



电子科技大学  
University of Electronic Science and Technology of China



# Recent Advances of Continual Learning

Wei Han



Data Mining Lab,  
Big Data Research Center, UESTC  
Email: [weihan@std.uestc.edu.cn](mailto:weihan@std.uestc.edu.cn)

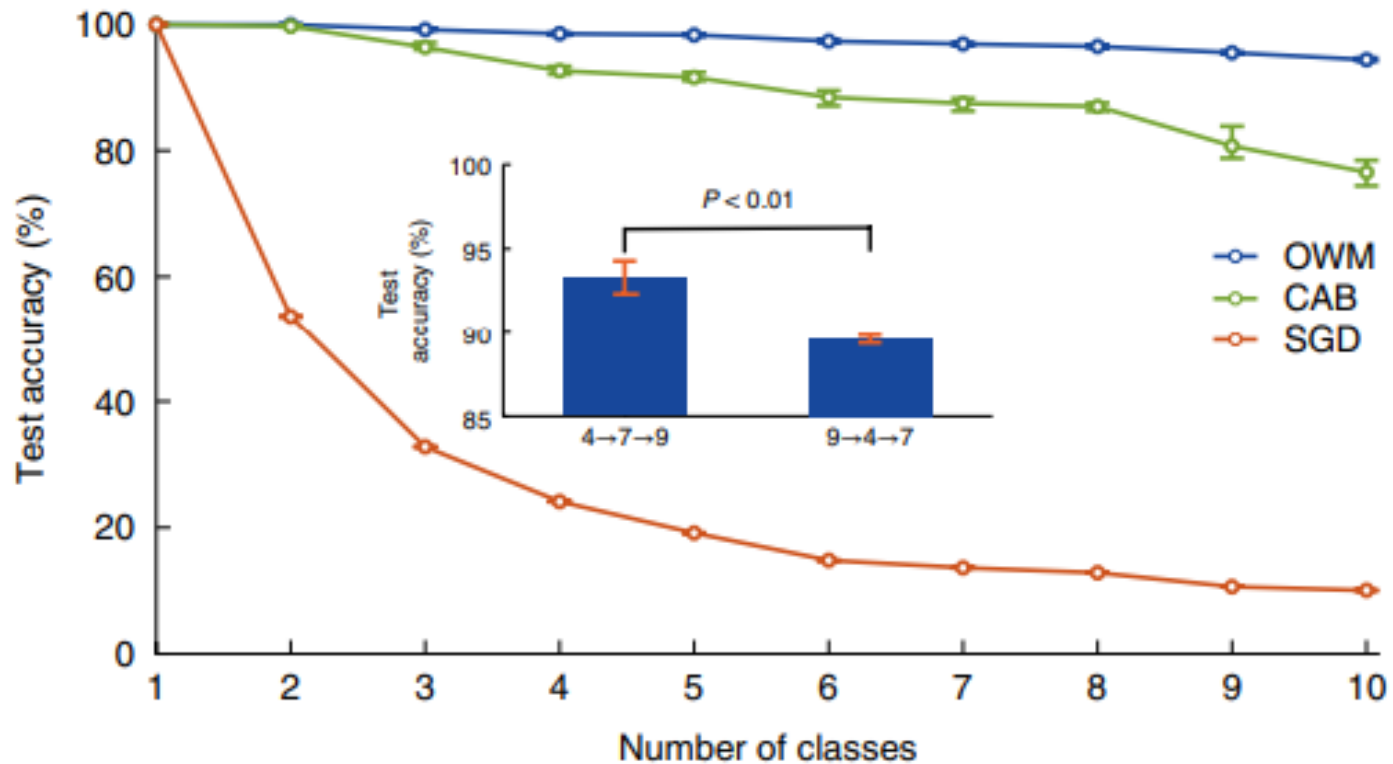
- Preliminary: Taxonomy of Continual Learning
- Supervised Continual Learning
- Unsupervised Continual Learning
- Our Proposal



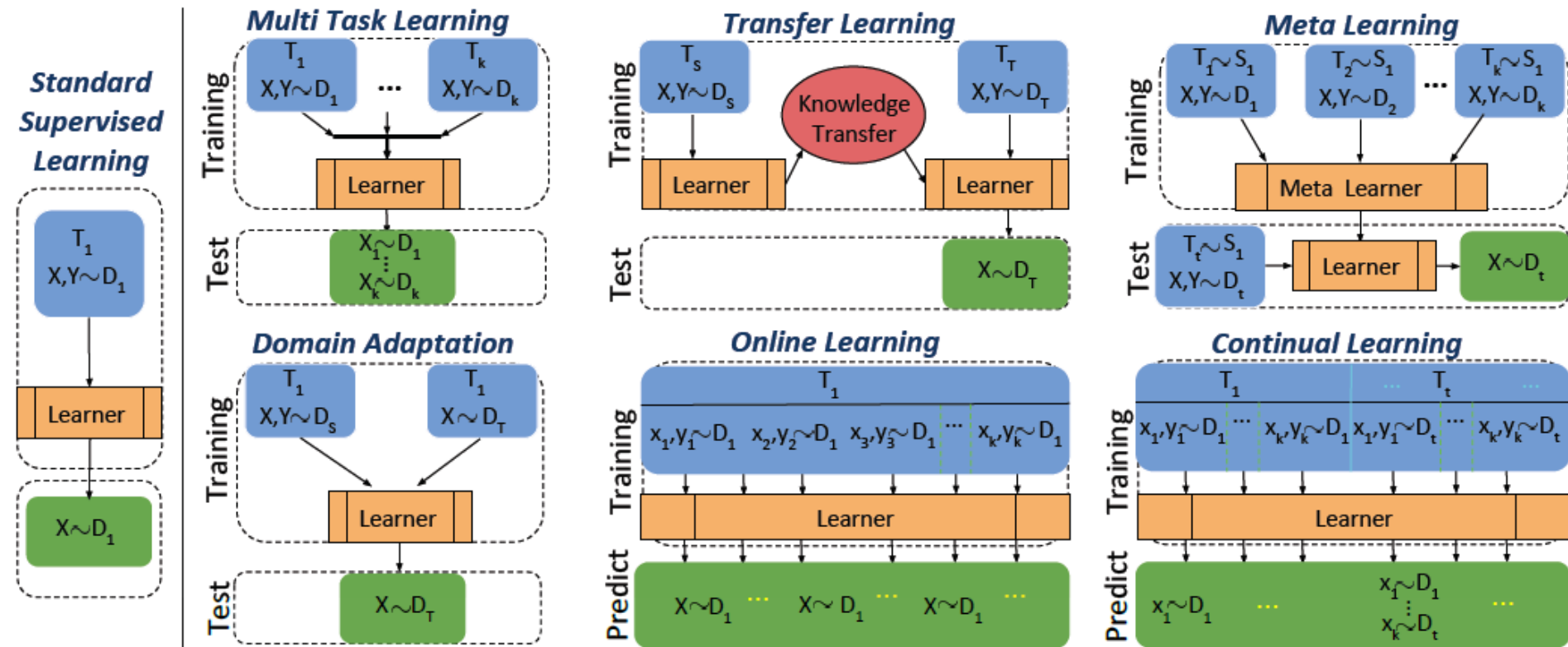
# Preliminary: Taxonomy of Continual Learning

## Continual Learning

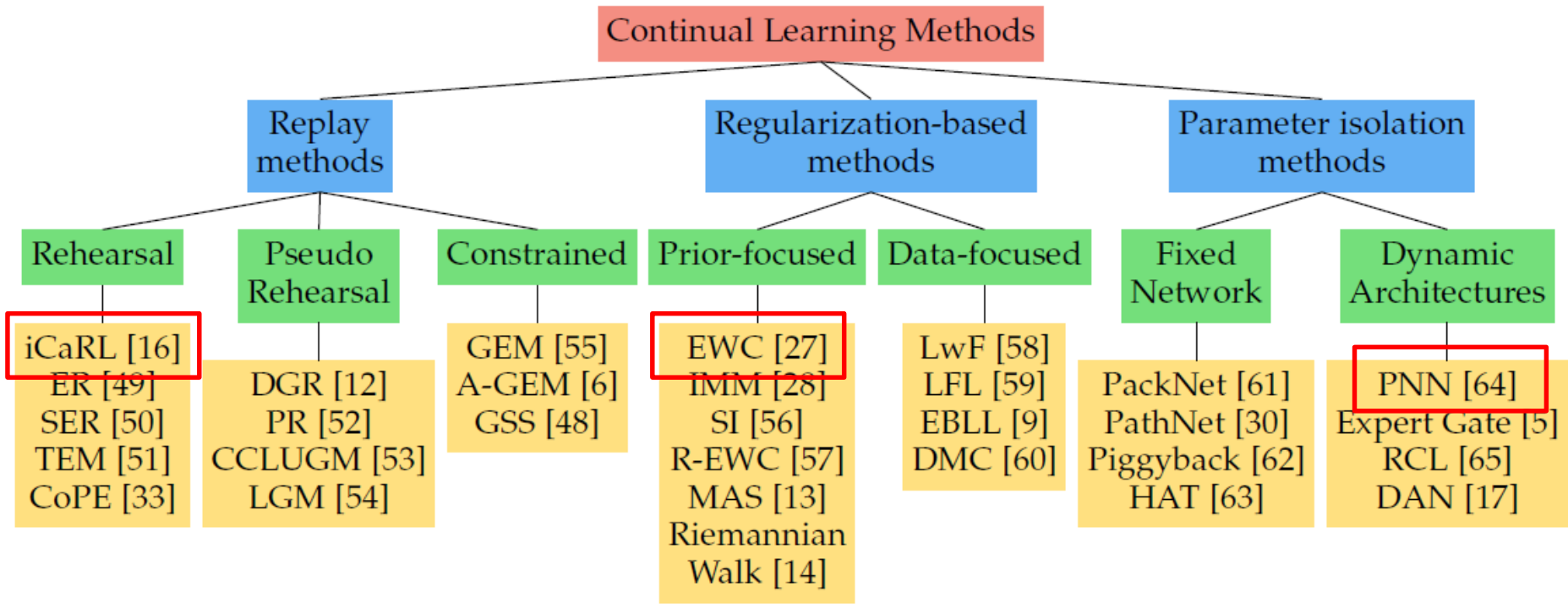
Train each class one by one



## Continual Learning



## ■ Taxonomy of Continual Learning



## ■ iCaRL (Replay-based)

---

### Algorithm 2 iCaRL INCREMENTALTRAIN

---

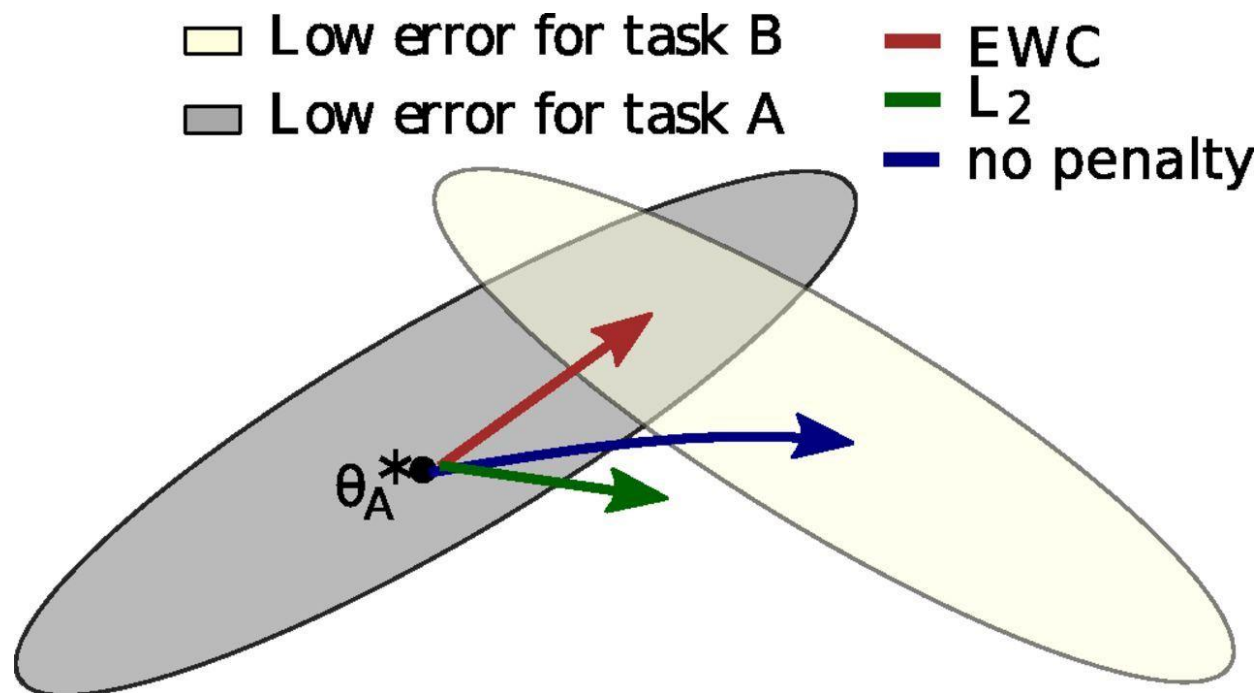
```

input  $X^s, \dots, X^t$  // training examples in per-class sets
input  $K$  // memory size
require  $\Theta$  // current model parameters
require  $\mathcal{P} = (P_1, \dots, P_{s-1})$  // current exemplar sets
 $\Theta \leftarrow \text{UPDATEREPRESENTATION}(X^s, \dots, X^t; \mathcal{P}, \Theta)$ 
 $m \leftarrow K/t$  // number of exemplars per class
for  $y = 1, \dots, s-1$  do
     $P_y \leftarrow \text{REDUCEEXEMPLARSET}(P_y, m)$ 
end for
for  $y = s, \dots, t$  do
     $P_y \leftarrow \text{CONSTRUCTEXEMPLARSET}(X_y, m, \Theta)$ 
end for
 $\mathcal{P} \leftarrow (P_1, \dots, P_t)$  // new exemplar sets

```

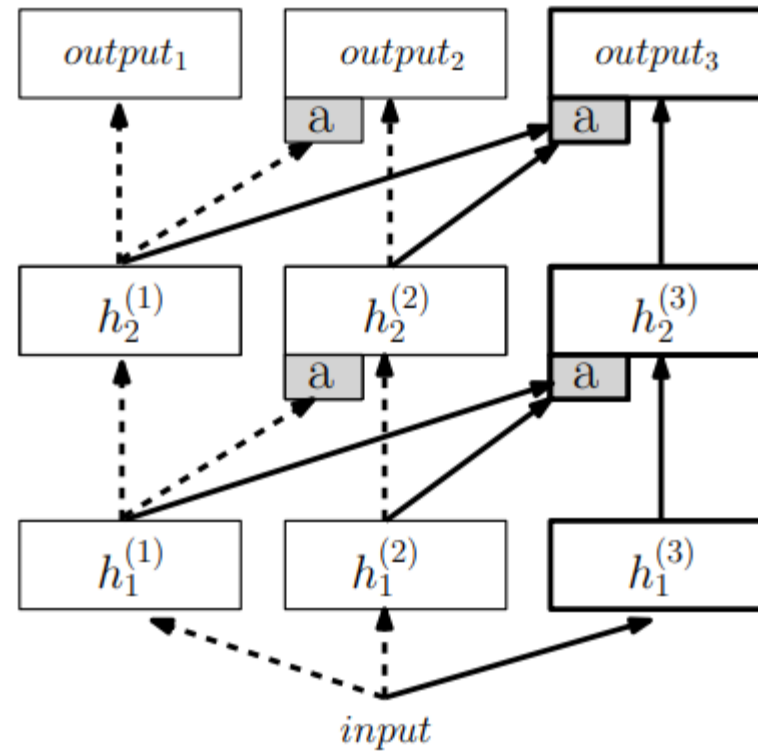
---

## ■ EWC (Regularization-based)

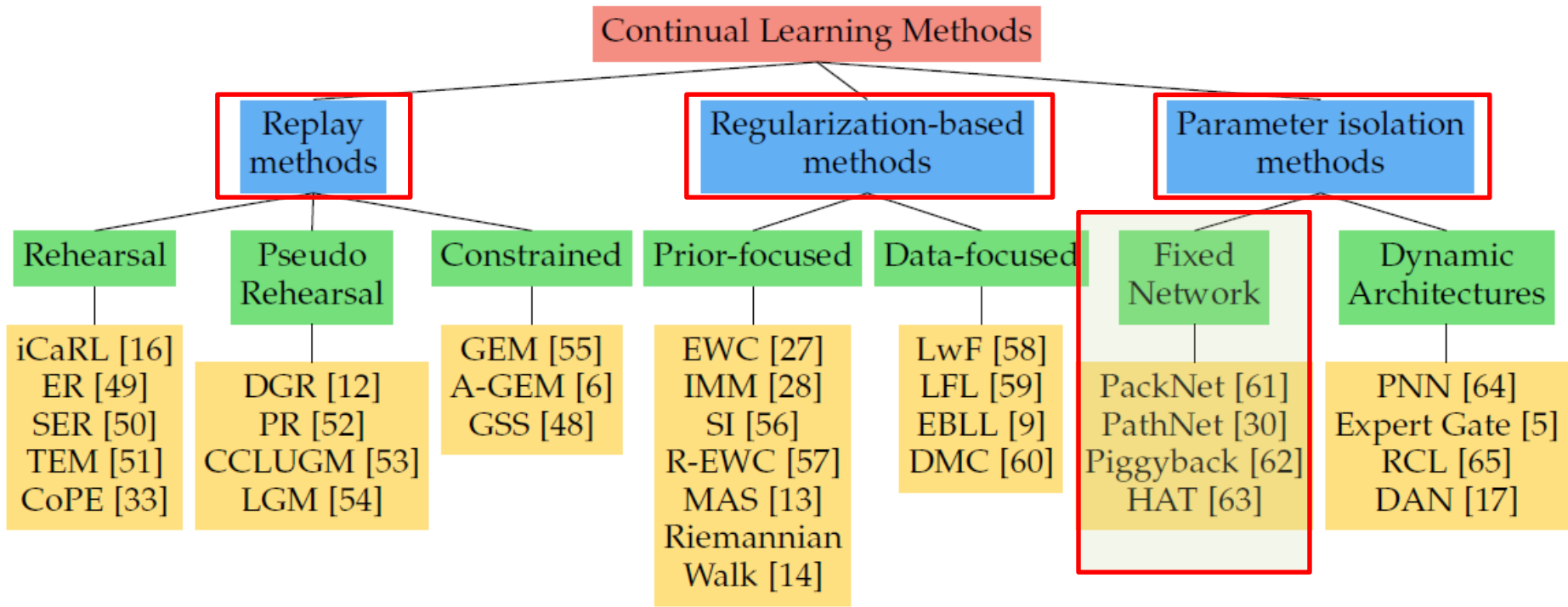




## ■ PNN (Parameter Isolation)

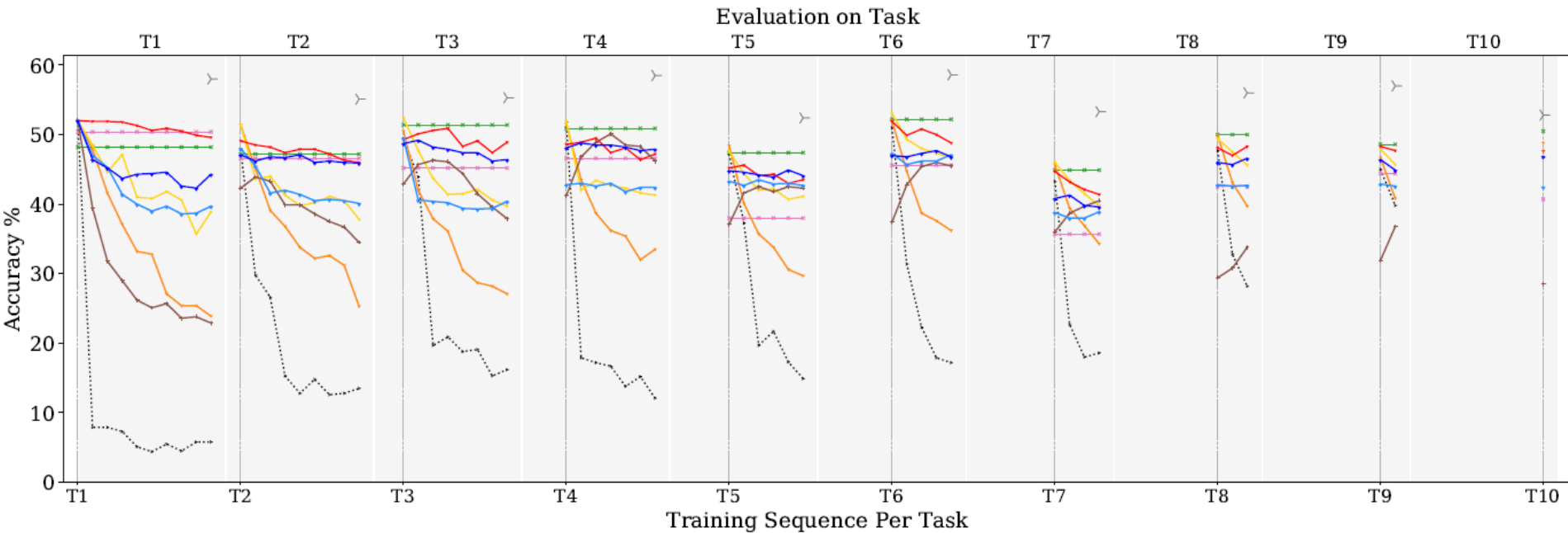


## ■ Taxonomy of Continual Learning



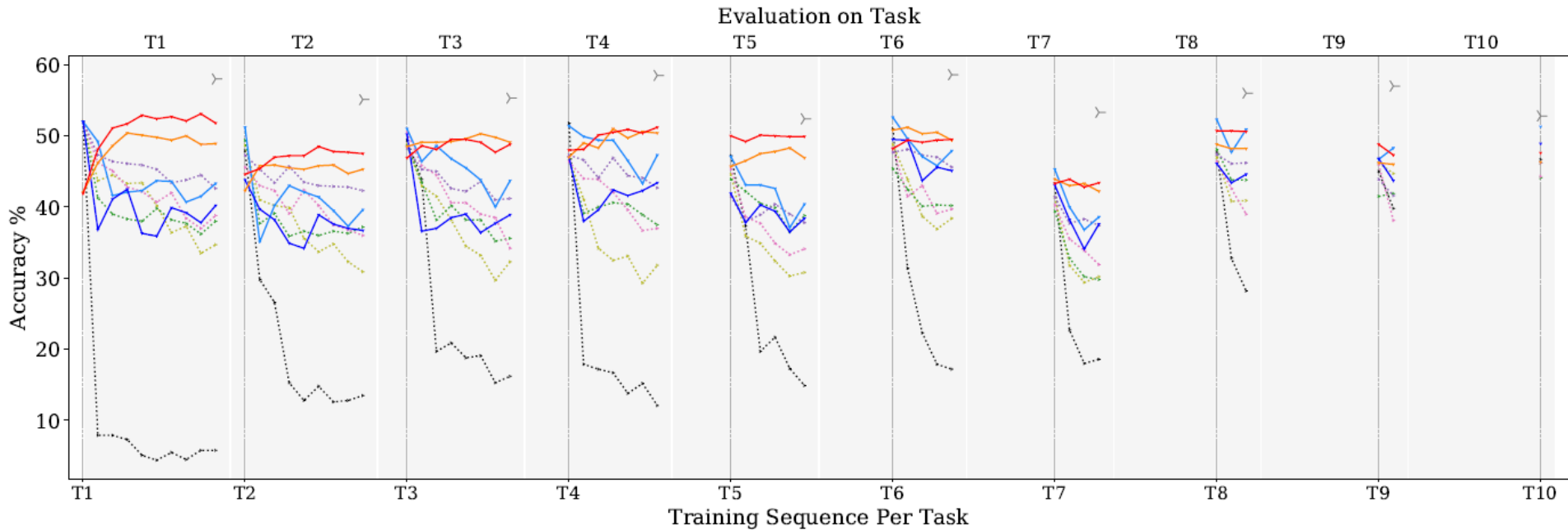
## Results of Continual Learning

finetuning: 21.30 (26.90)	PackNet: 49.13 (0.00)	SI: 33.93 (15.77)	MAS: 46.90 (1.58)	LwF: 41.91 (3.08)
joint*: 55.70 (n/a)	HAT: 43.57 (0.00)	EWC: 42.43 (7.51)	mode-IMM: 36.89 (0.98)	EBLL: 45.34 (1.44)



## Results of Continual Learning

--- finetuning: 21.30 (26.90)	--- R-PM 4.5k: 36.09 (10.96)	--- R-FM 4.5k: 37.31 (9.21)	--- GEM 4.5k: 45.13 (4.96)	--- iCaRL 4.5k: 47.27 (-1.11)
> joint*: 55.70 (n/a)	--- R-PM 9k: 38.69 (7.23)	--- R-FM 9k: 42.36 (3.94)	--- GEM 9k: 41.75 (5.18)	--- iCaRL 9k: 48.76 (-1.76)

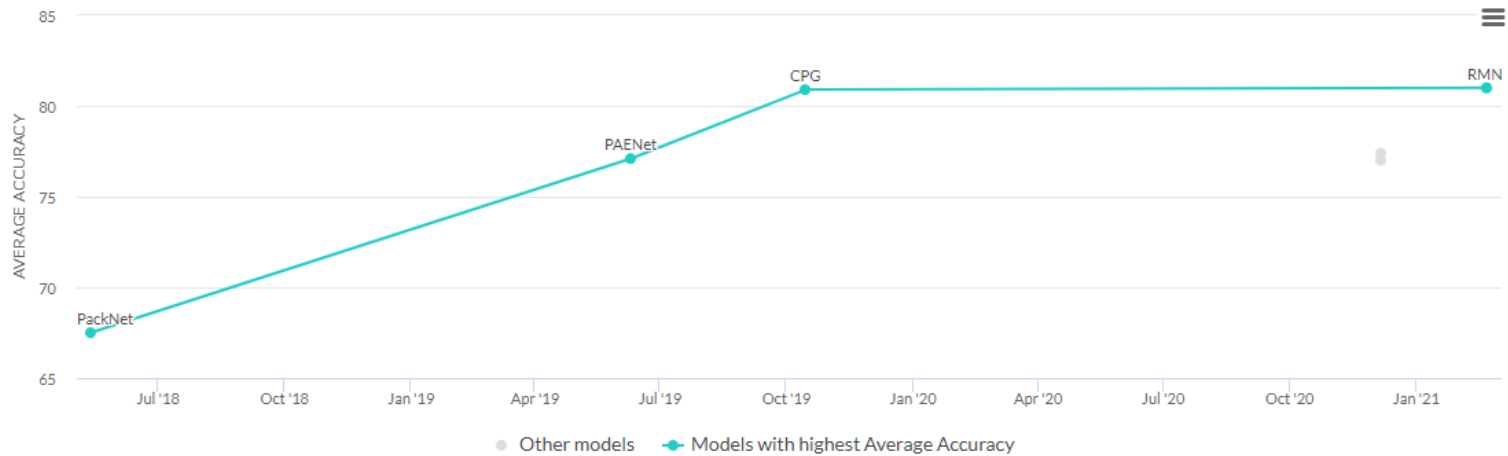


## Continual Learning on Cifar100 (20 tasks)

Leaderboard

Dataset

View Average Accuracy by Date



Filter: Strict Continual Learning untagged

Edit Leaderboard

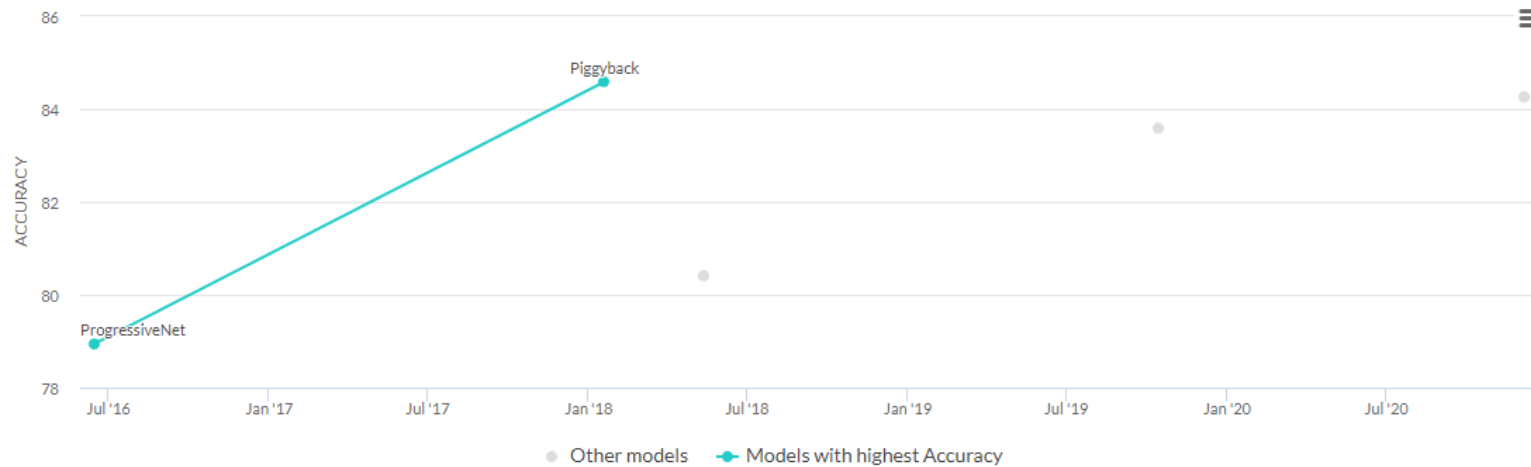
Rank	Model	Average ↑ Accuracy	Paper	Code	Result	Year	Tags
1	RMN	81	Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping			2021	<span>Strict Continual Learning</span>
2	CPG	80.9	Compacting, Picking and Growing for Unforgetting Continual Learning			2019	
3	CondConvContinual	77.4	EXTENDING CONDITIONAL CONVOLUTION STRUCTURES FOR ENHANCING MULTITASKING CONTINUAL LEARNING			2020	

## Continual Learning on CUBS (Fine-grained 6 Tasks)

Leaderboard

Dataset

View Accuracy by Date



Filter: untagged

Edit Leaderboard

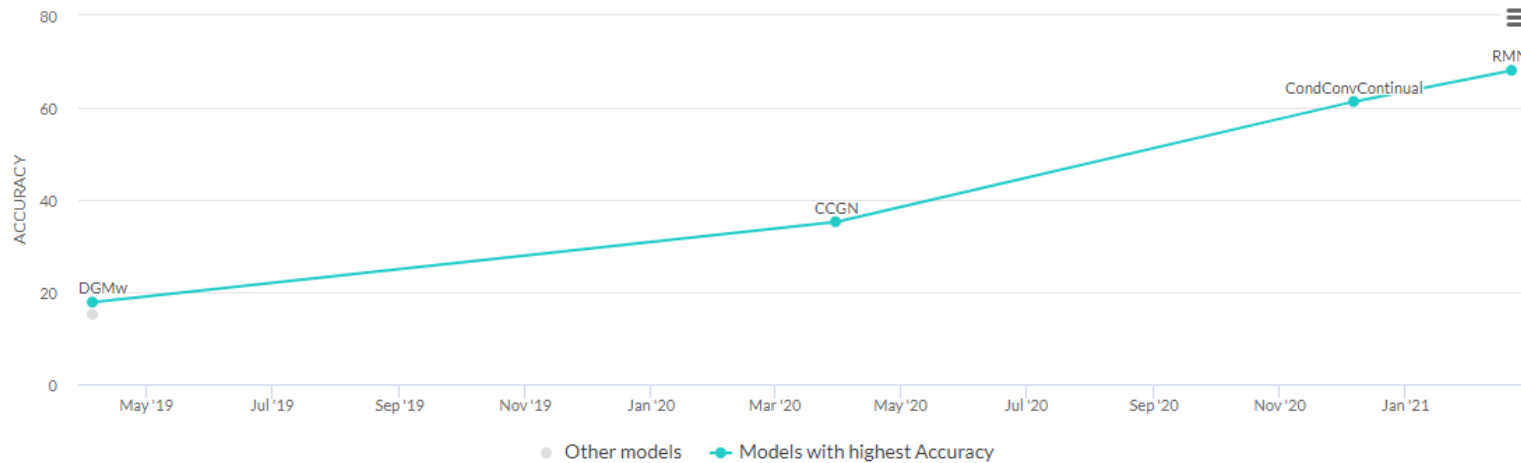
Rank	Model	Accuracy↑	Pretrained	Paper	Code	Result	Year	Tags
1	<b>Piggyback</b>	84.59	Yes	<a href="#">Piggyback: Adapting a Single Network to Multiple Tasks by Learning to Mask Weights</a>	<a href="#">Code</a>	<a href="#">Result</a>	2018	
2	<b>CondConvContinual</b>	84.26	Yes	<a href="#">EXTENDING CONDITIONAL CONVOLUTION STRUCTURES FOR ENHANCING MULTITASKING CONTINUAL LEARNING</a>	<a href="#">Code</a>	<a href="#">Result</a>	2020	
3	<b>CPG</b>	83.59	Yes	<a href="#">Compacting, Picking and Growing for Unforgetting Continual Learning</a>	<a href="#">Code</a>	<a href="#">Result</a>	2019	

## Continual Learning on ImageNet-50 (5 tasks)

Leaderboard

Dataset

View Accuracy by Date



Filter: untagged

Edit Leaderboard

Rank	Model	Accuracy↑	Paper	Code	Result	Year	Tags
1	RMN	68.1	<a href="#">Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping</a>	<a href="#">Code</a>	<a href="#">Result</a>	2021	
2	CondConvContinual	61.32	<a href="#">EXTENDING CONDITIONAL CONVOLUTION STRUCTURES FOR ENHANCING MULTITASKING CONTINUAL LEARNING</a>	<a href="#">Code</a>	<a href="#">Result</a>	2020	
3	CCGN	35.24	<a href="#">Conditional Channel Gated Networks for Task-Aware Continual Learning</a>	<a href="#">Code</a>	<a href="#">Result</a>	2020	

Oxford 102 Flowers					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	10.0 (-20.3)	5.1 (-17.1)	6.7 (-13.6)	17.3 (0.0)	9.8
Freezing	30.3 (0.0)	39.8 (0.0)	32.0 (0.0)	33.1 (0.0)	33.8
Joint	54.6 (+24.3)	58.9 (+11.5)	57.7 (+4.5)	47.0 (0.0)	54.6
EWC [14]	12.1 (-18.2)	11.6 (-38.1)	9.3 (-24.4)	25.8 (0.0)	14.7
HAT [41]	17.2 (-12.7)	19.3 (-28.5)	28.6 (+1.4)	31.6 (0.0)	24.2
PackNet [26]	32.0 (0.0)	53.7 (0.0)	43.6 (0.0)	37.9 (0.0)	41.8
TFM w/o FN	36.4 (0.0)	54.1 (0.0)	38.6 (0.0)	39.0 (0.0)	42.0
TFM	36.4 (0.0)	53.8 (0.0)	45.5 (0.0)	37.6 (0.0)	<b>43.3</b>

CUBS 200 Birds					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	7.4 (-30.2)	2.6 (-30.0)	29.7 (-3.4)	43.1 (0.0)	20.7
Freezing	37.6 (0.0)	35.1 (0.0)	35.4 (0.0)	38.4 (0.0)	36.6
Joint	48.7 (+11.1)	52.1 (+6.0)	50.7 (+1.5)	51.9 (0.0)	50.8
EWC [14]	16.2 (-21.4)	19.0 (-21.2)	24.2 (-14.0)	41.7 (0.0)	25.3
HAT [41]	18.7 (-1.8)	19.4 (-0.4)	28.5 (-0.6)	31.2 (0.0)	24.4
PackNet [26]	35.3 (0.0)	42.8 (0.0)	44.4 (0.0)	45.9 (0.0)	42.1
TFM w/o FN	42.9 (0.0)	44.1 (0.0)	48.3 (0.0)	49.1 (0.0)	46.1
TFM	42.9 (0.0)	43.1 (0.0)	49.9 (0.0)	48.8 (0.0)	<b>46.2</b>

Stanford 40 Actions					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	24.4 (-10.5)	26.5 (-7.7)	17.6 (-16.8)	28.9 (0.0)	24.4
Freezing	34.9 (0.0)	29.4 (0.0)	30.1 (0.0)	30.5 (0.0)	31.2
Joint	45.7 (+10.8)	40.3 (+4.8)	43.2 (-1.1)	40.2 (0.0)	42.4
EWC [14]	24.2 (-10.7)	28.2 (-2.0)	25.2 (-5.6)	34.3 (0.0)	28.0
HAT [41]	25.7 (-1.0)	25.5 (-2.7)	30.1 (-2.1)	34.4 (0.0)	28.9
PackNet [26]	32.5 (0.0)	32.9 (0.0)	36.7 (0.0)	34.3 (0.0)	34.1
TFM w/o FN	35.3 (0.0)	38.3 (0.0)	39.2 (0.0)	38.0 (0.0)	37.7
TFM	35.3 (0.0)	37.2 (0.0)	42.0 (0.0)	37.2 (0.0)	<b>38.0</b>





# Supervised Continual Learning

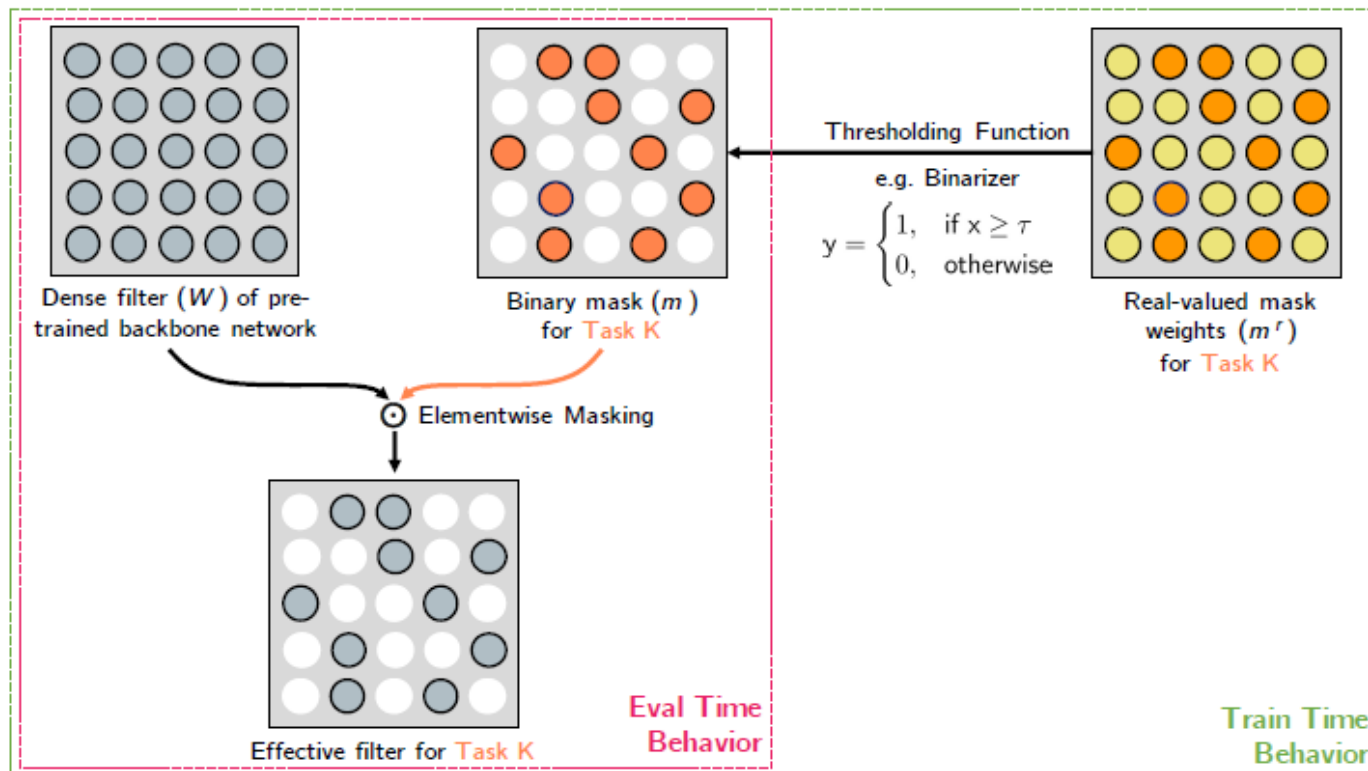
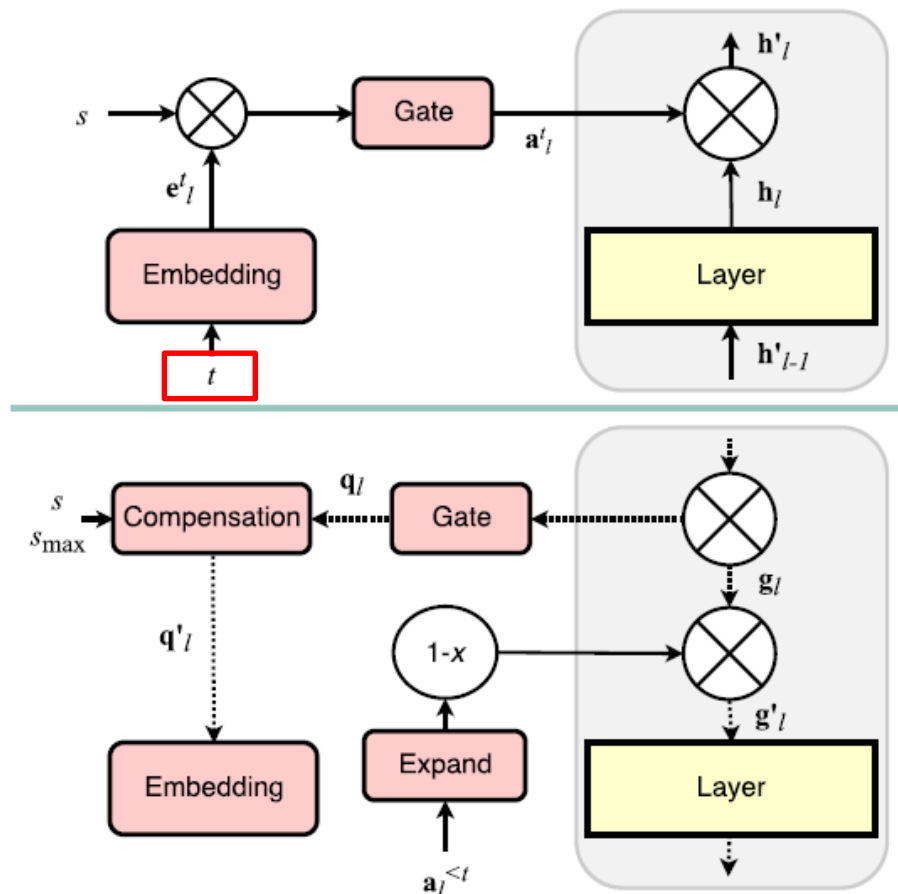


Fig. 1: Overview of our method for learning piggyback masks for fixed backbone networks. During training, we maintain a set of real-valued weights  $m^r$  which are passed through a thresholding function to obtain binary-valued masks  $m$ . These masks are applied to the weights  $W$  of the backbone network in an elementwise fashion, keeping individual weights active, or masked out. **The gradients obtained through backpropagation of the task-specific loss are used to update the real-valued mask weights.** After training, the real-valued mask weights are discarded and only the thresholded mask is retained, giving one network mask per task.

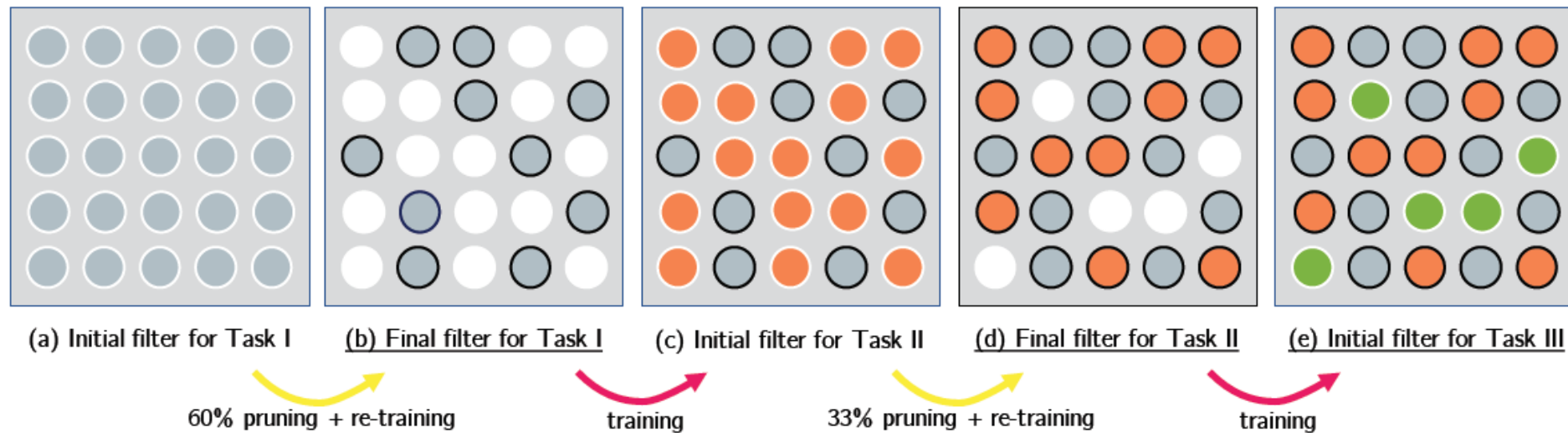
## Learnable hard attention



$$s = \frac{1}{s_{\max}} + \left( s_{\max} - \frac{1}{s_{\max}} \right) \frac{b-1}{B-1},$$

Figure 1. Schematic diagram of the proposed approach: forward (top) and backward (bottom) passes.

## ■ Pruning as Masks



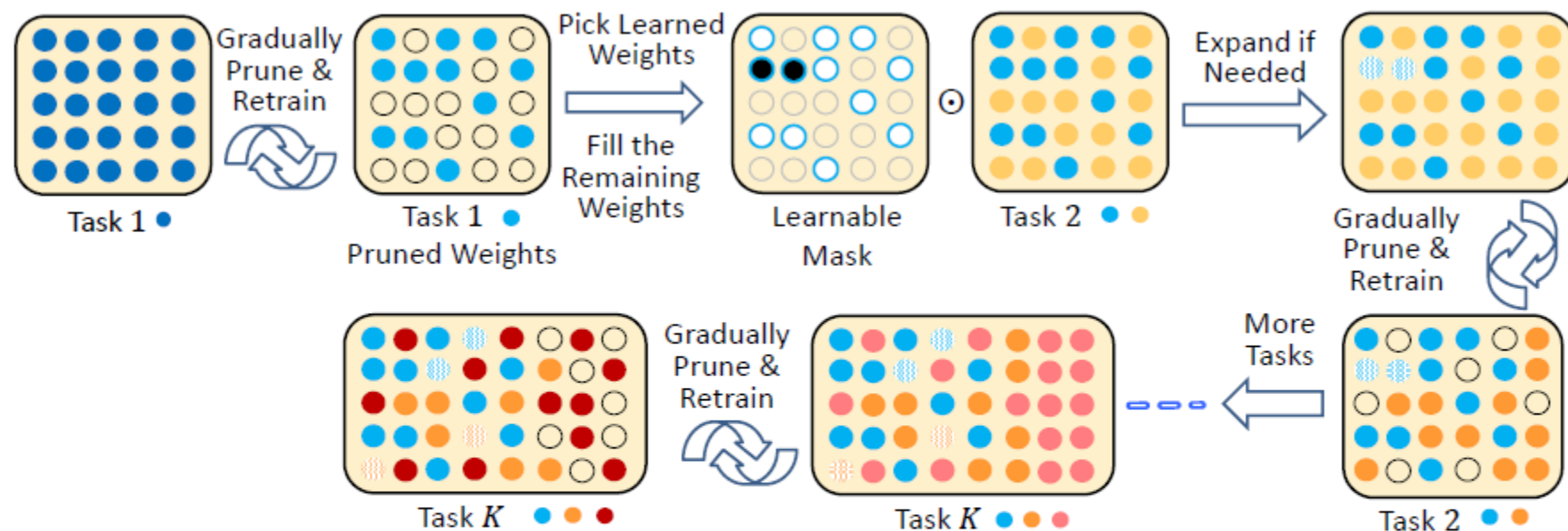


Figure 1: Compacting, Picking, and Growing (CPG) continual learning. Given a **well-trained** model, gradual **pruning** is applied to compact the model to release redundant weights. The compact model weights are kept to avoid forgetting. Then a learnable binary weight-picking mask is **trained** along with previously released space for new tasks to effectively **reuse** the knowledge of previous tasks. The model can be **expanded** for new tasks if it does not meet the performance goal. Best viewed in color.

## ■ Ternary state: ‘used’, ‘learnable’, ‘unused’

- Binary mask

$$y = (Wx) \odot m^{t,l}$$

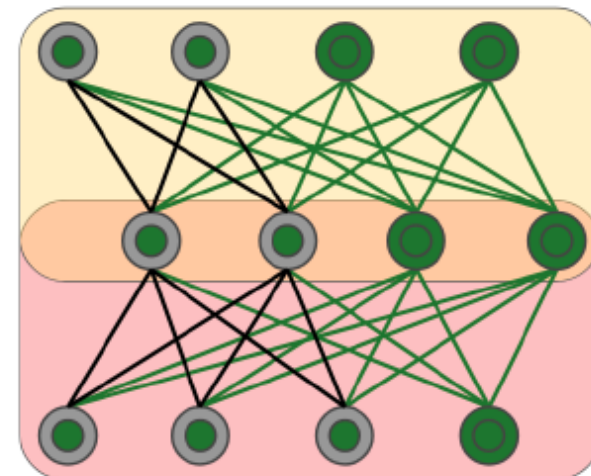
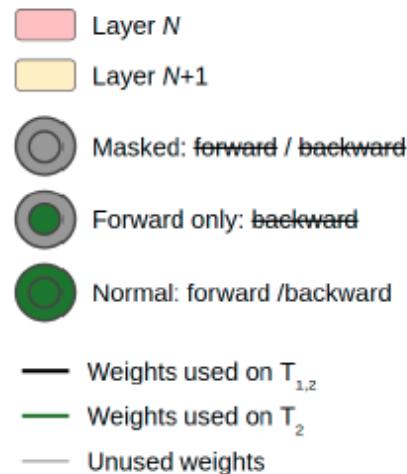
$$\frac{\partial \mathcal{L}}{\partial W_{ij}} = (m_i^{t,l} \wedge m_j^{t,l-1}) x_j \frac{\partial \mathcal{L}}{\partial y_i}$$

- Ternary mask

$$y = (Wx) \odot n^{t,l} \quad n_i^{t,l} = \begin{cases} 1, & \text{if } \exists s \leq t : m_i^{s,l} = 1 \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}} = (n_j^{t,l-1} n_i^{t,l} - n_j^{t-1,l-1} n_i^{t-1,l}) x_j \frac{\partial \mathcal{L}}{\partial y_i}$$

➔ 
$$\frac{\partial \mathcal{L}}{\partial W_{ij}} = (m_i^{t,l} \vee m_j^{t,l-1}) x_j \frac{\partial \mathcal{L}}{\partial y_i}$$



## ■ Task-specific feature normalization

$$\hat{x}_{l,i} = \gamma_{t,l,i} x_{l,i} + \beta_{t,l,i}$$

## ■ Growing Ternary Feature Map

$$n_j^{t,l} = \begin{cases} 1, & \text{for current task, if } 1 \leq j \leq I+N \\ 0, & \text{for previous tasks, if } I < j \leq I+N \end{cases}$$

$$m_j^{t,l} = \begin{cases} 0, & \text{for current task, if } 1 \leq j \leq I \\ 1, & \text{for current task, if } I < j \leq I+N \\ 0, & \text{for previous tasks, if } I < j \leq I+N \end{cases}$$

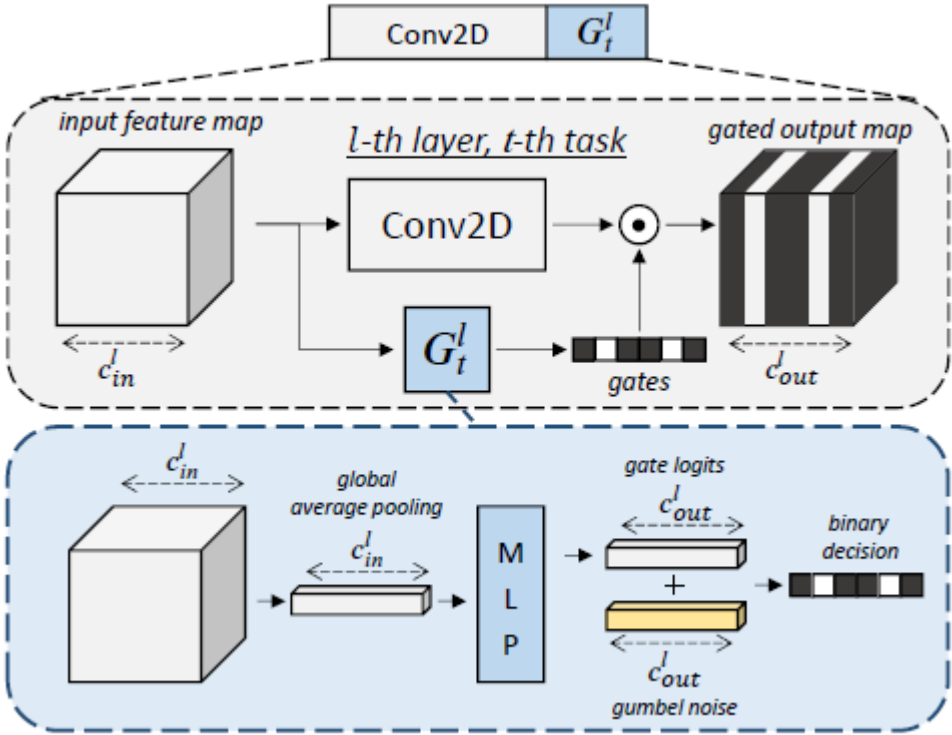
- Differentiable Mask?
- Knowledge Reuse?
- Expandable?
- Computational Complexity (upon *weights/features*)



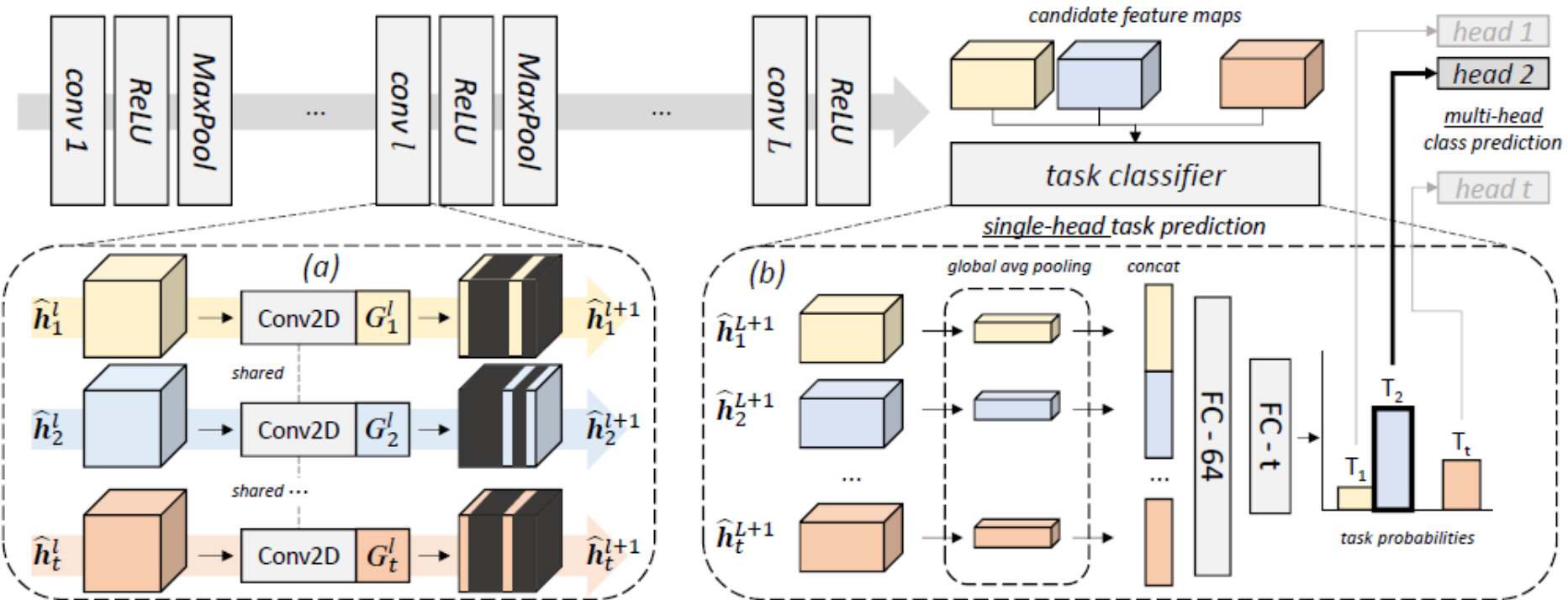


# Unsupervised Continual Learning

## Conditional Channel Gate



## Task Classifier



## ■ Conditional Convolution

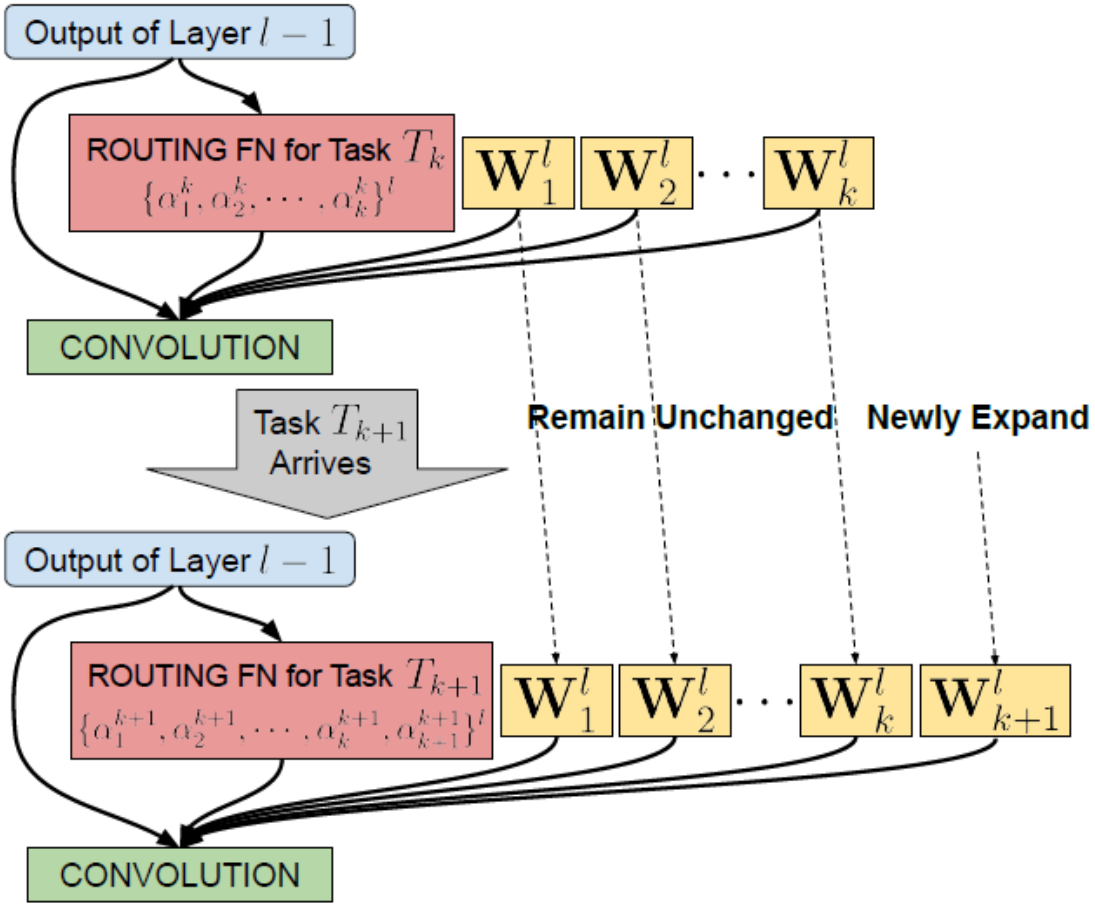
$$\text{CondConv}_n(\mathbf{x}) = \sigma((\alpha_1 \mathbf{W}_1 + \alpha_2 \mathbf{W}_2 + \cdots + \alpha_n \mathbf{W}_n) * \mathbf{x}),$$

$$\alpha_i = \text{Sigmoid}(\text{GlobalAveragePool}(\mathbf{x}) \mathbf{R}_i + \mathbf{b}_i),$$

$$M_{k+1}(\mathbf{x}) = \underline{C_{k+1}}(\underline{\mathcal{B}}(\mathbf{x}; \underline{\mathcal{W}}_{1:(k+1)}, \underline{\mathcal{A}}_{k+1})),$$

$$\underline{\arg \min}_{\underline{\mathcal{W}}_{k+1}, \underline{\mathcal{A}}_{k+1}, \underline{C_{k+1}}} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{k+1}} l_{k+1}(M_{k+1}(\mathbf{x}), \mathbf{y})$$

## Conditional Convolution



## ■ What have they done?

$$\mathbb{M}_{\mathbb{P}_k} \approx \prod_k \mathcal{L}_{R\mathcal{N}}(\underline{\mu}_k, \underline{\sigma}_k^2)$$

$$\mathcal{L}_R(x_k; \beta) = \frac{1}{1 + \exp(-(\beta(x_k - 0.5)))}$$

## ■ What did they claim?

Our proposed approach builds on the following **Optimal Overlap Hypothesis**: *For a strictly continually trained deep neural network, catastrophic forgetting and remembering can be minimized, without additional memory or data, by learning optimal representational overlap, such that the representational overlap is reduced for unrelated tasks and increased for tasks that are similar.*



## ■ What have they said?

- Understanding Catastrophic Forgetting

Task  $\log \mathcal{P}(\theta_1|D_1) = \overset{\text{Loss}}{\log \mathcal{P}(D_1|\theta_1)} + \log \mathcal{P}(\theta_1) - \log \mathcal{P}(D_1).$

$$\begin{aligned} \log \mathcal{P}(\theta_{1:2}|D_{1:2}) &= \log \mathcal{P}(D_2|\theta_2) + \overset{\text{Prior info}}{\log \mathcal{P}(\theta_1|D_1)} \\ &\quad - \log \mathcal{P}(D_2) \\ &= \log \mathcal{P}(D_2|\theta_2) + \log \mathcal{P}(D_1|\theta_1) \\ &\quad + \log \mathcal{P}(\theta_1) - \log \mathcal{P}(D_1) - \log \mathcal{P}(D_2). \end{aligned}$$

- How is it proved?

$$\mathcal{P}(\theta_1, \mathbb{M}_{P_1}|D_1) \propto \mathcal{P}(D_1|\theta_{\mathbb{M}_{P_1}})\mathcal{P}(\theta_{\mathbb{M}_{P_1}})$$

$$\begin{aligned} \mathcal{P}(\theta_{1:2}, \mathbb{M}_{P_2}|D_{1:2}) &\propto \mathcal{P}(D_2|\theta_{\mathbb{M}_{P_2}})\mathcal{P}(\theta_1, \mathbb{M}_{P_1}|D_1) \\ &\propto \mathcal{P}(D_2|\theta_{\mathbb{M}_{P_2}})\mathcal{P}(\theta_2'') \end{aligned}$$

## ■ Supervised Continual Learning

---

### Algorithm 1 *RMN* Supervised Continual Learning

---

- 1: **Input:** data  $x$ , ground truth  $y$  for  $n$  tasks, prune parameter  $\mu$ , corresponding task labels  $i$  paired with all  $x$
  - 2: **Given:** parameters  $\mathbf{W}$  & initialized relevance mappings  $\mathbb{M}_{\mathcal{P}}$
  - 3: **for** each task  $i$  **do**
  - 4:    $f(x_i; \mathbf{W}, \mathbb{M}_{\mathcal{P}_i}) \Rightarrow \hat{y}_i = \sigma((\mathbf{W} \odot \mathbb{M}_{\mathcal{P}_i}) \odot x_i)$
  - 5:   Compute Loss :  $L(\hat{y}_i, y_i)$
  - 6:   Optional: Add Sparsity Loss :  $L(\hat{y}_i, y_i) + (\mathbb{M}_{\mathcal{P}_i})_{l_0}$
  - 7:   Backpropagate and optimize
  - 8:   Prune  $\mathbb{M}_{\mathcal{P}} \leq \mu$  only.
  - 9:   Stabilize (fix) parameters in  $f$  where  $\mathbb{M}_{\mathcal{P}} = 1$
  - 10: **end for**
  - 11: **Inference:** For data  $x$  and ground-truth task label  $i$ :
  - 12: **Output:**  $f(x, i; \mathbf{W})$
-



# Optimal Relevance

## ■ Unsupervised

### Algorithm 2 RMN Unsupervised Continual Learning

- 1: **Input:** data  $x$ , ground truth  $y$ , prune parameter  $\mu$
- 2: **Given:** parameters  $\mathbf{W}$ ,  $\mathbb{M}_{P_{est\_j}}$  with  $est\_j = 0$ , Task Switch Detection Method  $TSD$
- 3: **for** each task  $j$  **do**
- 4:   Filter input  $x$  on  $f(x; \mathbf{W}, \mathbb{M}_{P_{0\dots j-1}})$
- 5:    $f(x; \mathbf{W}, \mathbb{M}_{P_{est\_j}}) \Rightarrow \hat{y}$
- 6:   Compute Loss :  $L(\hat{y}, y)$                    a Relevance modified Welsh's
- 7:   **if**  $TSD(x)$  is True **then**                   **t-test** on the KL divergence
- 8:      $est\_j ++$
- 9:     Add  $\mathbb{M}_{P_{est\_j}}$  to learn-able parameter list
- 10:     $f(x; \mathbf{W}, \mathbb{M}_{P_{est\_j}}) \Rightarrow \hat{y}$
- 11:    Re-Compute Loss :  $L(\hat{y}, y)$
- 12:    **end if**
- 13:    Backpropagate and optimize
- 14:    Sample  $x_g$  from standard Gaussian distribution with same shape as  $x$
- 15:     $f(x_g; \mathbf{W}, \mathbb{M}_{P_{est\_j}}) \Rightarrow \hat{y}$
- 16:    Compute Loss :  $\|\hat{y} - 0\|_2^2$
- 17:    Backpropagate and optimize
- 18:    Prune  $\mathbb{M}_P \leq \mu$  only.
- 19:    Stabilize (fix) parameters in  $f$  where  $\mathbb{M}_P \approx 1$
- 20: **end for**
- 21: **Inference:** For data  $x$ :
- 22: **Output:**  $max_k f(x, k; W)$

## ■ Is that real?

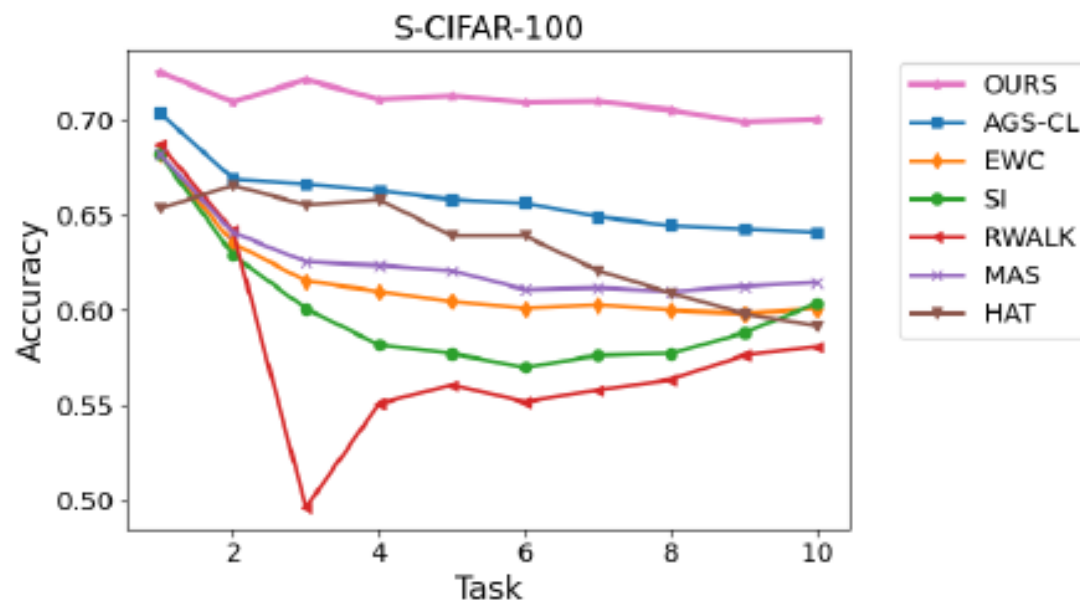


Figure 1. Average accuracy results on CIFAR-100 (10 tasks)

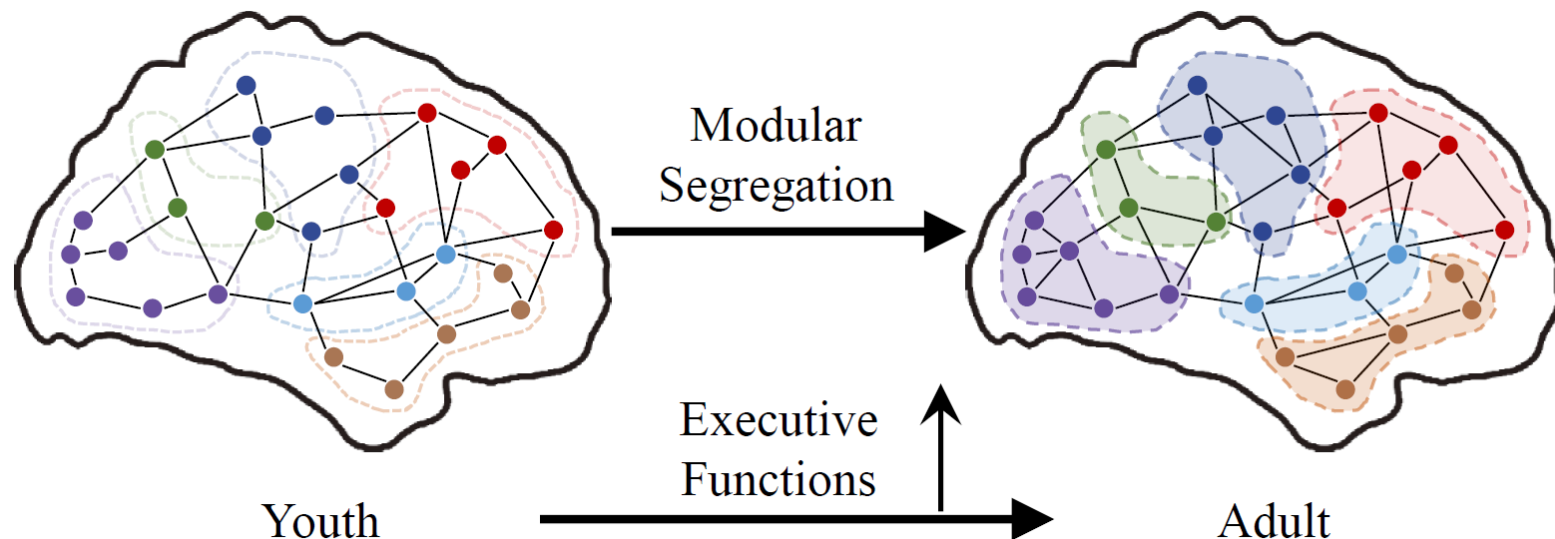
- Based on  $X$  (fc layer or t-test)
- Fix previous params?  $\checkmark$
- Joint optimization?  $\times$



# Our proposals

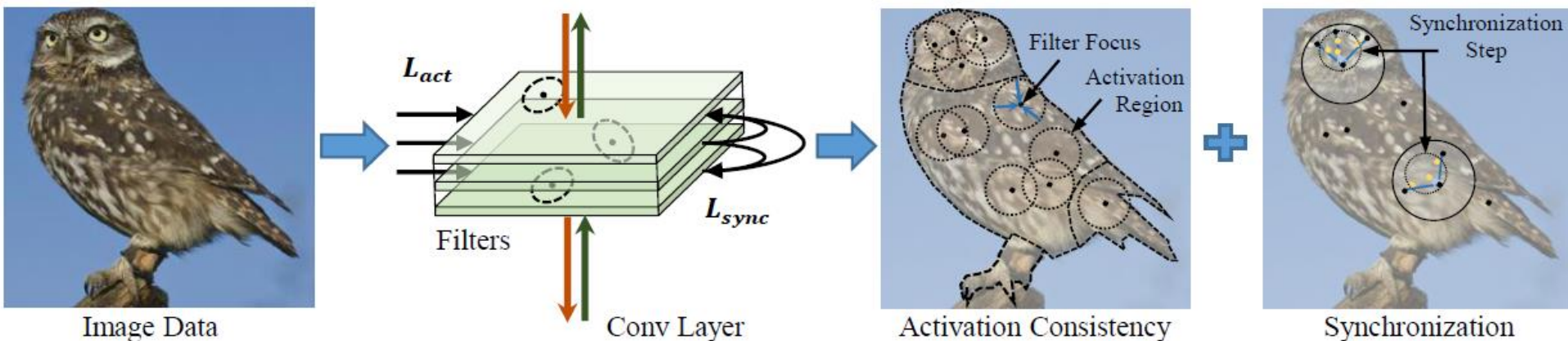
## ■ Motivation

- Interpretation via Synchronization as Localization
- A Feature Map as an Entity  $\rightarrow$  Sync as Modules



## ■ Idea

- Minimize Entropy from several levels
- Sync (unsupervised) for Module Segmentation



## ■ Functional Module Level

- Sync in clustering

$$\frac{d\theta_i}{dt} = \omega_i + \frac{S}{N} \sum_{j=1}^N \sin(\theta_j - \theta_i), (i = 1, \dots, N)$$

$$Nb_\epsilon(x) = \{y \in \mathcal{D} | \text{dist}(y, x) \leq \epsilon\}$$

$$x_i(t+1) = x_i(t) + \frac{1}{|Nb_\epsilon(x(t))|} \cdot \sum_{y \in Nb_\epsilon(x(t))} \sin(y_i(t) - x_i(t))$$

- Sync in mini-batch

$$L_{sync} = \frac{1}{N} \sum_{x_i} \frac{1}{|Nb_\epsilon(x_i)|} \sum_{x_j \in Nb_\epsilon(x_i)} \text{dist}(x_i, x_j)$$

## ■ Sync in mini-batch

$$L_{sync} = \frac{1}{N} \sum_{x_i} \frac{1}{|Nb_{\epsilon}(x_i)|} \sum_{x_j \in Nb_{\epsilon}(x_i)} dist(x_i, x_j)$$

### • Understanding

$$L_{sync} = \frac{1}{N} \sum_{x_i} \frac{1}{|Nb_{\epsilon}(x_i)|} \sum_{x_j \in Nb_{\epsilon}(x_i)} dist(x_i, x_j)$$

Word Embedding!

$$\cong - \sum_{i,j} p(i)p(j|i) \log q(j|i) = H(j|i)$$

$$p(i) = \frac{1}{N}$$

$$q(j|i) = \frac{e^{-dist(x_i, x_j)/\tau_2}}{\sum_j e^{-dist(x_i, x_j)/\tau_2}}$$

$$p(j|i) = \frac{Nb_{\epsilon}(x_i, x_j)}{\sum_j Nb_{\epsilon}(x_i, x_j)} = \begin{cases} \frac{1}{|Nb_{\epsilon}(x_i)|}, & dist(x_i, x_j) < \epsilon \\ 0, & otherwise \end{cases}$$



## ■ Motivation

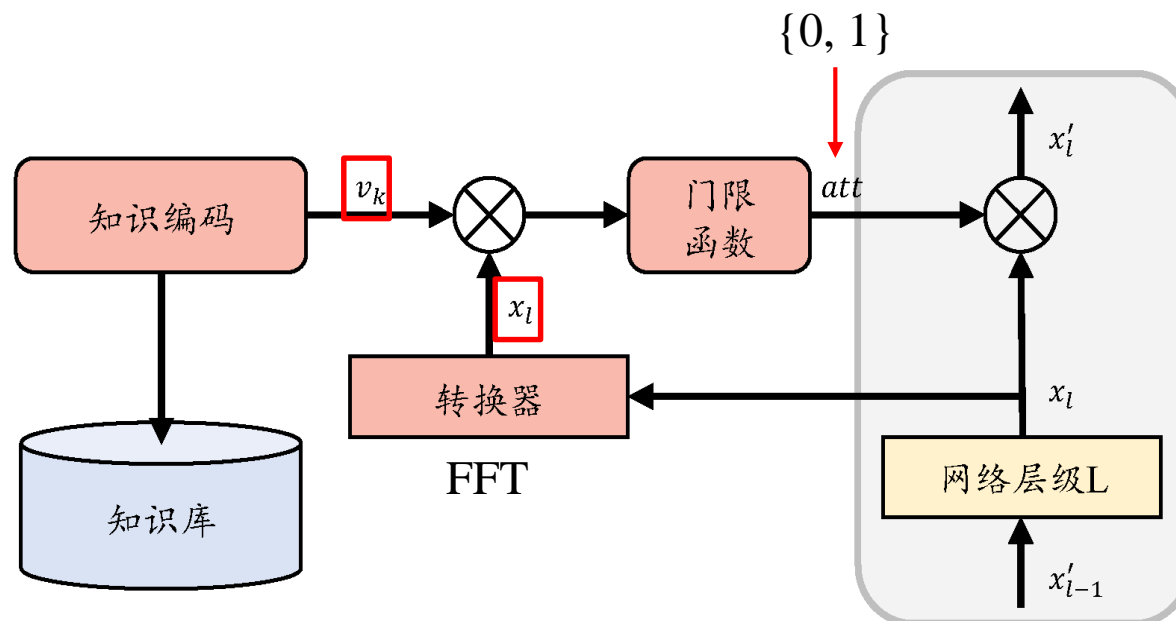
- Functional modules on **the whole structure**
- Unsupervised knowledge embedding
- Knowledge embedding management

## ■ Idea

- Fourier Transform
- Sync to regular representation space
- Attention mechanism

## ■ Attention Structure

- Knowledge embedding as Key vector
- Hard attention -> activate the knowledge or not



## ■ What's the cons?

- 'Soft'-mask
- No absolutely fixed parameter
- Joint optimization

## ■ So, Sync could be a promising way to it?

- **Stable** functional module

## ■ What's the pros?

1. **Constant memory.** To avoid unbounded systems, the consumed memory should be constant w.r.t. the number of tasks or length of the data stream.
6. **Problem agnostic** continual learning is not limited to a specific setting (e.g. only classification).
7. **Adaptive** systems learn from available unlabeled data as well, opening doors for adaptation to specific user data.
8. **No test time oracle** providing the task label should be required for prediction.
9. **Task revisiting** of previously seen tasks should enable enhancement of the corresponding task knowledge.
10. **Graceful forgetting.** Given an unbounded system and infinite stream of data, selective forgetting of trivial information is an important mechanism to achieve balance between stability and plasticity.

## ■ What's the pros?

- Local, Disentangled, Interpretable, Comparable
- Knowledge management for Continual learning
- Multi-source, partial, model-based Transfer learning
- Zero-shot via compositional generalization
- Match with existing knowledge systems
- Interact with human beings
- .....

## ■ Plan

- Exar
- Adap
- Expe
- Gene
- comp
- Thec
- Perf

Oxford 102 Flowers					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	10.0 (-20.3)	5.1 (-17.1)	6.7 (-13.6)	17.3 (0.0)	9.8
Freezing	30.3 (0.0)	39.8 (0.0)	32.0 (0.0)	33.1 (0.0)	33.8
Joint	54.6 (+24.3)	58.9 (+11.5)	57.7 (+4.5)	47.0 (0.0)	54.6
EWC [14]	12.1 (-18.2)	11.6 (-38.1)	9.3 (-24.4)	25.8 (0.0)	14.7
HAT [41]	17.2 (-12.7)	19.3 (-28.5)	28.6 (+1.4)	31.6 (0.0)	24.2
PackNet [26]	32.0 (0.0)	53.7 (0.0)	43.6 (0.0)	37.9 (0.0)	41.8
TFM w/o FN	36.4 (0.0)	54.1 (0.0)	38.6 (0.0)	39.0 (0.0)	42.0
TFM	36.4 (0.0)	53.8 (0.0)	45.5 (0.0)	37.6 (0.0)	<b>43.3</b>

CUBS 200 Birds					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	7.4 (-30.2)	2.6 (-30.0)	29.7 (-3.4)	43.1 (0.0)	20.7
Freezing	37.6 (0.0)	35.1 (0.0)	35.4 (0.0)	38.4 (0.0)	36.6
Joint	48.7 (+11.1)	52.1 (+6.0)	50.7 (+1.5)	51.9 (0.0)	50.8
EWC [14]	16.2 (-21.4)	19.0 (-21.2)	24.2 (-14.0)	41.7 (0.0)	25.3
HAT [41]	18.7 (-1.8)	19.4 (-0.4)	28.5 (-0.6)	31.2 (0.0)	24.4
PackNet [26]	35.3 (0.0)	42.8 (0.0)	44.4 (0.0)	45.9 (0.0)	42.1
TFM w/o FN	42.9 (0.0)	44.1 (0.0)	48.3 (0.0)	49.1 (0.0)	46.1
TFM	42.9 (0.0)	43.1 (0.0)	49.9 (0.0)	48.8 (0.0)	<b>46.2</b>

Stanford 40 Actions					
Method	Task 1	Task 2	Task 3	Task 4	Avg.
Finetuning	24.4 (-10.5)	26.5 (-7.7)	17.6 (-16.8)	28.9 (0.0)	24.4
Freezing	34.9 (0.0)	29.4 (0.0)	30.1 (0.0)	30.5 (0.0)	31.2
Joint	45.7 (+10.8)	40.3 (+4.8)	43.2 (-1.1)	40.2 (0.0)	42.4
EWC [14]	24.2 (-10.7)	28.2 (-2.0)	25.2 (-5.6)	34.3 (0.0)	28.0
HAT [41]	25.7 (-1.0)	25.5 (-2.7)	30.1 (-2.1)	34.4 (0.0)	28.9
PackNet [26]	32.5 (0.0)	32.9 (0.0)	36.7 (0.0)	34.3 (0.0)	34.1
TFM w/o FN	35.3 (0.0)	38.3 (0.0)	39.2 (0.0)	38.0 (0.0)	37.7
TFM	35.3 (0.0)	37.2 (0.0)	42.0 (0.0)	37.2 (0.0)	<b>38.0</b>

g,  
, etc.)

tional)

## ■ Plan

- Exam existing CL methods in the same setting
- Adapt from FcaNet, add HAT and Sync
- Experiments of CL and Interpretability(?)
- General Applications (transfer/zero-shot learning, compact existing networks & efficient networks, etc.)
- Theoretical guarantee of accuracy boundary (optional)
- Perform as a conference paper

Thank  
you

