

#### **Forget-free Continual Learning with Winning Subnetworks**

Haeyong Kang<sup>\*1</sup> Rusty John Lloyd Mina<sup>\*1</sup> Sultan Rizky Hikmawan Madjid<sup>1</sup> Jaehong Yoon<sup>1</sup> Mark Hasegawa-Johnson<sup>2</sup> Sung Ju Hwang<sup>13</sup> Chang D. Yoo<sup>1</sup>

ICML 2022

# **Incremental learning**





### Incremental learning (Continual

learning, lifelong learning):

- 独立同分布(i.i.d.)假设
- 稳定性-可塑性困境

Task incremental learning:  $\{\mathcal{Y}^{(t)}\} \neq \{\mathcal{Y}^{(t+1)}\}$ Task known & Task unknown Methods:

- Replay
- Regularization-based
- Parameter isolation

# **Motivation**

Limitation: memory usage increases as new tasks arrive.

- Replay: additional room to store the replay buffer.
- Parameter isolation: addition room for models to deal with new task.

How can we build a memory-efficient CL model that does not exceed the capacity of the backbone network?

#### Lottery Ticket Hypothesis(LTH):

- Model-purning method
- Over-parameter
- Winner ticket: subnetworks
- "Subnetworks of large network—when trained in isolation—converge in a comparable number of iterations to comparable accuracy"



[1] The Lottery Ticket Hypothesis : Finding Sparse Trainable Neural Networks. ICLR 2019 (best paper)

# **Motivation**



*Figure 1.* **Concept Comparison:** (a) Piggyback (Mallya et al., 2018), and SupSup (Wortsman et al., 2020) find the optimal binary mask on a fixed backbone network a given task (b) PackNet (Mallya & Lazebnik, 2018) and CLNP (Golkar et al., 2019) forces the model to reuse all features and weights from previous subnetworks which causes bias in the transfer of knowledge (c) APD (Yoon et al., 2020) selectively reuse and dynamically expand the dense network (d) Our WSN selectively reuse and dynamically expand subnetworks within a dense network. Green edges are reused weights

# Method



Figure 2. An illustration of Winning SubNetworks (WSN): (a) The top-c% weights  $\hat{\theta}_{t-1}$  at prior task are obtained, (b) In the forward pass of a new task, WSN reuses weights selected from prior tasks, (c) In the backward pass, WSN updates only non-used weights, and (d) after several iterations of (b) and (c), we acquire again the top-c% weights  $\hat{\theta}_t$  including subsets of reused weights. for the new task.

# **Overview & Details**

#### Algorithm 1 Winning SubNetworks (WSN)

- input  $\{\mathcal{D}_t\}_{t=1}^{\mathcal{T}}$ , model weights  $\theta$ , score weights s, binary mask  $\mathbf{M}_0 = \mathbf{0}^{|\boldsymbol{\theta}|}$ , layer-wise capacity c 1: Randomly initialize  $\theta$  and s. 2: for task  $t = 1, ..., \mathcal{T}$  do for batch  $\mathbf{b}_t \sim \mathcal{D}_t$  do 3: Obtain mask  $\mathbf{m}_t$  of the top-c% scores s at each layer 4: 5: Compute  $\mathcal{L}(\boldsymbol{\theta} \odot \mathbf{m}_t; \mathbf{b}_t)$ 6:  $\theta \leftarrow \theta - \eta \left( \frac{\partial \mathcal{L}}{\partial \theta} \odot (\mathbf{1} - \mathbf{M}_{t-1}) \right)$   $\triangleright$  Weight update 7:  $\mathbf{s} \leftarrow \mathbf{s} - \eta(\frac{\partial \mathcal{L}}{\partial \mathbf{s}})$ Weight score update 8: end for 9:  $\mathbf{M}_t \leftarrow \mathbf{M}_{t-1} \lor \mathbf{m}_t$ ▷ Accumulate binary mask 10: end for
- Freeze parameters selected by pervious steps.
- Compress binary mask M with **Huffman encoder**.

```
output = model(data)
loss = criterion(output, target)
loss.backward()
```

getattr(getattr(model, module\_name), module\_attr).grad[consolidated\_masks[ke y] == 1] = 0

optimizer.step()

## **Experiments**

Method	Permuted MNIST			5 Datasets			Omniglot Rotation		
	ACC (%)	CAP (%)	BWT	ACC (%)	CAP (%)	BWT	ACC (%)	CAP (%)	BWT
STL	$97.37 (\pm 0.01)$	1,000.0	-	$93.44 (\pm 0.12)$	500.0	-	$82.13 (\pm 0.08)^*$	10,000.0	-
FINETUNE	$78.22 (\pm 0.84)$	100.0	$-0.21 (\pm 0.01)$	$80.06 (\pm 0.74)$	100.0	$-0.17 (\pm 0.01)$	44.48 (± 1.68)	100.0	$-0.45 (\pm 0.02)$
EWC (Kirkpatrick et al., 2017)	92.01 (± 0.56)	100.0	$-0.03 (\pm 0.00)$	$88.64 (\pm 0.26)^{\dagger}$	$100.0^{+}$	$-0.04 \ (\pm \ 0.01)^{\dagger}$	68.66 (± 1.92)*	$100.0^{*}$	-
HAT (Serrà et al., 2018)	-	-	-	$91.32~(\pm 0.18)^{\dagger}$	$100.0^{\dagger}$	$-0.03 (\pm 0.00)^{\dagger}$	-	-	-
GPM (Saha et al., 2021)	$94.96 (\pm 0.07)$	100.0	$-0.02 (\pm 0.01)$	$91.22~(\pm 0.20)^{\dagger}$	100.0	$-0.01~(\pm 0.00)^{\dagger}$	$85.24 (\pm 0.37)$	100.0	$-0.01 (\pm 0.00)$
PackNet (Mallya & Lazebnik, 2018)	96.37 (± 0.04)	96.38	0.0	92.81 (± 0.12)	82.86	0.0	30.70 (± 1.50)	399.2	0.0
SupSup (Wortsman et al., 2020)	96.31 (± 0.09)	$122.89 (\pm 0.07)$	0.0	$93.28 (\pm 0.21)$	$104.27 (\pm 0.21)$	0.0	$58.14 (\pm 2.42)$	$407.12 (\pm 0.17)$	0.0
WSN, $c = 0.03$	94.84 (± 0.11)	<b>19.87</b> (± 0.16)	0.0	90.57 (± 0.65)	$12.11 (\pm 0.06)$	0.0	80.68 (± 2.60)	75.87 (± 1.24)	0.0
WSN, $c = 0.05$	$95.65 (\pm 0.03)$	$26.49 (\pm 0.16)$	0.0	91.61 (± 0.21)	17.26 (± 0.25)	0.0	$87.28 (\pm 0.72)$	79.85 (± 1.19)	0.0
WSN, $c = 0.1$	$96.14 (\pm 0.03)$	$40.41 (\pm 0.54)$	0.0	92.67 (± 0.12)	$28.01 (\pm 0.28)$	0.0	83.10 (± 1.56)	83.08 (± 1.61)	0.0
WSN, $c = 0.3$	<b>96.41</b> ( $\pm$ 0.07)	77.73 (± 0.36)	0.0	93.22 (± 0.32)	$62.30 (\pm 0.69)$	0.0	81.89 (± 1.15)	$102.2 (\pm 0.89)$	0.0
WSN, $c = 0.5$	$96.24 (\pm 0.11)$	98.10 (± 0.25)	0.0	$93.41 (\pm 0.13)$	86.10 (± 0.57)	0.0	79.80 (± 2.16)	121.2 (± 0.50)	0.0
MTL	$96.70 \ (\pm \ 0.02)^\dagger$	100.0	-	$91.54~(\pm 0.28)^{\dagger}$	100.0	-	81.23 (± 0.52)	100.0	-

- ACC: accuracy
- CAP: capacity of parameters(%)
- BWT: the ACC drop during continual learning(i.e., forgetting)

### **Experiments**

*Table 2.* Performance comparisons of the proposed method and other state-of-the-art including baselines - PackNet (Mallya & Lazebnik, 2018) and SupSup (Wortsman et al., 2020) - on various benchmark datasets. We report the mean and standard deviation of the average accuracy (ACC), average capacity (CAP), and average backward transfer (BWT) across 5 independent runs with 5 seeds under the same experimental setup (Deng et al., 2021). The best results are highlighted in bold. Also, † denotes results reported from Deng et al. (2021).

Method	CIFAR-100 Split			CIFAR-100 Superclass			TinyImageNet		
	ACC (%)	CAP (%)	BWT (%)	ACC (%)	CAP (%)	BWT (%)	ACC (%)	CAP (%)	BWT (%)
La-MaML (Gupta et al., 2020) GPM (Saha et al., 2021) FS-DGPM (Deng et al., 2021)	$\begin{array}{l} 71.37 \ (\pm \ 0.67)^\dagger \\ 73.18 \ (\pm \ 0.52)^\dagger \\ 74.33 \ (\pm \ 0.31)^\dagger \end{array}$	100.0 100.0 100.0	$\begin{array}{l} \textbf{-5.39} \ (\pm \ 0.53)^\dagger \\ \textbf{-1.17} \ (\pm \ 0.27)^\dagger \\ \textbf{-2.71} \ (\pm \ 0.17)^\dagger \end{array}$	$\begin{array}{l} 54.44\ (\pm\ 1.36)^{\dagger}\\ 57.33\ (\pm\ 0.37)^{\dagger}\\ 58.81\ (\pm\ 0.34)^{\dagger}\end{array}$	100.0 100.0 100.0	$\begin{array}{l} \textbf{-6.65} \ (\pm \ 0.85)^\dagger \\ \textbf{-0.37} \ (\pm \ 0.12)^\dagger \\ \textbf{-2.97} \ (\pm \ 0.35)^\dagger \end{array}$	$\begin{array}{l} 66.90 \ (\pm \ 1.65)^\dagger \\ 67.39 \ (\pm \ 0.47)^\dagger \\ 70.41 \ (\pm \ 1.30)^\dagger \end{array}$	100.0 100.0 100.0	$\begin{array}{c} \textbf{-9.13} \ (\pm \ 0.90)^\dagger \\ \textbf{1.45} \ (\pm \ \textbf{0.22})^\dagger \\ \textbf{-2.11} \ (\pm \ 0.84)^\dagger \end{array}$
PackNet (Mallya & Lazebnik, 2018) SupSup (Wortsman et al., 2020)	$\begin{array}{c} 72.39 \ (\pm \ 0.37) \\ 75.47 \ (\pm \ 0.30) \end{array}$	$\begin{array}{c} 96.38 \ (\pm \ 0.00) \\ 129.00 \ (\pm \ 0.03) \end{array}$	0.0 0.0	$\begin{array}{c} 58.78 \ (\pm \ 0.52) \\ 61.70 \ (\pm \ 0.31) \end{array}$	$\begin{array}{c} 126.65 \ (\pm \ 0.00) \\ 162.49 \ (\pm \ 0.00) \end{array}$	0.0 0.0	$\begin{array}{c} 55.46 \ (\pm \ 1.22) \\ 59.60 \ (\pm \ 1.05) \end{array}$	$\begin{array}{c} 188.67 \ (\pm \ 0.00) \\ 214.52 \ (\pm \ 0.89) \end{array}$	0.0 0.0
WSN, $c = 0.03$ WSN, $c = 0.05$ WSN, $c = 0.1$ WSN, $c = 0.3$ WSN, $c = 0.5$	$\begin{array}{l} 70.65 \ (\pm \ 0.36) \\ 72.44 \ (\pm \ 0.27) \\ 74.55 \ (\pm \ 0.47) \\ 75.98 \ (\pm \ 0.68) \\ \textbf{76.38} \ (\pm \ \textbf{0.34}) \end{array}$	$\begin{array}{c} \textbf{18.56} \ (\pm \ \textbf{0.25}) \\ 25.09 \ (\pm \ \textbf{0.42}) \\ 39.87 \ (\pm \ \textbf{0.62}) \\ \textbf{80.26} \ (\pm \ \textbf{1.53}) \\ \textbf{99.13} \ (\pm \ \textbf{0.48}) \end{array}$	0.0 0.0 0.0 0.0 0.0	$\begin{array}{l} 54.99\ (\pm\ 0.71)\\ 57.99\ (\pm\ 1.34)\\ 60.45\ (\pm\ 0.37)\\ 61.47\ (\pm\ 0.30)\\ \textbf{61.79}\ (\pm\ \textbf{0.23}) \end{array}$	$\begin{array}{c} \textbf{22.30} \ (\pm \ \textbf{0.22}) \\ 27.37 \ (\pm \ \textbf{0.33}) \\ 38.55 \ (\pm \ \textbf{0.20}) \\ \textbf{63.47} \ (\pm \ \textbf{1.33}) \\ 80.93 \ (\pm \ \textbf{1.58}) \end{array}$	0.0 0.0 0.0 0.0 0.0	$\begin{array}{l} 68.72 \ (\pm \ 1.63) \\ 71.22 \ (\pm \ 0.94) \\ \textbf{71.96} \ (\pm \ \textbf{1.41}) \\ 70.92 \ (\pm \ 1.37) \\ 69.06 \ (\pm \ 0.82) \end{array}$	$\begin{array}{c} \textbf{37.19} \ (\pm \ \textbf{0.21}) \\ 41.98 \ (\pm \ \textbf{0.52}) \\ 48.65 \ (\pm \ \textbf{3.03}) \\ 73.44 \ (\pm \ \textbf{2.35}) \\ 92.03 \ (\pm \ \textbf{1.80}) \end{array}$	0.0 0.0 0.0 0.0 0.0
Multitask	$79.75~(\pm 0.38)^{\dagger}$	100.0	-	$61.00 (\pm 0.20)^{\dagger}$	100.0	-	$77.10 (\pm 1.06)^{\dagger}$	100.0	-

## **Experiment**



- "reused for all tasks" keeps in a low level, while "reused per tasks" is more.
- WSN tends uses diverse weights to solve the sequential tasks within all used weights.



Paramenter Correlations between tasks.