# Split-and-Bridge: Adaptive Class Incremental Learning within a Single Neural Network

Jong-Yeong Kim, Dong-Wan Choi

Department of Computer Science and Engineering, Inha University, South Korea

kjy93217@naver.com, dchoi@inha.ac.kr

https://github.com/bigdata-inha/Split-and-Bridge

## Standard KD-based Class Incremental Learning

- **Class incremental learning**
  Learning tasks arrive in a sequence and deep neural network $\theta$ must continually learn to increment already acquired knowledge.

- **Rehearsal method**
  Store a subset of previous samples $M_t$, and train them together with samples of a new task $D_t$ to prevent forgetting previous knowledge.

- **Knowledge Distillation based method**
  Try to mitigate forgetting by transferring the previous knowledge distilled from the pre-trained model.

- **Standard KD-based Method loss function**

$$\lambda \mathcal{L}_{kd}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t) + (1-\lambda) \mathcal{L}_{ce}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t)$$

Softened probability(reference model)

$$\mathcal{L}_{kd}(\mathcal{D}, \Theta) = -\sum_{(\mathbf{x},\mathbf{y})\in\mathcal{D}} \overline{\hat{q}(\mathbf{x})} \log q(\mathbf{x}) \qquad \mathcal{L}_{ce}(\mathcal{D}, \Theta) = -\sum_{(\mathbf{x},\mathbf{y})\in\mathcal{D}} \mathbf{y} \log p(\mathbf{x})$$

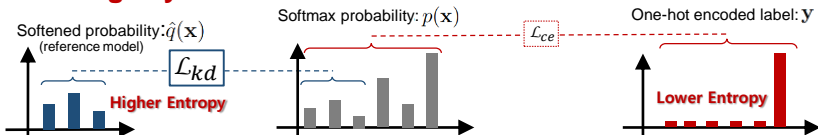Softened probability  Softmax probability

## Motivation

- **Our observation**: we can think of class incremental learning as the problem of learning 3 types of knowledge.

  Intra-old / Intra-new / Cross-task

- We can further identify **which part of the loss function is utilized** to acquire each type of knowledge in Standard KD-based method.
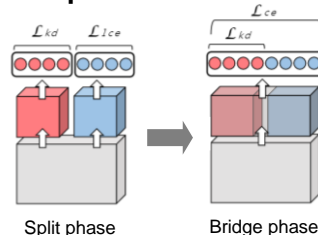
  - Intra-old knowledge: $\mathcal{L}_{kd}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t) + \mathcal{L}_{ce}(\mathcal{M}_t, \Theta_t)$
  - Intra-new knowledge: $\mathcal{L}_{ce}(\mathcal{D}_t, \Theta_t)$
  - Cross-task knowledge: $\mathcal{L}_{ce}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t)$

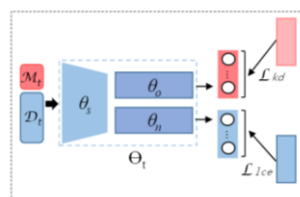- KD-based method **suffer from learning intra-new and cross-task knowledge by CE loss**

Softened probability: $\hat{q}(\mathbf{x})$ (reference model)  $\mathcal{L}_{kd}$  **Higher Entropy**

Softmax probability: $p(\mathbf{x})$  $\mathcal{L}_{ce}$

One-hot encoded label: $\mathbf{y}$  **Lower Entropy**

## Overview

- We propose a two phase learning method within a single network to learn **without any competition between losses**



Split phase    Bridge phase

## Proposed Adaptable Incremental Learning

- **Separated learning within a single network**

  To learn the intra-new knowledge as independently as possible from the task of preserving the intra-old knowledge.



  Separated partition: Old partition

  $$\Theta_t \longrightarrow \langle \theta_s, [\overline{\theta_o, \theta_n}] \rangle_t$$

  New partition

  Loss function:

  $$\mathcal{L}_{kd}(\mathcal{D}_t \cup \mathcal{M}_t, \langle \theta_s, \theta_o \rangle_t) + \mathcal{L}_{lce}(\mathcal{D}_t, \langle \theta_s, \theta_n \rangle_t)$$
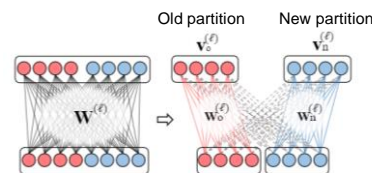
  Localized Cross Entropy:

  $$\mathcal{L}_{lce}(\mathcal{D}_t, \Theta) = -\sum_{(\mathbf{x},\mathbf{y})\in\mathcal{D}_t} \mathbf{y_t} \log p_t(\mathbf{x})$$

  Local Softmax probability

- **Weight sparsification across tasks**

  We gradually remove inter-connected weights $W_{o,n}$ and $W_{n,o}$ to get a separated network with less previous knowledge loss.
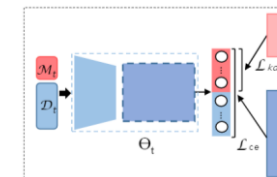


  Old partition  New partition

  Loss function:

  $$\mathcal{L}_{kd}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t) + \mathcal{L}_{lce}(\mathcal{D}_t, \Theta_t)$$
  $$+ \gamma \sum_{\ell=S+1}^{L} (||\mathbf{W}_{o,n}^{(\ell)}||_2 + ||\mathbf{W}_{n,o}^{(\ell)}||_2),$$

## Bridge phase

We re-connect two partitions $\theta_o$ and $\theta_n$ in order to learn the cross-task knowledge between them.



Re-connect separated partitions:

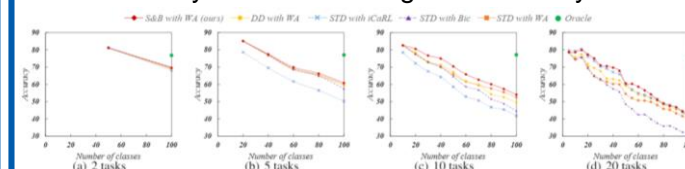$$\langle \theta_s, [\theta_o, \theta_n] \rangle_t \longrightarrow \Theta_t$$

connect through zero-initialized weight

$$\lambda \mathcal{L}_{kd}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t) + (1-\lambda) \mathcal{L}_{ce}(\mathcal{D}_t \cup \mathcal{M}_t, \Theta_t)$$
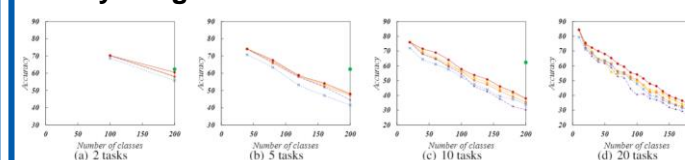
## Experimental Results

- **CIFAR-100** in ResNet-18
  consistently achieves the highest accuracy



(a) 2 tasks  (b) 5 tasks  (c) 10 tasks  (d) 20 tasks

- **Tiny-ImageNet** in ResNet-18 (highest accuracy)



(a) 2 tasks  (b) 5 tasks  (c) 10 tasks  (d) 20 tasks

- Average **intra-new and intra-old accuracy**

  highest intra-new accuracy & preserving intra-old accuracy



(a) Intra-new (CIFAR-100)  (b) Intra-old (CIFAR-100)  (c) Intra-new (Tiny-ImageNet)  (d) Intra-old (Tiny-ImageNet)