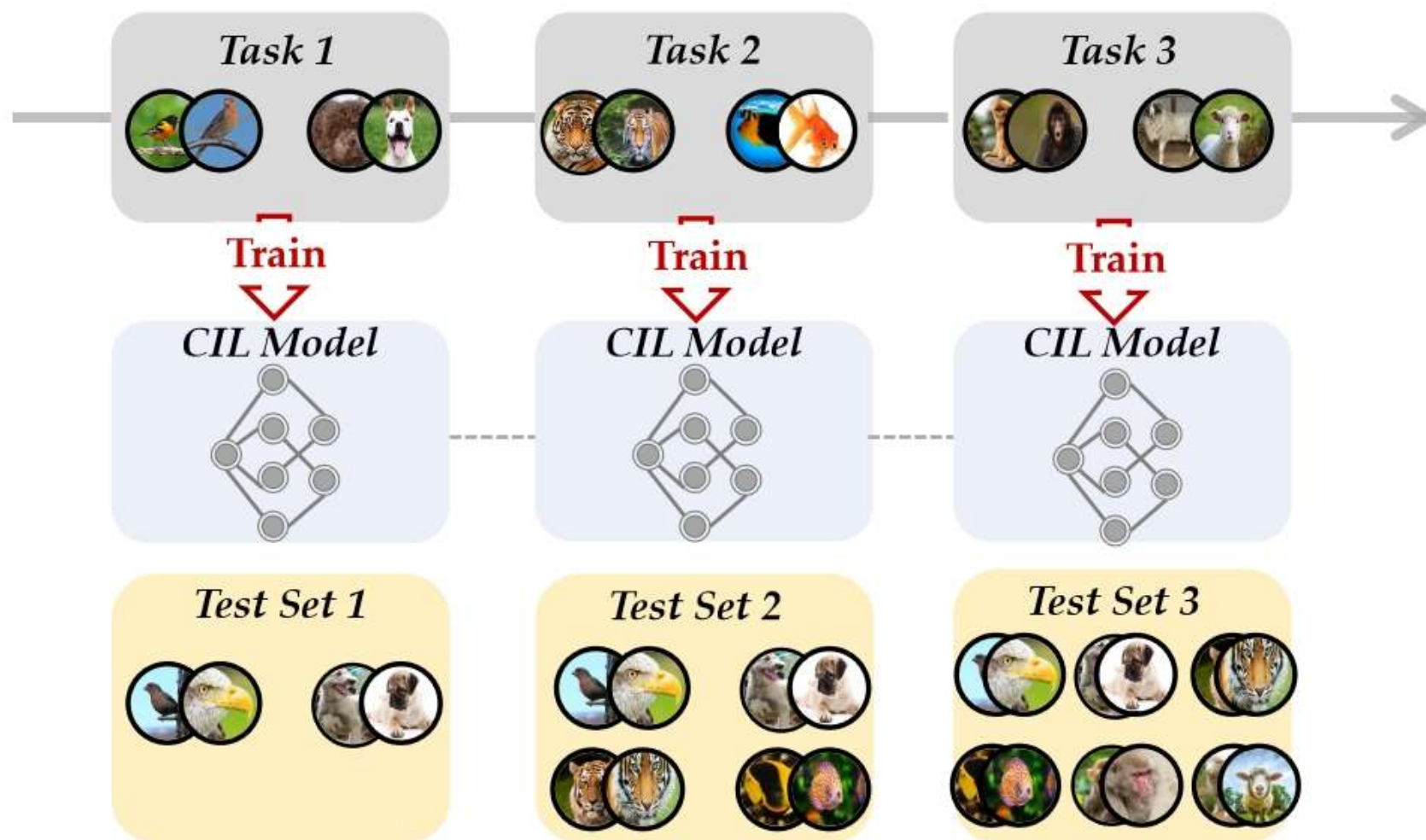


Class-Incremental Learning

Introduction



Catastrophic forgetting

- Data-Centric Class-Incremental Learning
- Model-Centric Class-Incremental Learning
- Algorithm-Centric Class-Incremental Learning

Preliminaries

- There is a sequence of B training tasks

$$\{\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^B\}$$

- b-th incremental step with n_b training instances

$$\mathcal{D}^b = \left\{ (\mathbf{x}_i^b, y_i^b) \right\}_{i=1}^{n_b}$$

- Every label space haven't overlap

$$Y_b \cap Y_{b'} = \emptyset \text{ for } b \neq b'.$$

- After each task, the model is evaluated over all seen class.

$$\mathcal{Y}_b = Y_1 \cup \dots \cup Y_b. \quad f^* = \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_t^1 \cup \dots \cup \mathcal{D}_t^b} \mathbb{I}(y \neq f(\mathbf{x}))$$

Exemplars and exemplar set

- Exemplar Set is an extra collection of instances from former tasks

$$\mathcal{E} = \{(\mathbf{x}_j, y_j)\}_{j=1}^M, y_j \in \mathcal{Y}_{b-1}$$

- Exemplar Set Management

- keep a fixed number of exemplars per class

$$R|\mathcal{Y}_b|$$

- saving a fixed number of exemplars

$$\lceil \frac{M}{|\mathcal{Y}_b|} \rceil$$

- Exemplar Selection

- Randomly sample exemplars
 - Herding

$$\mu_y \leftarrow \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i). \quad \|\mu_y - \phi(\mathbf{x}_i)\| \quad \text{top-}\lceil \frac{M}{|\mathcal{Y}_b|} \rceil$$

Data-Centric Class-Incremental Learning

- Data Replay
- Data Regularization
- Discussions about Data-Centric Methods

Data Replay

- Revisit exemplars
- Exemplar selection
 - Uncertainty
 - Data center
- Exemplars set efficiency
 - Low fidelity feature
- Generative replay
 - The quality of generated data
 - GAN also suffers the catastrophic forgetting

$$\mathcal{L} = \sum_{(\mathbf{x}, y) \in (\mathcal{D}^b \cup \mathcal{E})} \ell(f(\mathbf{x}), y).$$

Rainbow Memory: Continual Learning with a Memory of Diverse Samples

Memory Efficient Class-Incremental Learning for Image Classification

Data Regularization

- Optimizing the model for new classes will not hurt former ones
 - GEM

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{(\mathbf{x}, y) \in \mathcal{D}^b} \ell(f(\mathbf{x}), y)$$

$$\text{s.t.} \quad \sum_{(\mathbf{x}_j, y_j) \in \mathcal{E}} \ell(f(\mathbf{x}_j), y_j) \leq \sum_{(\mathbf{x}_j, y_j) \in \mathcal{E}} \ell(f^{b-1}(\mathbf{x}_j), y_j)$$

$$\langle g, g_{old} \rangle := \left\langle \frac{\partial \ell(f(\mathbf{x}), y)}{\partial \theta}, \frac{\partial \ell(f, \mathcal{E})}{\partial \theta} \right\rangle \geq 0,$$

Discussions about Data-Centric Methods

- Data replay may suffer the overfitting problem
- Data-imbalance problem
- Few-shot exemplars VS. Many-shot task
- How to select data as exemplar

Algorithm 1 Diversity-Aware Memory Update

```
1: Input:  $K$  denotes memory size,  $N_t$  denotes the number of seen classes until task  $t$ ,  $\mathcal{D}_t^S$  denotes stream data at task  $t$ ,  $\mathcal{D}_{t-1}^M$  denotes exemplars stored in a episodic memory after task  $t - 1$ .
2: Output:  $\mathcal{D}_t^M$  exemplars after learning task  $t$ .
3:  $\mathcal{D}_t^M = \{\}$  ▷ New exemplars from scratch
4:  $k_c = \text{floor}(K/N_t)$  ▷ Class-balanced sampling
5: for  $c = 1, 2, \dots, N_t$  do
6:    $\mathcal{D}_c = \{(x, y) | y = c, (x, y) \in \mathcal{D}_t^S \cup \mathcal{D}_{t-1}^M\}$ 
7:   Sort  $\mathcal{D}_c$  by  $u(x)$  computed by (4)
8:   for  $j = 1, 2, \dots, k_c$  do
9:      $i = j * |\mathcal{D}_c|/k_c$  ▷  $|\mathcal{D}_c|/k_c$  step-size indexing
10:     $\mathcal{D}_t^M += \mathcal{D}_c[i]$ 
11:   end for
12: end for
```

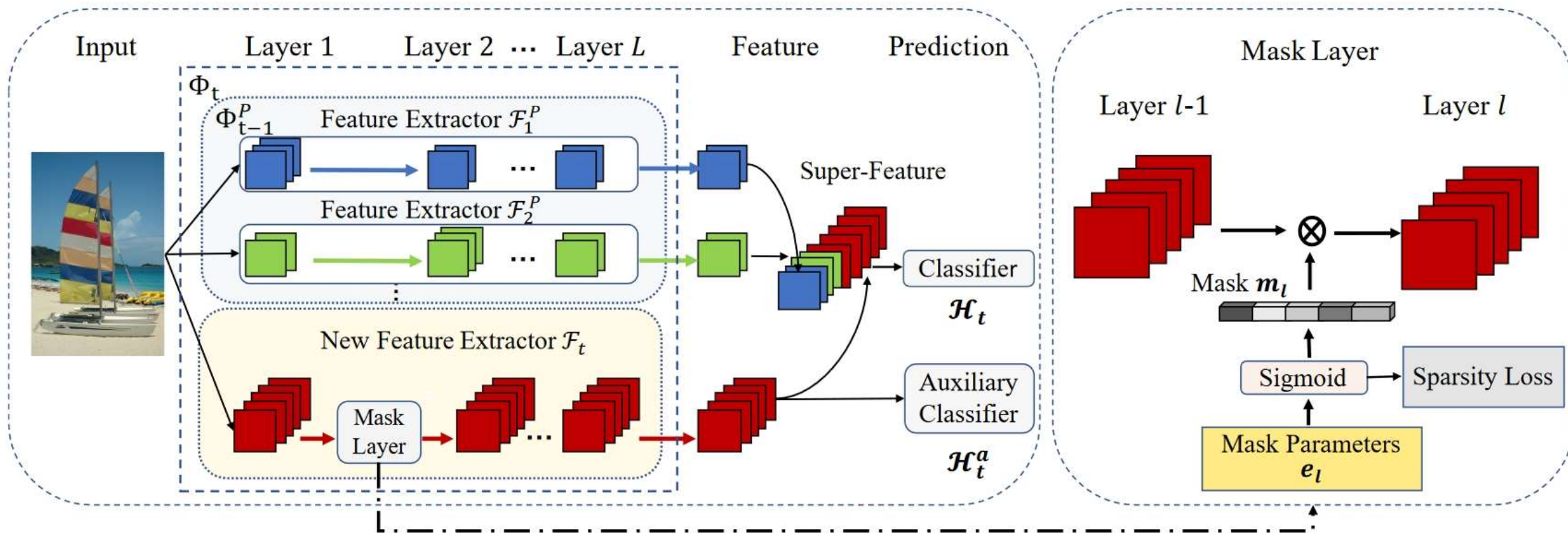
Model-Centric Class-Incremental Learning

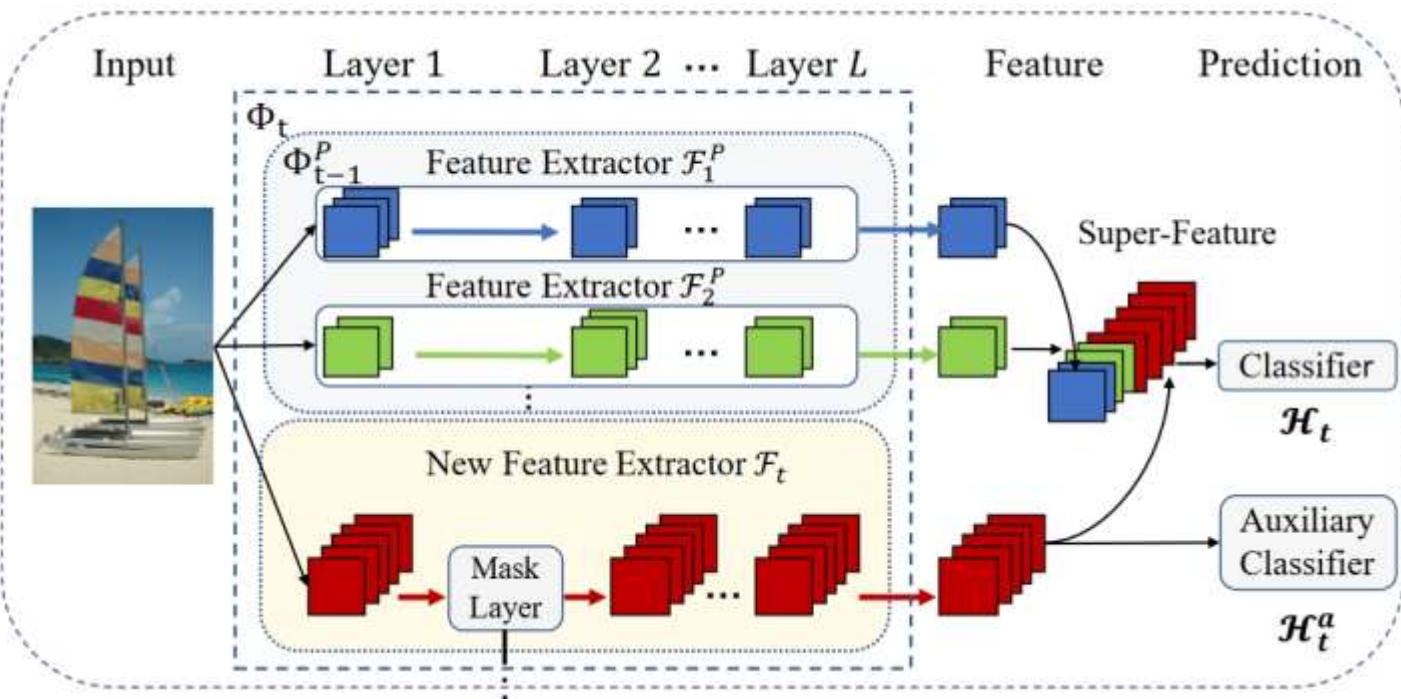
- Dynamic Networks
- Parameter Regularization

Dynamic Networks

- Expands a new backbone when facing new tasks and aggregates the features with a larger FC layer DER

DER





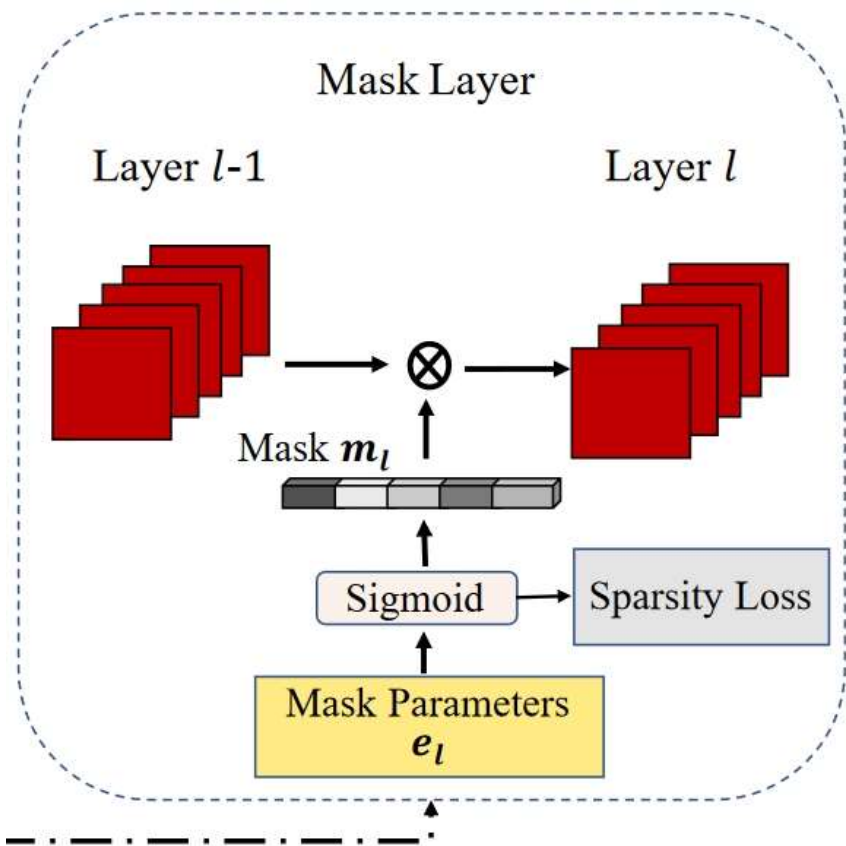
$$\mathcal{L}_{\mathcal{H}_t} = -\frac{1}{|\tilde{\mathcal{D}}_t|} \sum_{i=1}^{|\tilde{\mathcal{D}}_t|} \log(p_{\mathcal{H}_t}(y = y_i | \mathbf{x}_i))$$

$$|\mathcal{Y}_t| + 1$$

$$\mathcal{L}_{\text{ER}} = \mathcal{L}_{\mathcal{H}_t} + \lambda_a \mathcal{L}_{\mathcal{H}_t^a}$$

$$\mathbf{u} = \Phi_t(\mathbf{x}) = [\Phi_{t-1}(\mathbf{x}), \mathcal{F}_t(\mathbf{x})]$$

$$p_{\mathcal{H}_t}(\mathbf{y} | \mathbf{x}) = \text{Softmax}(\mathcal{H}_t(\mathbf{u}))$$



$$\tilde{\mathbf{u}} = \Phi_t^P(\mathbf{x}) = [\mathcal{F}_1^P(\mathbf{x}), \mathcal{F}_2^P(\mathbf{x}), \dots, \phi_t(\mathbf{x})]$$

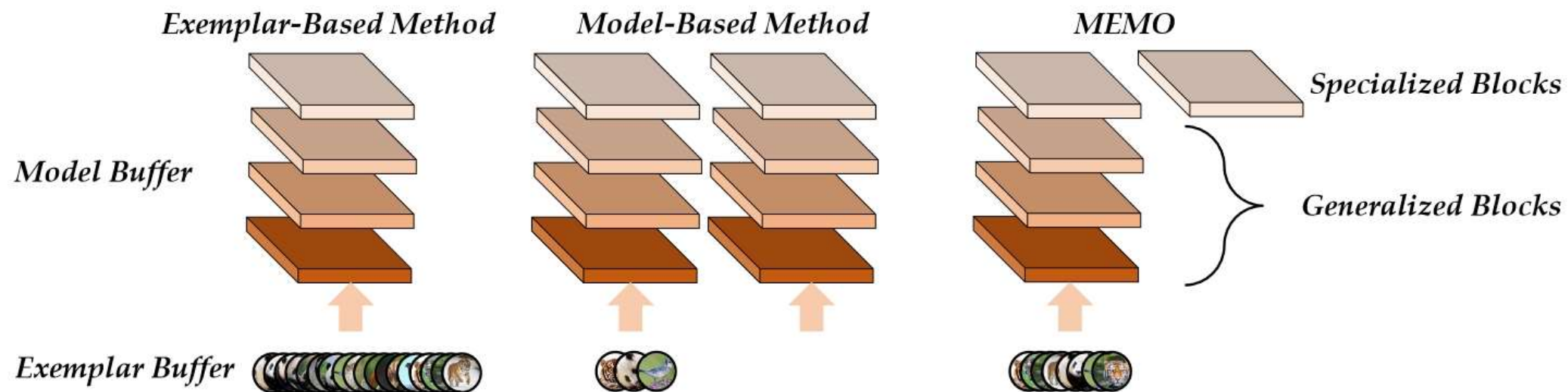
$$\mathbf{m}_l = \sigma(se_l)$$

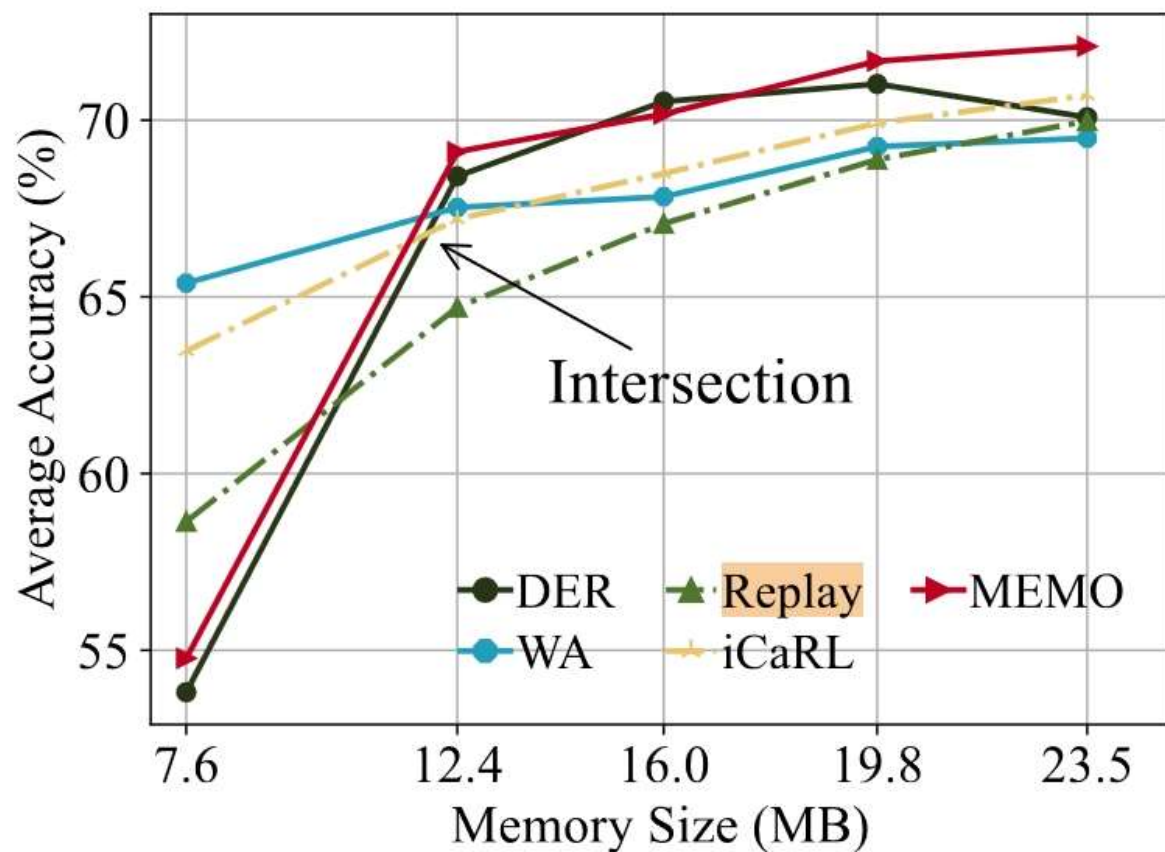
$$s = \frac{1}{s_{\max}} + \left(s_{\max} - \frac{1}{s_{\max}}\right) \frac{b-1}{B-1}$$

$$\mathcal{L}_S = \frac{\sum_{l=1}^L K_l \|\mathbf{m}_{l-1}\|_1 \|\mathbf{m}_l\|_1}{\sum_{l=1}^L K_l c_{l-1} c_l}$$

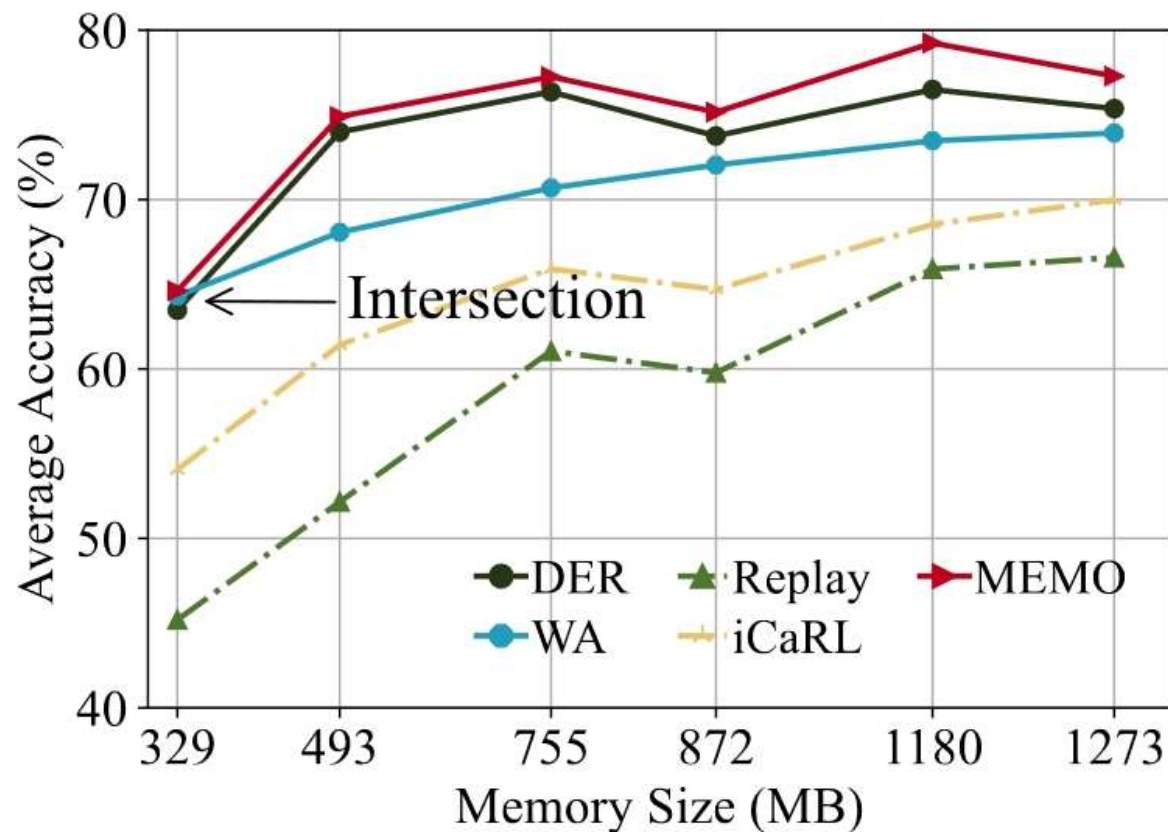
$$\mathcal{L}_{\text{DER}} = \mathcal{L}_{\mathcal{H}_t} + \lambda_a \mathcal{L}_{\mathcal{H}_t^a} + \lambda_s \mathcal{L}_S$$

MEMO

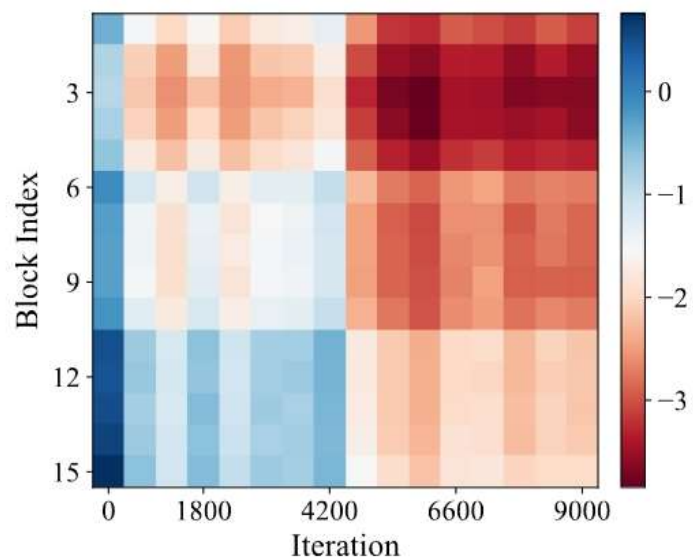




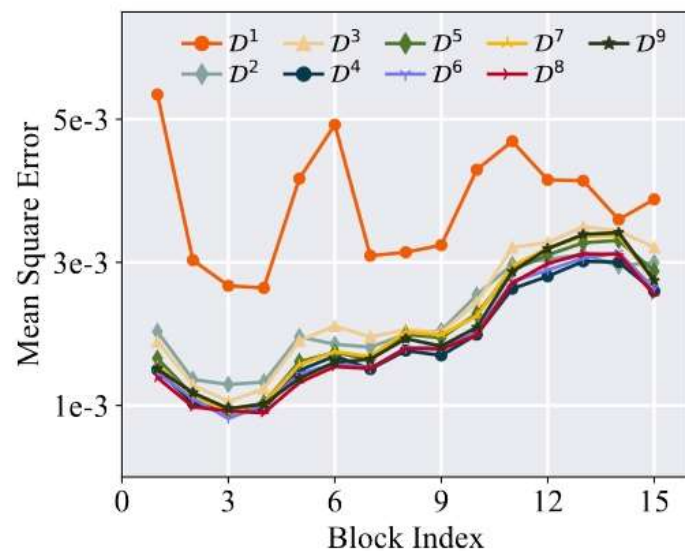
(a) CIFAR100, Base0 Inc10



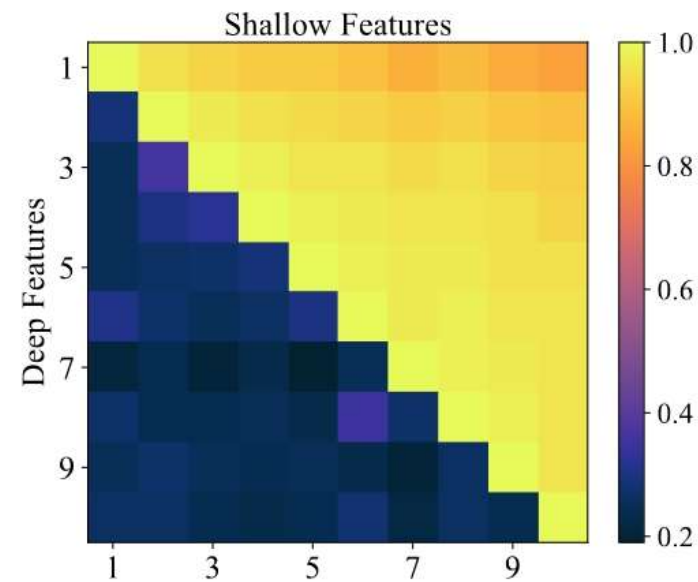
(b) ImageNet100, Base50 Inc5



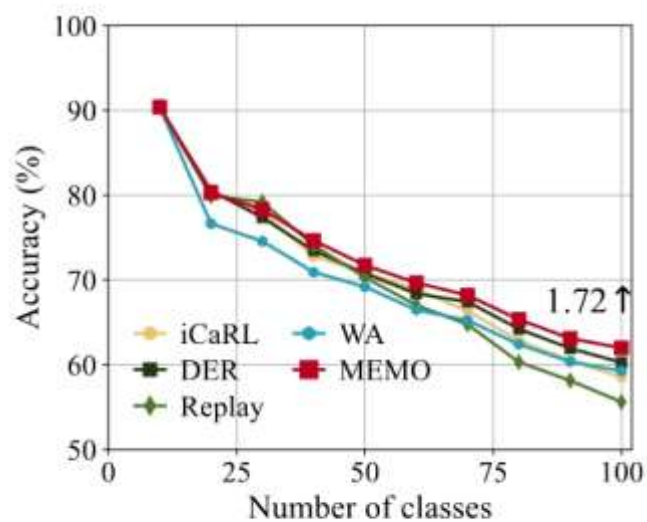
(a) Gradient norm (log scale)



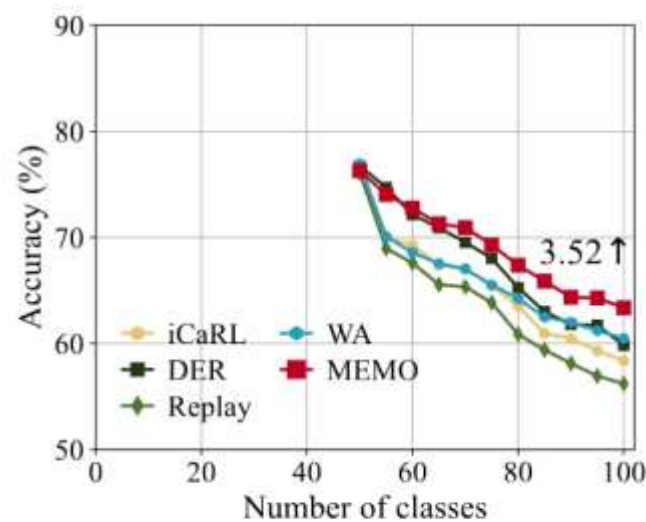
(b) MSE of different blocks



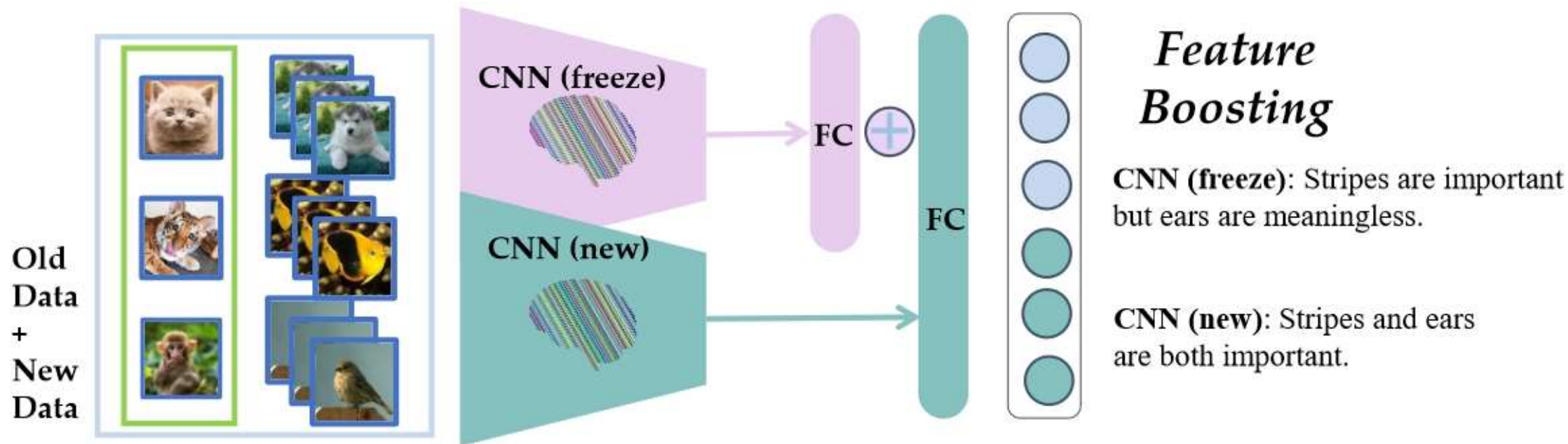
(c) CKA between backbones



(b) CIFAR100 Base0 Inc10



(c) CIFAR100 Base50 Inc5



Foster: Feature Boosting and Compression for Class-Incremental Learning

Parameter Regularization

- keep the important parameters static to maintain former knowledge

$$\mathcal{L} = \ell(f(\mathbf{x}), y) + \frac{1}{2}\lambda \sum_k \Omega_k (\theta_k^{b-1} - \theta_k)^2 .$$

- Fisher matrix
- Prompt

Algorithm-Centric Class-Incremental Learning

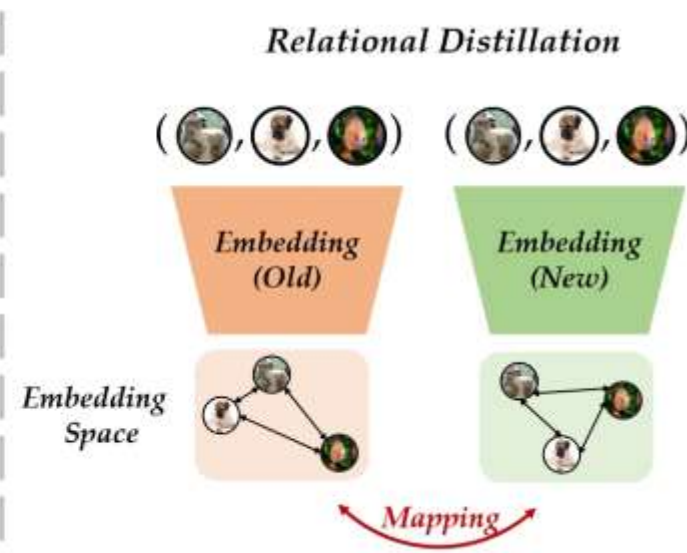
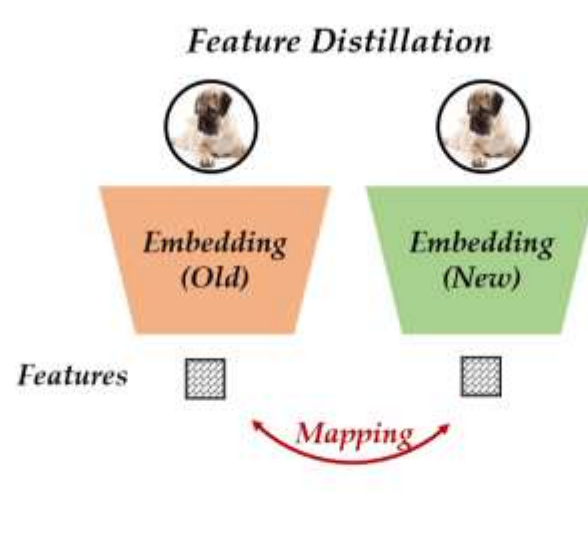
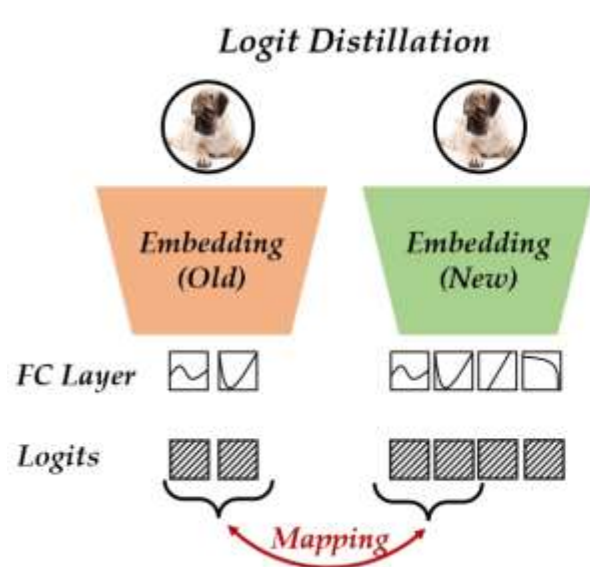
- Knowledge Distillation
- Model Rectify

Knowledge Distillation

$$\mathcal{L} = \underbrace{\ell(f(\mathbf{x}), y)}_{\text{Learning New Classes}} + \underbrace{\sum_{k=1}^{|\mathcal{Y}_{b-1}|} -\mathcal{S}_k(f^{b-1}(\mathbf{x})) \log \mathcal{S}_k(f(\mathbf{x}))}_{\text{Remembering Old Classes}}, \quad (14)$$

$$\mathcal{L} = \ell(f(\mathbf{x}), y) + (1 - \langle \frac{\phi^{b-1}(\mathbf{x})}{\|\phi^{b-1}(\mathbf{x})\|}, \frac{\phi(\mathbf{x})}{\|\phi(\mathbf{x})\|} \rangle).$$

$$\sum_{\{\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k\} \in \mathcal{D}^b} \|\cos \angle \mathbf{t}_i \mathbf{t}_j \mathbf{t}_k - \cos \angle \mathbf{s}_i \mathbf{s}_j \mathbf{s}_k\|,$$



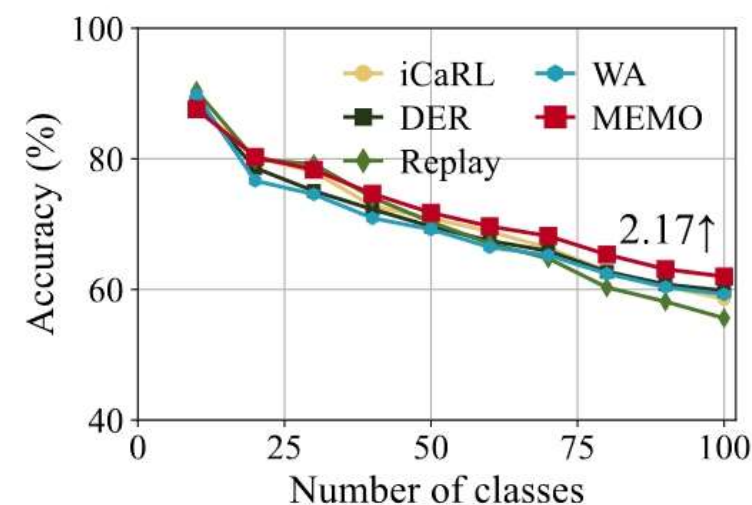
Model Rectify

- Try to find the abnormal behaviors in CIL models and rectify them like the oracle model
 - the weight norm of new classes is significantly larger than old ones
 - the model tends to predict instances as the new classes with larger weights
 - weight drifting
 - logits of new classes are much larger than old ones

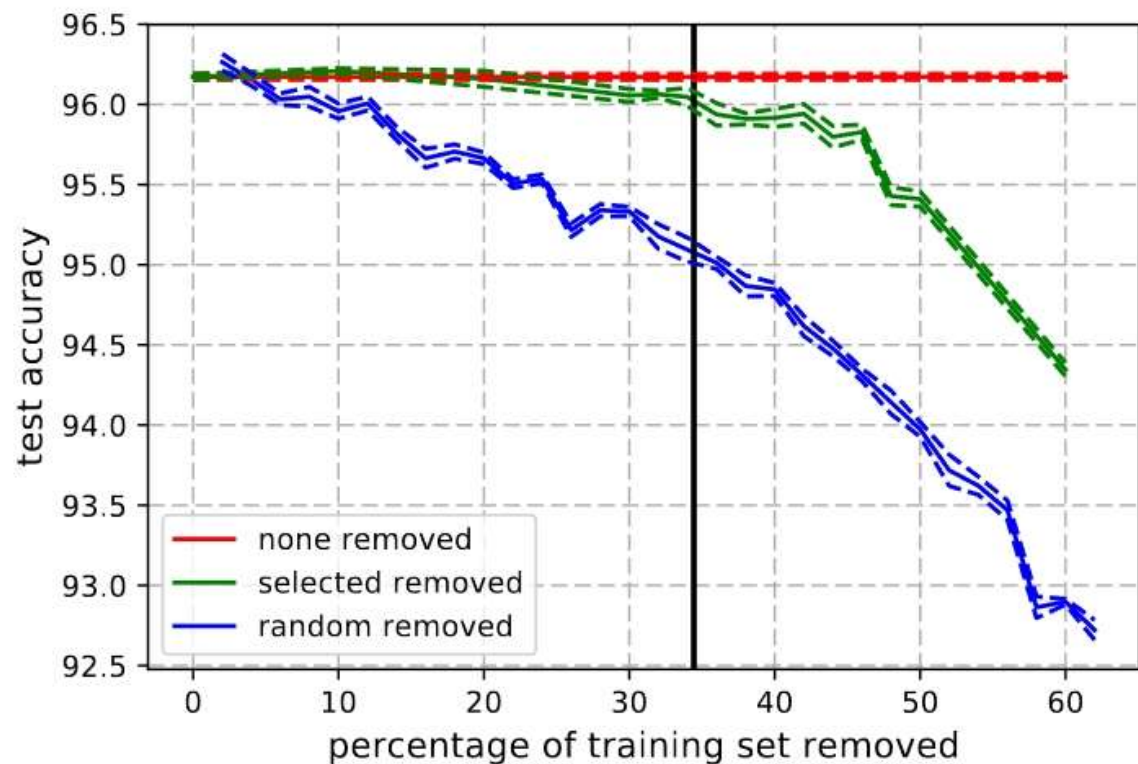
- Data imbalance
- How to select exemplar

Methods	5 steps		10 steps		20 steps		50 steps	
	#Paras	Avg	#Paras	Avg	#Paras	Avg	#Paras	Avg
Bound	11.2	80.40	11.2	80.41	11.2	81.49	11.2	81.74
iCaRL[27]	11.2	71.14 \pm 0.34	11.2	65.27 \pm 1.02	11.2	61.20 \pm 0.83	11.2	56.08 \pm 0.83
UCIR[12]	11.2	62.77 \pm 0.82	11.2	58.66 \pm 0.71	11.2	58.17 \pm 0.30	11.2	56.86 \pm 3.74
BiC[12]	11.2	73.10 \pm 0.55	11.2	68.80 \pm 1.20	11.2	66.48 \pm 0.32	11.2	62.09 \pm 0.85
WA[39]	11.2	72.81 \pm 0.28	11.2	69.46 \pm 0.29	11.2	67.33 \pm 0.15	11.2	64.32 \pm 0.28
PODNet[6]	11.2	66.70 \pm 0.64	11.2	58.03 \pm 1.27	11.2	53.97 \pm 0.85	11.2	51.19 \pm 1.02
RPSNet[26]	60.6	70.5	56.5	68.6	-	-	-	-
Ours(w/o P)	33.6	76.80 \pm 0.79 (+3.7)	61.6	75.36 \pm 0.36 (+5.9)	117.6	74.09 \pm 0.33 (+6.76)	285.6	72.41 \pm 0.36 (+8.09)
Ours	2.89	75.55 \pm 0.65 (+2.45)	4.96	74.64 \pm 0.28 (+5.18)	7.21	73.98 \pm 0.36 (+6.65)	10.15	72.05 \pm 0.55 (+7.73)

23.5MB	$ \mathcal{E} $	$S(\mathcal{E})$	Model Type	# Parameters	Model Size
Replay	7400	21.67MB	ResNet32	0.46M	1.76MB
iCaRL	7400	21.67MB	ResNet32	0.46M	1.76MB
WA	7400	21.67MB	ResNet32	0.46M	1.76MB
DER	2000	5.85MB	ResNet32	4.60M	17.6MB
MEMO	3300	9.66MB	ResNet32	3.62M	13.83MB



AN EMPIRICAL STUDY OF EXAMPLE FORGETTING DURING DEEP NEURAL NETWORK LEARNING



Algorithm 1 Computing forgetting statistics.

```
initialize  $\text{prev\_acc}_i = 0, i \in \mathcal{D}$ 
initialize forgetting  $T[i] = 0, i \in \mathcal{D}$ 
while not training done do
     $B \sim \mathcal{D}$  # sample a minibatch
    for example  $i \in B$  do
        compute  $\text{acc}_i$ 
        if  $\text{prev\_acc}_i > \text{acc}_i$  then
             $T[i] = T[i] + 1$ 
         $\text{prev\_acc}_i = \text{acc}_i$ 
    gradient update classifier on  $B$ 
return  $T$ 
```

