

ASIA SOUTH PACIFIC DESIGN AUTOMATION CONFERENCE

Date: January 20-23, Tokyo

### LightCL: Compact Continual Learning with Low Memory Footprint For Edge Device

#### Zeqing Wang, Fei Cheng\*, Kangye Ji, Bohu Huang

School of Computer Science and Technology Xidian University

2025/1/21



#### Outline

- Background
- Challenge
- Method
- Experimental Results
- Conclusion

# 1 Background

Edge devices are everywhere







#### **Continual Learning !**

Background



# 2 Challenge

#### Catastrophic Forgetting

- Model training on a new task tends to "forget" the knowledge learned from previous tasks
- Key idea is to make trade-off between learning plasticity and memory stability for gaining generalizability



A Comprehensive Survey of Forgetting in Deep Learning Beyond Continual Learning [Wang et al., TPAMI 2024]

## 2 Challenge

- Limited Resource in Edge Device
  - Limited computational power and memory
  - Training consumes more than inference
  - Memory has become the primary bottleneck in AI applications



AI and Memory Wall [GHOLAMI et al., Micro 2024]

#### Motivation

• CL process has the potential to leverage previous knowledge when training on

new tasks **Area Predundancy** in training



CL Setting



- Analysis of Generalizability
  - Memory stability (MS) and Learning plasticity (LP) are two different characteristics of generalizability
  - MS denotes the loss of previous knowledge
  - LP denotes the adaptation to new knowledge

$$\theta^{new_i} = \begin{cases} \theta_j^{new_i} = \theta_j^T & \text{if } j \neq i, \\ \theta_i^{new_i} = \theta_i^{T-1} & \text{otherwise.} \end{cases}$$
(1) 
$$\begin{cases} MS = a_{T-1}(\theta^{new_i}) - a_{T-1}(\theta^T), \\ LP = a_T(\theta^{new_i}) - a_T(\theta^T). \end{cases}$$
(2)

A Comprehensive Survey of Forgetting in Deep Learning Beyond Continual Learning [Wang et al., TPAMI 2024]

#### Analysis of Generalizability

• During CL, lower and middle layers have strong generalizability, and deeper

layers have less generalizability



Proposed Architecture



#### Maintain Generalizability

- Freezing lower and middle layers to maintain generalizability during CL
- Reduce resource consumption and maintain generalized knowledge

#### • Memorize Feature Patterns

- Feature maps consume lower consumption in deeper layers
- The output mode of feature mapping is not related to sample information, but to the location of feature map
- Regulate features important for previous tasks during learning new task without accessing to previous samples

- Memorize Feature Patterns
  - Select references
    - Recognize important positions by *l*1-norm with certain ratio
    - All important positions are stored in  $I = \{ (i, j) \}$
    - Select a few samples from the new task to run it on the model and get feature map standards, called *FS*, as references
  - Regulation

$$\mathcal{L}_f(FM, FS) = \beta \sum_{(i,j) \in I} (FM_{i,j} - FS_{i,j})^2.$$

Whole loss function

$$\mathcal{L}(\theta^t) = \mathcal{L}_c(f(\theta^t, x_t), y_t) + \mathcal{L}_f(FM, FS).$$

#### Scenarios

- Task incremental learning (TIL) and Class incremental learning (CIL) scenarios
- Models
  - ResNet-18 and ResNet-50 (pretrained on the ImageNet32×32)

#### Datasets

- Split CIFAR-10 (5 tasks) and Split Tiny-ImageNet (10 tasks) datasets
- Metrics

$$AA = \frac{1}{T} \sum_{i=1}^{T} a_i \left( \theta^T \right)$$

Average Accuracy

Forward & Backward Training FLOPs Parameters Gradients Feature Maps Memory Footprint

#### Main Result

Method	Buffer size	Sparsity	Split CIFAR-10				Split TinyImageNet			
			CIL(↑)	TIL(↑)	$  FLOPs \\ \times 10^{15} (\downarrow) $	Mem (MB)(↓)	CIL(↑)	TIL(↑)	$FLOPs \times 10^{16} (\downarrow)$	Mem (MB)(↓)
JOINT SGD	/	0.0	$94.67 \pm 0.15$ $34.21 \pm 3.06$	$99.04{\scriptstyle\pm0.07}\\86.19{\scriptstyle\pm1.32}$	37.5 7.5	153.4 153.4	${}^{61.38 \pm 0.12}_{10.63 \pm 0.36}$	$\begin{array}{c} 83.14{\pm}0.52\\ 34.24{\pm}0.84\end{array}$	120 12	357.9 357.9
EWC [14] PNN [25]	/	0.0	40.42±3.81 /	$\begin{array}{c} 89.77 {\pm} 3.01 \\ 94.62 {\pm} 0.52 \end{array}$	7.5 10.5	196.0 323.9	11.17±1.55 /	$34.70 \pm 1.55$ $73.94 \pm 3.56$	12 26.4	400.5 528.4
FDR [2] ER [5] DER++ [3]	15	0.0	25.64±2.52 54.44±2.82 40.75±2.55	85.54±3.08 92.25±1.92 85.68±4.35	11 11 14.5	185.4 185.4 217.3	$\begin{array}{c} 10.80 {\pm} 0.71 \\ 11.82 {\pm} 0.31 \\ 7.98 {\pm} 1.84 \end{array}$	$\begin{array}{c} 34.42{\pm}0.38\\ 36.76{\pm}0.57\\ 35.88{\pm}0.41\end{array}$	17.7 17.7 23.3	485.7 485.7 613.5
SparCLer [30] SparCLder++ [30]	15	0.9	$\begin{array}{c} 36.98 {\pm} 2.42 \\ 35.58 {\pm} 0.84 \end{array}$	$\begin{array}{c} 85.33 {\pm} 2.49 \\ 85.28 {\pm} 1.00 \end{array}$	1.1 1.5	106.6 140.6	$\frac{10.01 \pm 0.55}{9.12 \pm 1.02}$	$\begin{array}{c} 33.01 {\pm} 0.78 \\ 31.86 {\pm} 1.77 \end{array}$	1.8 2.3	407.0 536.8
LightCL	15	0.0 0.9	${\begin{array}{c} 55.53 \pm 0.35 \\ 38.31 \pm 3.93 \end{array}}$	$96.28{\scriptstyle\pm1.08}\atop86.88{\scriptstyle\pm1.69}$	5.4 2.3	92.8 25.0	${\begin{array}{c} 12.79 \pm 0.30 \\ 8.95 \pm 0.35 \end{array}}$	$37.39{\scriptstyle\pm 0.63}\\33.94{\scriptstyle\pm 0.52}$	8.6 3.7	$\begin{pmatrix}128.0\\87.1\end{pmatrix}$

'/' indicates the corresponding item is not needed (Buffer is not required in JOINT, SGD, EWC, and PNN) or cannot run (PNN cannot run under CIL settings). FLOPs and Mem are calculated theoretically with the precision to one decimal place.

#### reduce at most 6.16×

#### Other Results

- Training from scratch
- First task as pretraining process and LightCL is conducted on the subsequent tasks
- Still outstanding in ResNet-18 and ResNet-50

Method	Split Cl	FAR-10	Split TinyImageNet								
	$ $ CIL( $\uparrow$ )	TIL(↑)	$\operatorname{CIL}(\uparrow)$	TIL(↑)							
ResNet-18											
JOINT	92.20±0.15	98.31±0.12	59.99±0.19	$82.04 \pm 0.10$							
SGD	$  21.00 \pm 1.61$	$65.78 \pm 3.52$	$7.28 \pm 0.04$	$24.99 \pm 0.27$							
EWC	21.32±0.17	$66.13 \pm 0.17$	$8.58 \pm 0.22$	$24.47 \pm 0.30$							
PNN	/	$95.13 \pm 0.72$	/	$61.88 \pm 1.00$							
FDR	26.19±3.14	$77.61 \pm 0.68$	$8.58 \pm 0.21$	$25.20 \pm 0.29$							
ER	$35.90 \pm 1.53$	$82.40 \pm 1.64$	$8.76 \pm 0.23$	$26.29 \pm 1.18$							
DER++	36.25±1.99	$82.05 \pm 1.63$	$6.53 \pm 1.44$	$27.21 \pm 1.85$							
SparCLer	30.09±1.60	$75.12 \pm 0.63$	7.70±0.13	$25.24 \pm 0.55$							
SparCL <sub>DER++</sub>	$34.36 \pm 1.32$	$78.56 \pm 2.88$	$8.16 \pm 0.58$	$26.53 \pm 0.54$							
LightCL	37.86±2.16	85.69±0.98	9.03±0.16	33.74±0.68							
LightCL <sub>0.9</sub>	$35.72 \pm 2.67$	$83.37{\pm}1.10$	$8.76{\scriptstyle\pm0.63}$	$36.86{\scriptstyle \pm 0.29}$							
ResNet-50											
JOINT	94.66±0.15	98.93±0.09	61.04±0.42	$82.70 \pm 0.42$							
SGD	14.16±0.39	$72.05 \pm 1.52$	$8.55 \pm 0.20$	$29.99 \pm 0.23$							
ER	24.19±1.43	$76.65 \pm 1.50$	$8.81 \pm 0.62$	$31.38 \pm 0.06$							
DER++	$26.42 \pm 1.67$	$75.13 \pm 0.69$	$6.23 \pm 1.23$	$27.22 \pm 2.70$							
LightCL	35.08±0.50	$81.49{\scriptstyle\pm0.92}$	9.23±0.67	33.81±1.47							

Other Results



Training on Jetson Nano B01



#### Ablation Study

- All CL methods can enhance performance with the pre-trained model
- *Maintain Generalizability* can even further enhance performance
- *Memorize Feature Patterns* can further increase accuracy



### **5** Conclusion

- It is the **first** study to propose two new metrics of learning plasticity (LP) and memory stability (**MS**) to quantitatively evaluate generalizability
- Lower and middle layers have more generalizability, while deeper layers hold the opposite.
- We propose a new method, LightCL, to overcome catastrophic forgetting through *Maintain Generalizability* and *Memorize Feature Patterns*
- LightCL shows great improvement in delaying forgetting and memory efficiency

### **Q & A**

### Thank you for your attention!

Presenter: Zeqing Wang

Email: wang.zeqing@stu.xidian.edu.cn



Github