg p$	$\neg q$	$\neg p \wedge \neg q$	$r$	$\neg p \wedge \neg q \rightarrow r$
$p$	$q$	$p \wedge q$	$\neg r$	$p \wedge q \rightarrow \neg r$
$p$	$\neg q$	$p \wedge \neg q$	$\neg r$	$p \wedge \neg q \rightarrow \neg r$
$\neg p$	$q$	$\neg p \wedge q$	$\neg r$	$\neg p \wedge q \rightarrow \neg r$

Table 6: Compound propositions.

Let us make  $p, q, r$  an realization of general propositions

$P$ : classifier-weights are normalized,

$Q$ : classifier-weights are frozen,

$R$ : freezing embedding-weights improves the performance, respectively. We want to construct a common proposition for all the four cases all to be true. Namely, we need to solve for a compositive proposition  $\mathcal{C}(P, Q) \rightarrow R$  that satisfies the truth table with ①, ②, ③, ④ ordered top-down.

$P$	$Q$	$\mathcal{C}(P, Q)$	$R$	$\mathcal{C}(P, Q) \rightarrow R$
0	0		1	1
1	1		0	1
1	0		0	1
0	1		0	1

Table 7: A truth table that is not completed.

Note that  $A \rightarrow B$  is 0 iff  $A$  is 1 and  $B$  is 0. Therefore, we want  $\mathcal{C}(P, Q)$ 's truth value of the 2, 3, 4 line never to be 1. Given the value pairs of  $P$  and  $Q$ , the only way to make that happen is to let  $\mathcal{C}(P, Q)$  be  $\neg P \wedge \neg Q \rightarrow R$ , which is a solution that satisfies all four cases, and thus is always true.

$P$	$Q$	$\neg P \wedge \neg Q$	$R$	$\neg P \wedge \neg Q \rightarrow R$
0	0	1	1	1
1	1	0	0	1
1	0	0	0	1
0	1	0	0	1

Table 8: The truth table is realized.

Namely, we have  $\neg P \wedge \neg Q \Rightarrow R$ , which is exactly Conjecture 2,  $\neg p \wedge \neg q \Rightarrow r$ , with a change of notations.

**Conjecture 4** (the 'if' part).  $p \vee q \Rightarrow \neg r$ .

**Analysis.** Given propositions ②, ③, ④, we will combine them and derive a logically-equivalent premise.

$$\begin{aligned}
 & (p \wedge q) \vee (p \wedge \neg q) \vee (\neg p \wedge q) \\
 \Leftrightarrow & (p \vee (p \wedge \neg q) \vee (\neg p \wedge q)) \wedge (q \vee (p \wedge \neg q) \vee (\neg p \wedge q)) \\
 \Leftrightarrow & ((p \vee (p \wedge \neg q) \vee \neg p) \wedge (p \vee (p \wedge \neg q) \vee q)) \\
 \wedge & ((q \vee (p \wedge \neg q) \vee \neg p) \wedge (q \vee (p \wedge \neg q) \vee q)) \\
 \Leftrightarrow & ((p \vee p \vee \neg p) \vee (p \vee \neg q \vee \neg p) \wedge (p \vee \vee q) \wedge (p \vee \neg q \vee q)) \\
 \wedge & ((q \vee p \vee \neg p) \wedge (q \vee \neg q \vee \neg p) \wedge (q \vee p \vee q) \wedge (q \vee \neg q \vee q)) \\
 \Leftrightarrow & (1 \wedge 1 \vee (p \vee p \vee q) \wedge 1) \wedge (1 \wedge 1 \wedge (q \vee p \vee q) \wedge 1)
 \end{aligned}$$

$$\begin{aligned} &\Leftrightarrow (p \vee p \vee q) \wedge (q \vee p \vee q) \\ &\Leftrightarrow (p \vee q) \wedge (p \vee q) \Leftrightarrow p \vee q. \end{aligned}$$

Similarly, we can derive ②  $\vee$  ③  $\vee$  ④ as  
 $(p \wedge q \rightarrow \neg r) \vee (p \wedge \neg q \rightarrow \neg r) \vee (\neg p \wedge q \rightarrow \neg r)$   
 $\Leftrightarrow (p \wedge q) \vee (p \wedge \neg q) \vee (\neg p \wedge q) \rightarrow \neg r.$

With the premise replaced, we have

$$\text{②} \vee \text{③} \vee \text{④} \Leftrightarrow p \vee q \rightarrow \neg r,$$

Given ②, ③, ④ are all always true. it holds that  $p \vee q \rightarrow \neg r$  is always true, Namely, we have  $p \vee q \Rightarrow \neg r.$