



**ICML**  
International Conference  
On Machine Learning

# Optimizing Mode Connectivity for Class Incremental Learning

Haitao Wen, Haoyang Cheng, Heqian Qiu, Lanxiao Wang, Lili Pan, Hongliang Li

University of Electronic Science and Technology of China

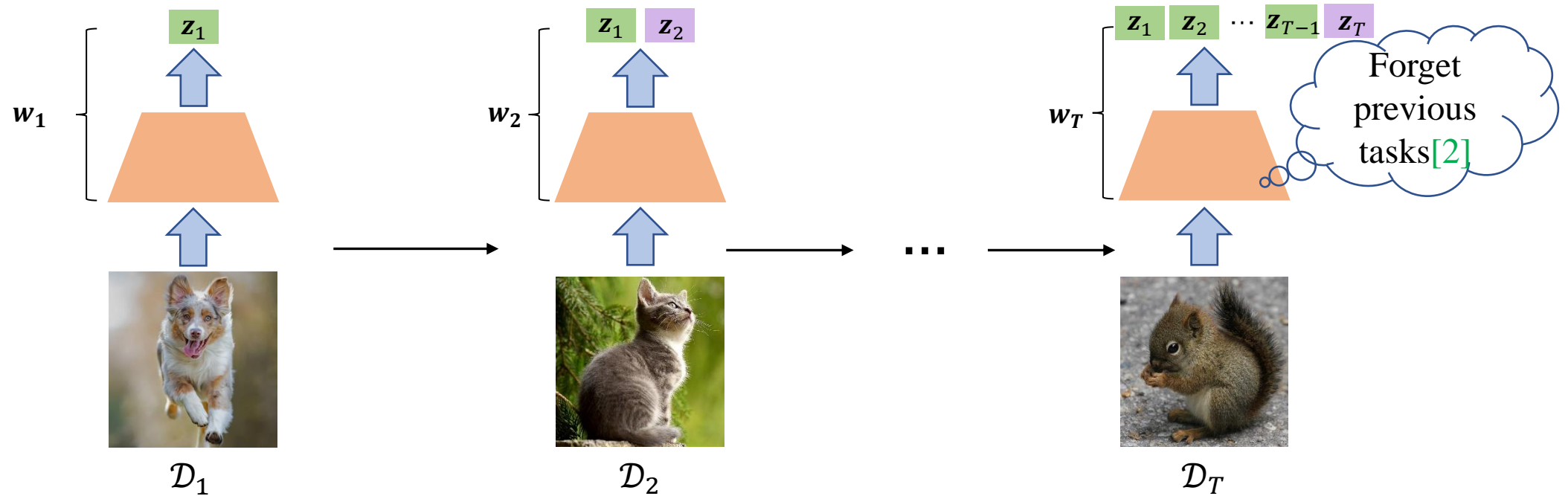
Email: [haitaowen@std.uestc.edu.cn](mailto:haitaowen@std.uestc.edu.cn)

<https://github.com/HaitaoWen/EOPC>

# Introduction

 Backbone     Classifier head     New classifier head

$\mathcal{D}_i$ : dataset of the  $i$ -th task     $w_i$ : model parameters the  $i$ -th task     $z_i$ : new classifier parameters for  $i$ -th task



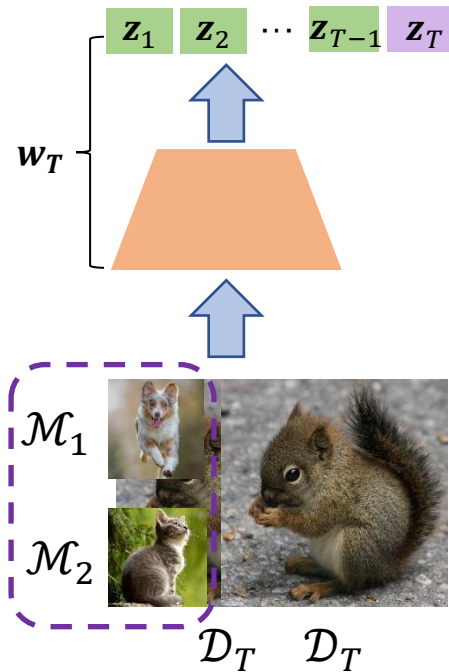
The process of class incremental learning (CIL)[1].

[1] Van de Ven, Gido M., and Andreas S. Tolias. "Three scenarios for continual learning." arXiv preprint arXiv:1904.07734 (2019).

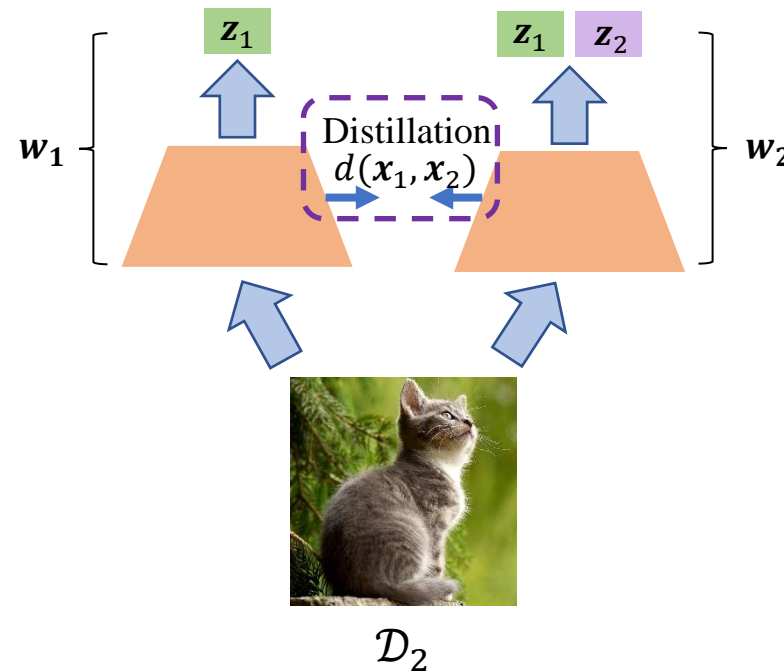
[2] McCloskey, Michael, and Neal J. Cohen. "Catastrophic interference in connectionist networks: The sequential learning problem." Psychology of learning and motivation. Vol. 24. Academic Press, 1989. 109-165.

# Existing Work

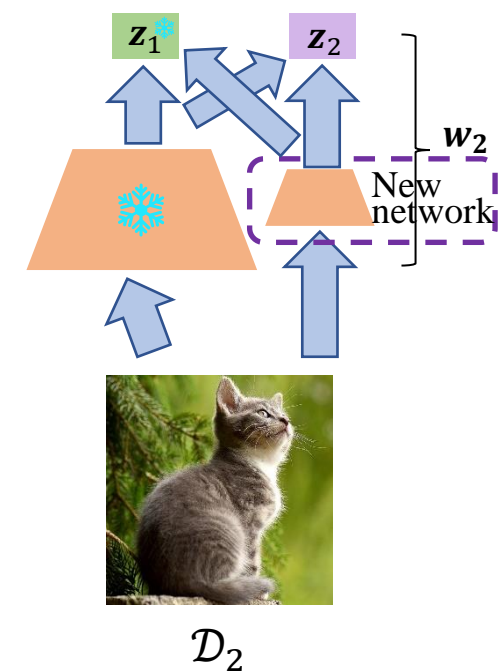
$\mathcal{M}_i$ : memory of the  $i$ -th task     $x_i$ : feature of the  $i$ -th task     $d(x, y)$ : similarity metric    ❄️: frozen



Memory Replay [3]



Regularization [4]



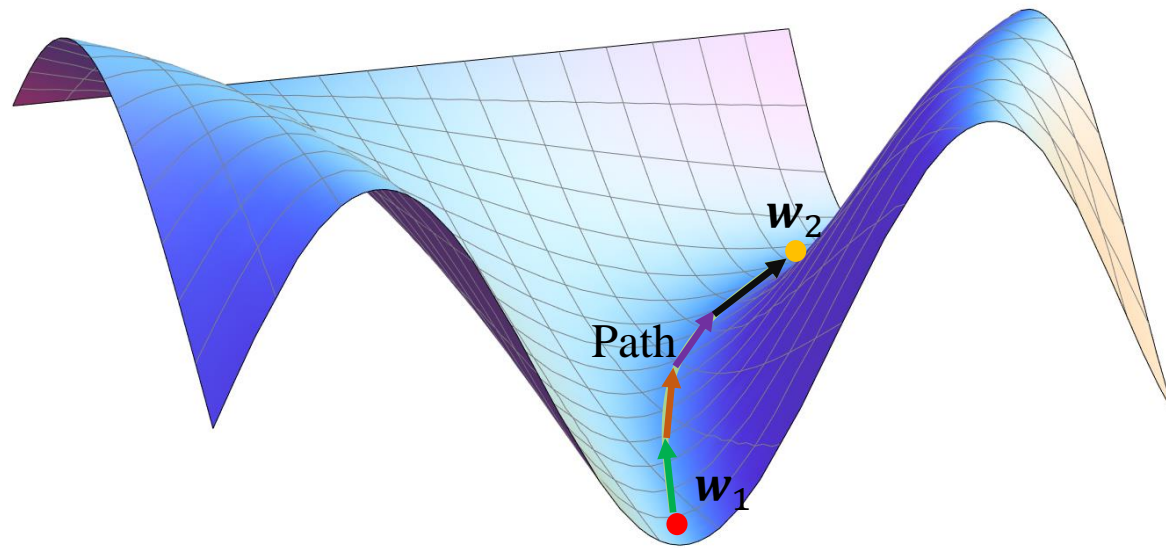
Dynamic Architecture [5]

[3] Liu, Yaoyao, et al. "Mnemonics training: Multi-class incremental learning without forgetting." Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition. 2020.

[4] Douillard, Arthur, et al. "Podnet: Pooled outputs distillation for small-tasks incremental learning." Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16. Springer International Publishing, 2020.

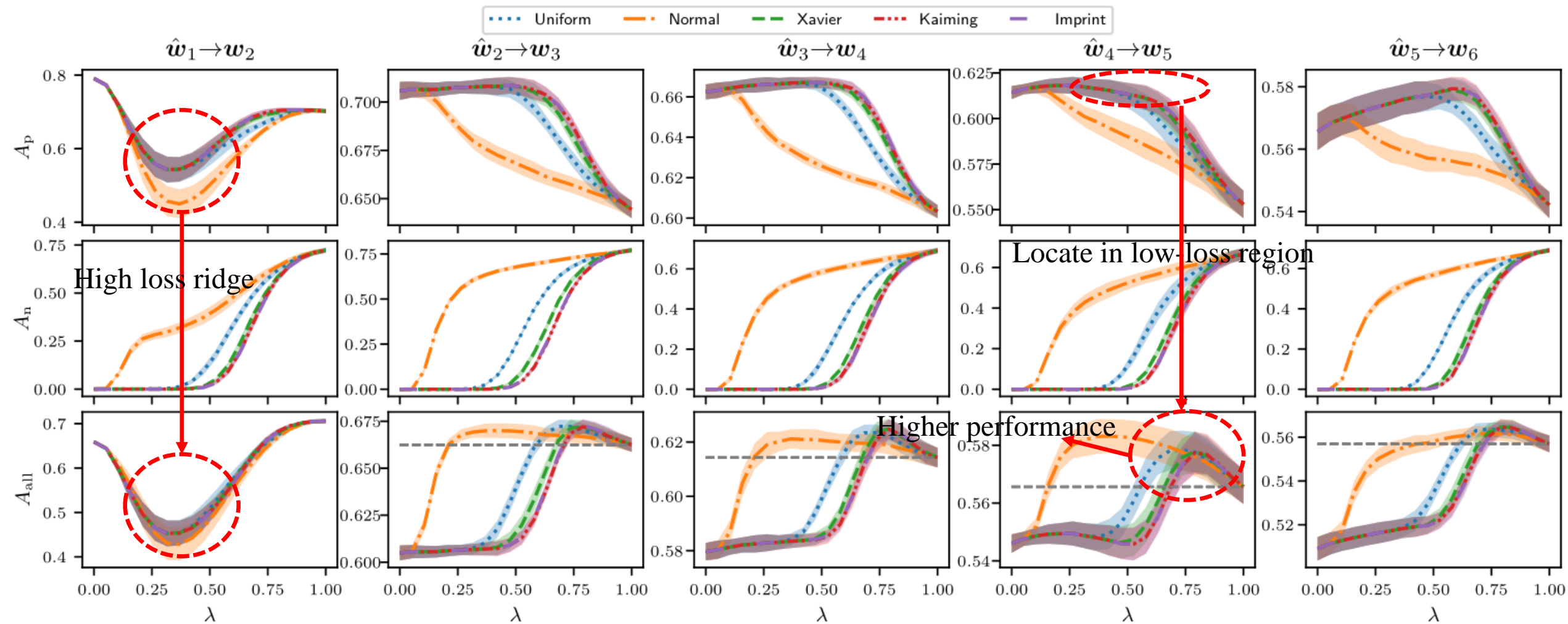
[5] Liu, Yaoyao, Bernt Schiele, and Qianru Sun. "Adaptive aggregation networks for class-incremental learning." Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition. 2021.

# Mode Connectivity



Loss landscape [6]

# Linear Connectivity in CIL



Testing accuracy curves along the linear connection between two adjacent continual minima of PODNet [4] for 5 steps of increments (i.e., 6 tasks in total) on CIFAR-100.  $A_p$ ,  $A_n$ , and  $A_{all}$  denote accuracy on previous tasks, on the new task, and on all learned tasks respectively.  $\lambda$  is the interpolation factor.

# OPC: Optimizing Connectivity between Minima

Let  $\mathbf{p}_\theta(\lambda): [0,1] \rightarrow \mathbb{R}^n$  be the parameterized arbitrary path between minima  $\hat{\mathbf{w}}_{t-1}$  and  $\mathbf{w}_t$ , such that

$$\mathbf{p}_\theta(0) = \hat{\mathbf{w}}_{t-1} \quad \text{and} \quad \mathbf{p}_\theta(1) = \mathbf{w}_t$$

We commonly use the expected loss  $\hat{\ell}(\boldsymbol{\theta})$  along the path to characterize its quality, i.e.,

$$\hat{\ell}(\boldsymbol{\theta}) = \int_0^1 \mathcal{L}(\mathbf{p}_\theta(\lambda)) d\lambda = \mathbb{E}_{\lambda \sim U(0,1)} [\mathcal{L}(\mathbf{p}_\theta(\lambda))],$$

where  $\mathcal{L}$  is the task loss, such as cross-entropy loss, NCA loss [4], or embedding loss [7],  $U(0,1)$  is the uniform distribution on the interval  $[0,1]$ .

we can randomly sample points  $\lambda$  between  $[0,1]$  and minimize loss  $\mathcal{L}(\mathbf{p}_\theta(\lambda))$  with respect to  $\boldsymbol{\theta}$  to optimize the path, i.e.,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{p}_\theta(\lambda)), \lambda \sim U[0,1]$$

However, there are significant differences between continual minima because of catastrophic forgetting, we can not directly connect them each other. We redefine a low-loss path taking named switching point (SP) as a bridge, where the part between the previous minimum and SP is for previous tasks, and the part between SP and the new minimum is for the new task.

# OPC: Optimizing Connectivity between Minima

The expected loss along the path can be reformulated for continual learning as follows,

$$\ell(\boldsymbol{\theta}) = \int_0^{\lambda^*} \mathcal{L}_{1:t-1}(\mathbf{p}_{\boldsymbol{\theta}}(\lambda)) d\lambda + \int_{\lambda^*}^1 \mathcal{L}_t(\mathbf{p}_{\boldsymbol{\theta}}(\lambda)) d\lambda,$$

where  $\lambda^*$  corresponds to SP in the interval  $[0,1]$ ,  $\mathcal{L}_{1:t-1}$  is loss on previous tasks and  $\mathcal{L}_t$  is loss on the new task.

# OPC: Optimizing Connectivity between Minima

## Connectivity Modeling:

Path modeled by Fourier series:

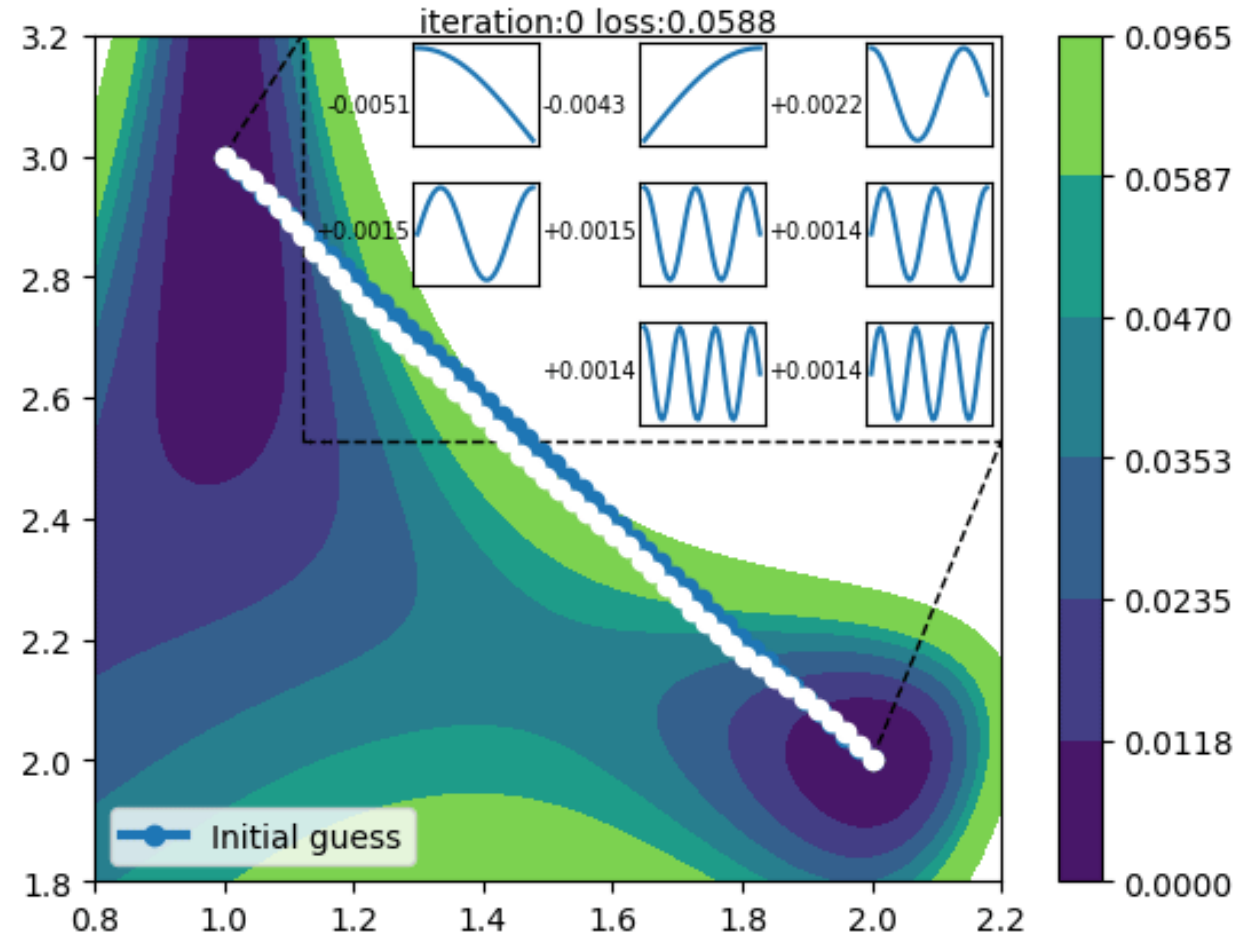
$$\mathbf{p}_\theta(\lambda) = (\mathbf{A}\mathbf{C} + (1 - \lambda)\mathbf{1}_L) \cdot \hat{\mathbf{w}}_{t-1} + (\mathbf{B}\mathbf{S} + \lambda\mathbf{1}_L) \cdot \mathbf{w}_t$$

$$\mathbf{A} = \begin{bmatrix} \alpha_{1,1} & \dots & \alpha_{1,N} \\ & \ddots & \\ \alpha_{L,1} & \dots & \alpha_{L,N} \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} \cos(\frac{\pi}{2}\lambda) \\ \vdots \\ \cos(\frac{(4N-3)\pi}{2}\lambda) \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} \beta_{1,1} & \dots & \beta_{1,N} \\ & \ddots & \\ \beta_{L,1} & \dots & \beta_{L,N} \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} \sin(\frac{\pi}{2}\lambda) \\ \vdots \\ \sin(\frac{(4N-3)\pi}{2}\lambda) \end{bmatrix}$$

where  $\theta = [\mathbf{A}^T, \mathbf{B}^T]^T$ , it must meet the following Condition to make path still pass through endpoints:

$$\theta \mathbf{1}_N = \mathbf{0}_{2L}$$



A toy example of optimizing connectivity between minima in the 2-dimensional plane.



# OPC: Optimizing Connectivity between Minima

## Connectivity Regularization:

- The tangent of the path is

$$\mathbf{p}'_{\theta}(\lambda) = (\mathbf{A}'\mathbf{S} - \mathbf{1}_L) \cdot \hat{\mathbf{w}}_{t-1} + (\mathbf{B}'\mathbf{C} + \mathbf{1}_L) \cdot \mathbf{w}_t$$

Where  $\mathbf{A}', \mathbf{B}' \in \mathbb{R}^{L \times N}$ , specifically,

$$\mathbf{A}' = \begin{bmatrix} -\frac{\pi}{2}\alpha_{1,1} & \dots & -\frac{(4N-3)\pi}{2}\alpha_{1,N} \\ & \vdots & \\ -\frac{\pi}{2}\alpha_{L,1} & \dots & -\frac{(4N-3)\pi}{2}\alpha_{L,N} \end{bmatrix}$$
$$\mathbf{B}' = \begin{bmatrix} \frac{\pi}{2}\beta_{1,1} & \dots & \frac{(4N-3)\pi}{2}\beta_{1,N} \\ & \vdots & \\ \frac{\pi}{2}\beta_{L,1} & \dots & \frac{(4N-3)\pi}{2}\beta_{L,N} \end{bmatrix}$$

- Then, randomly sample noise from the normal distribution and orthogonalize it with the tangent direction, i.e.,

$$\boldsymbol{\epsilon} = \hat{\boldsymbol{\epsilon}} - \frac{\mathbf{p}'_{\theta}(\lambda)^{\top} \hat{\boldsymbol{\epsilon}} \mathbf{p}'_{\theta}(\lambda)}{\|\mathbf{p}'_{\theta}(\lambda)\|^2}, \quad \hat{\boldsymbol{\epsilon}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Next, add normalized noise scaled by radius  $r$  to the path and obtain the point on the surface of the cylinder,

$$\tilde{\mathbf{p}}_{\theta}(\lambda) = \mathbf{p}_{\theta}(\lambda) + r \frac{\boldsymbol{\epsilon}}{\|\boldsymbol{\epsilon}\|}$$

- Replace  $\mathbf{p}_{\theta}(\lambda)$  with  $\tilde{\mathbf{p}}_{\theta}(\lambda)$  in expected loss along the path, we can get the flattening connectivity regularization, i.e.,

$$\ell(\boldsymbol{\theta}) = \int_0^{\lambda^*} \mathcal{L}_{1:t-1}(\tilde{\mathbf{p}}_{\theta}(\lambda)) d\lambda + \int_{\lambda^*}^1 \mathcal{L}_t(\tilde{\mathbf{p}}_{\theta}(\lambda)) d\lambda$$

# OPC: Optimizing Connectivity between Minima

## Connectivity Optimization:

- The parameters of the path  $\theta = [A^T, B^T]^T$  must meet the following condition to make path keep pass through endpoints:

$$\theta \mathbf{1}_N = \mathbf{0}_{2L}$$

- We adopt gradient projection to update parameters along the direction orthogonal to the normal of equation, i.e.,:

$$\begin{aligned} \Delta(\lambda) &= \nabla_{\theta} \mathcal{L}(\tilde{\mathbf{p}}_{\theta}(\lambda)) (\mathbf{I}_{N \times N} - \mathbf{1}_N (\mathbf{1}_N^T \mathbf{1}_N)^{-1} \mathbf{1}_N^T) \\ &= \nabla_{\theta} \mathcal{L}(\tilde{\mathbf{p}}_{\theta}(\lambda)) - \frac{1}{N} \nabla_{\theta} \mathcal{L}(\tilde{\mathbf{p}}_{\theta}(\lambda)) \mathbf{1}_N \mathbf{1}_N^T. \end{aligned}$$

- Then, the iterative rule for parameters of the path is as follows,

$$\theta \leftarrow \theta - \gamma \Delta(\lambda), \quad \lambda \sim U[0, 1]$$

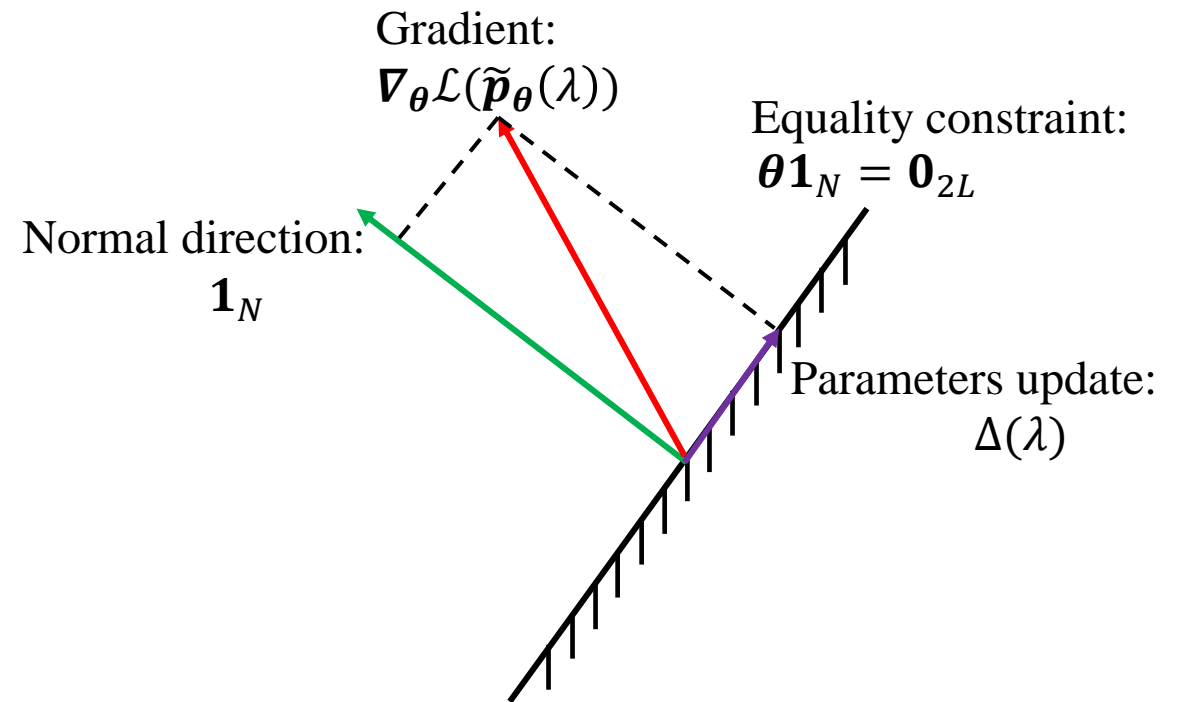


Diagram of gradient projection

# OPC: Optimizing Connectivity between Minima

## EOPC: Ensembling with OPC

- The optimized path provides infinite low-loss solutions on both sides of the switching point (SP), i.e.,  $\mathbf{p}_\theta(\lambda^*)$ . To further improve performance on learned tasks, we propose EOPC to ensemble points within a local bent cylinder around SP.
- We take  $\bar{\mathbf{w}}$  as the minimum of the current task and the initial parameters of model in the next task.

- The cylinder is constructed according to the tangent of the path, let  $S$  be the set of points within this cylinder and can be formulated as follows,

$$S = \left\{ \mathbf{w} \mid (\mathbf{w} - \mathbf{p}_\theta(\lambda))^T \mathbf{p}'_\theta(\lambda) = 0, \|\mathbf{w} - \mathbf{p}_\theta(\lambda)\|_2 \leq r; \right. \\ \left. \lambda \in [\lambda^* - \tau/2, \lambda^* + \tau/2] \right\},$$

- We adopt ensembling in parameter space by averaging points within  $S$ , i.e.,

$$\bar{\mathbf{w}} = \frac{1}{M} \sum_{i=1}^M \mathbf{w}_i, \quad \mathbf{w}_i \sim S$$

where  $M$  is the number of total sampling points.

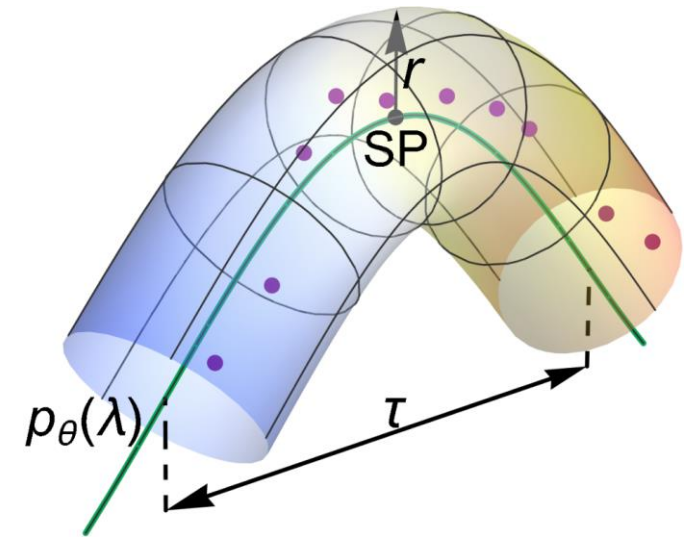
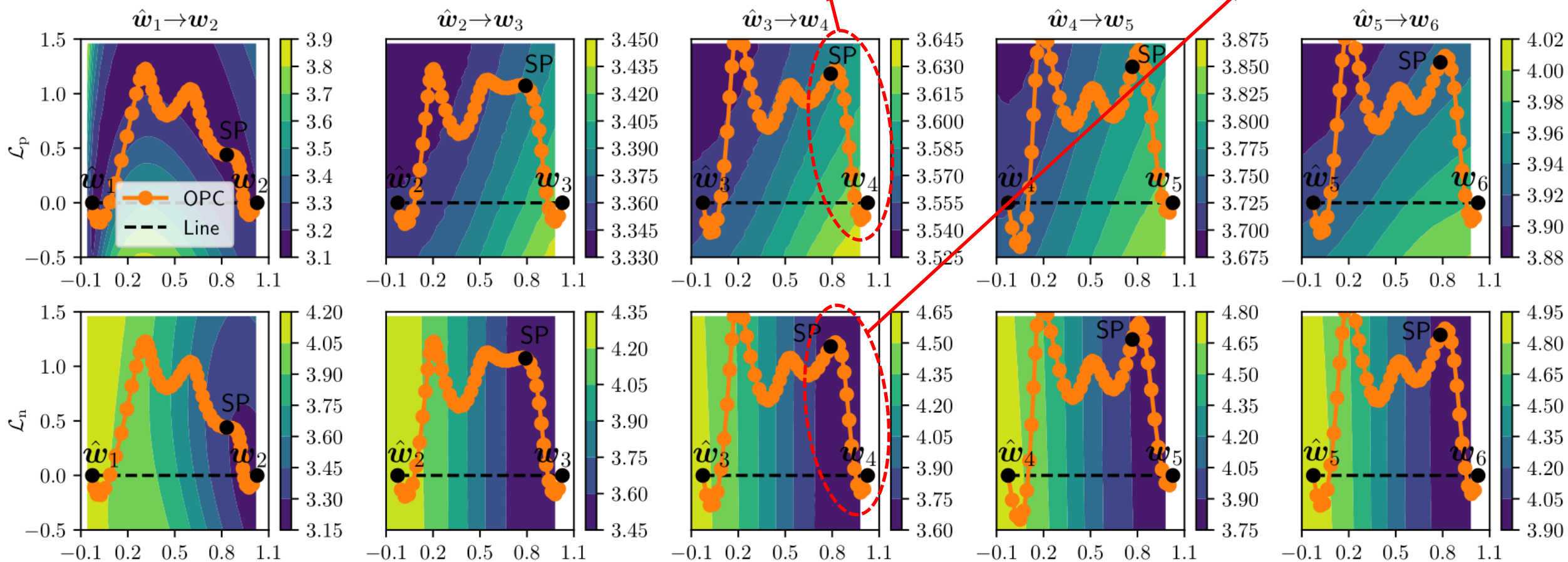


Diagram of local bent cylinder

# Experiments

Compared with  $w_4$ , SP locates in a lower-loss region in landscape of previous tasks.

Compared with  $w_4$ , SP locates in the same low-loss region in landscape of the new task.



Visualization of paths found by OPC in loss landscape of previous tasks ( $\mathcal{L}_p$ ) and the new task ( $\mathcal{L}_n$ ).

# Experiments

Method	CIFAR-100			ImageNet-100			ImageNet-1K			
	$\mathcal{A}$ (%) $\uparrow$	5	10	25	5	10	25	5	10	25
iCaRL (Rebuffi et al., 2017)	57.83	52.63	49.02	64.75	58.80	52.46	51.60	47.42	41.03	
BiC <sup>†</sup> (Wu et al., 2019)	59.36	54.20	50.00	70.07	64.96	57.73	62.65	58.72	53.47	
LUCIR (Hou et al., 2019)	63.62	60.95	57.79	71.93	69.43	63.51	66.13	61.63	54.05	
Mnemonics <sup>†</sup> (Liu et al., 2020)	63.34	62.28	60.96	72.58	71.37	69.74	64.63	63.01	61.00	
GeoDL <sup>†</sup> (Simon et al., 2021)	65.14	65.03	<b>63.12</b>	73.87	73.55	71.72	65.23	64.46	62.20	
AFC (Kang et al., 2022)	65.87	64.45	62.05	77.27	<b>75.47</b>	<b>72.41</b>	69.07	66.85	<b>63.40</b>	
PODNet (Douillard et al., 2020)	65.47	63.13	59.85	76.32	73.54	63.05	68.33	65.35	58.62	
w/ EOPC	66.68 <sup>+1.21</sup>	64.94 <sup>+1.81</sup>	62.36 <sup>+2.51</sup>	77.12 <sup>+0.8</sup>	74.53 <sup>+0.99</sup>	68.18 <sup>+5.13</sup>	<b>69.72</b> <sup>+1.39</sup>	<b>67.57</b> <sup>+2.22</sup>	62.35 <sup>+3.73</sup>	
AANet (Liu et al., 2021)	66.53	64.63	61.05	77.98	74.70	68.65	68.87	65.65	60.07	
w/ EOPC	<b>67.55</b> <sup>+1.02</sup>	<b>65.54</b> <sup>+0.91</sup>	61.82 <sup>+0.77</sup>	<b>78.95</b> <sup>+0.97</sup>	74.99 <sup>+0.29</sup>	70.10 <sup>+1.45</sup>	69.47 <sup>+0.6</sup>	67.35 <sup>+1.7</sup>	62.20 <sup>+2.13</sup>	

The adaptation results of EOPC and comparison results with existing incremental learning methods on CIFAR-100, ImageNet-100, and ImageNet-1K.

Thank you!