

# The Dual Nature of Plasticity Loss in Deep Continual Learning: Dissection and Mitigation

Haoyu Albert Wang, Wei P. Dai, Jiawei Zhang,  
Jialun Ma, Mingyi Huang, Yuguo Yu

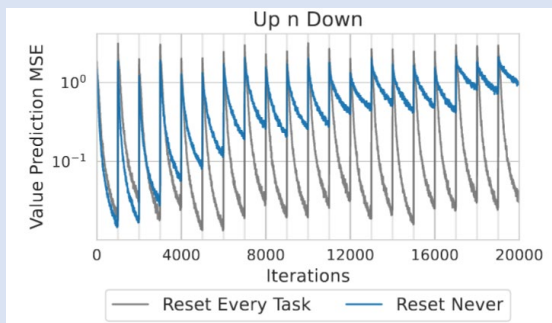
----- *Neural Information Processing Systems (2025)* -----

# Loss of Plasticity in Deep Continual Learning



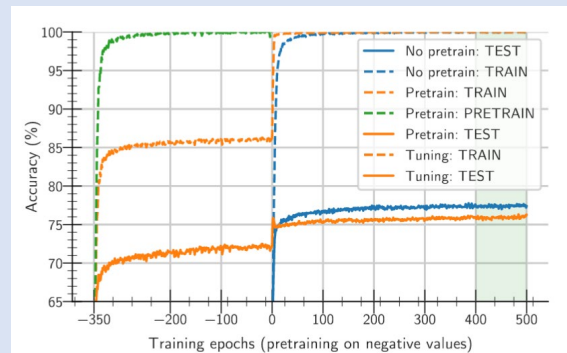
- Loss of plasticity is a widely observed phenomenon in both continual learning and reinforcement learning.
- It refers to the degradation of performance on new tasks, which eventually prevents the system from learning continuously.

## RL



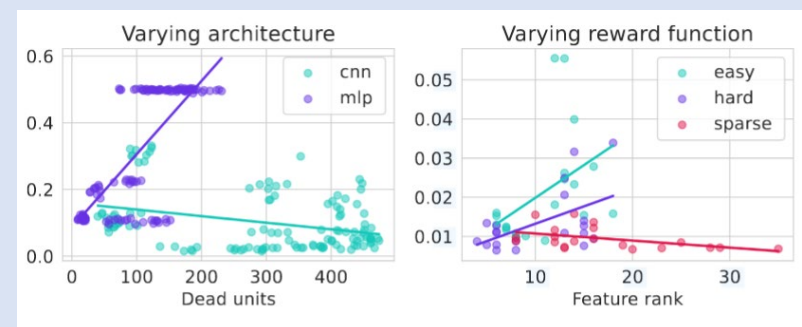
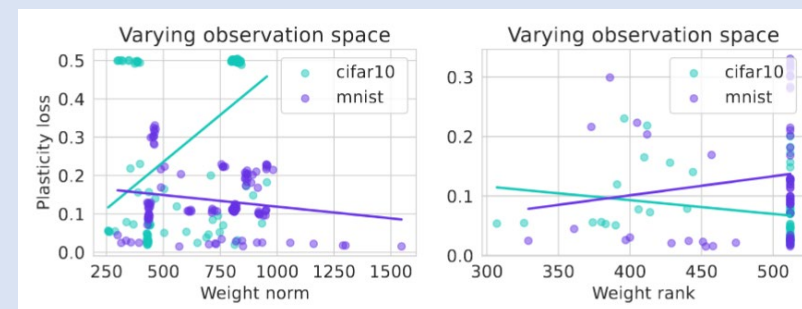
Evgenii Nikishin et al. 2023

## CL



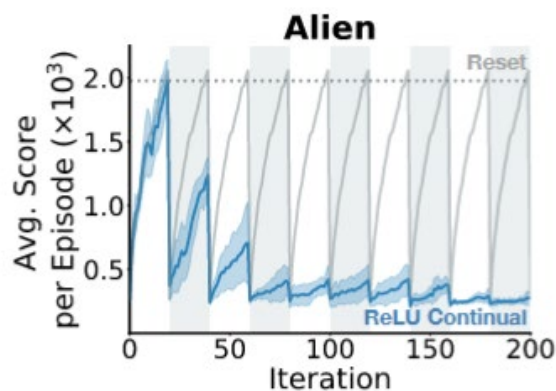
Tudor Berariu et al. 2023

## Potential factors of LoP

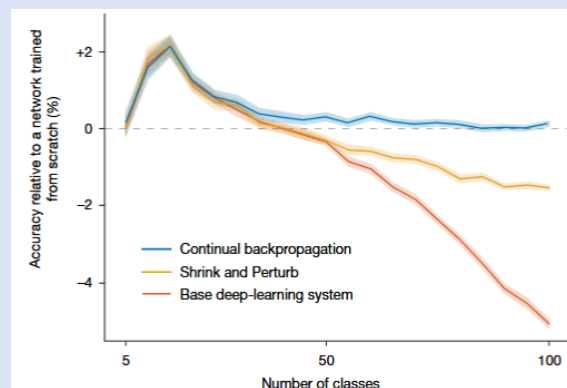


Clare Lyle et al. 2023

**Inconclusive and Indirect**



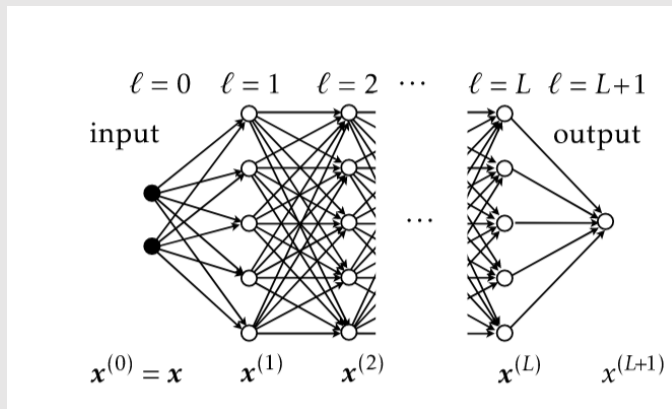
Zaheer Abbas et al. 2023



Shibhansh Dohare et al. 2024

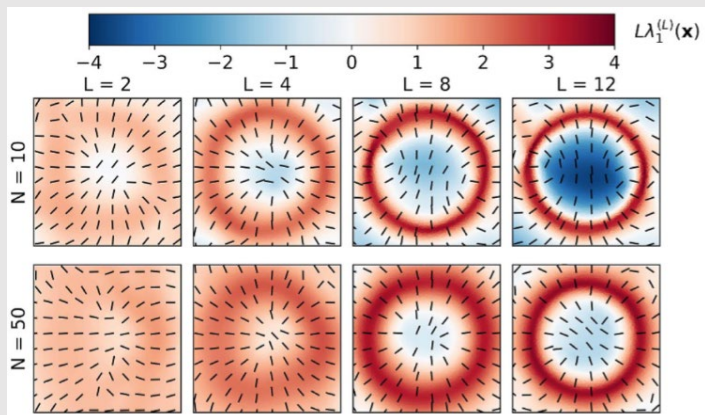
## Combining FTLE and Neural Collapse Reveals the Dynamics of Representation Space Underlying Continual Learning

### FTLE as a Framework for Understanding the Geometric Structure of Neural Network Mappings



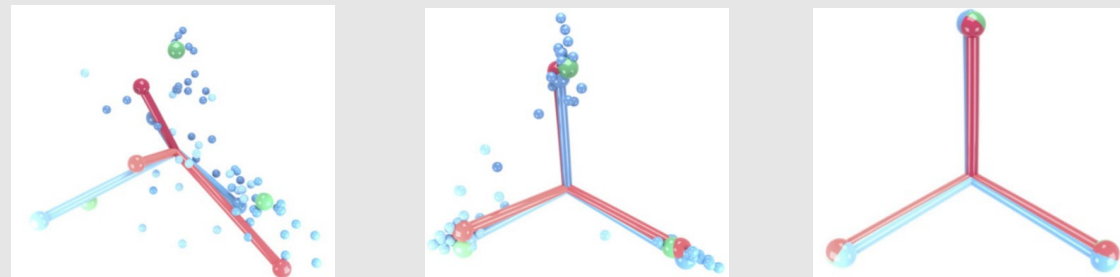
$$\delta \mathbf{x}^{(\ell)} = \mathbb{J}_{\ell} \delta \mathbf{x}$$

$$\lambda_1^{(\ell)}(\mathbf{x}) \equiv \ell^{-1} \log \Lambda_1^{(\ell)}(\mathbf{x})$$



L. Storm et al. 2024

### Prevalence of neural collapse during the terminal phase of deep learning training



(NC1) Variability collapse:  $\Sigma_W \rightarrow \mathbf{0}$ .

(NC2) Convergence to simplex ETF:

$$\begin{aligned} \left| \|\mu_c - \mu_G\|_2 - \|\mu_{c'} - \mu_G\|_2 \right| &\rightarrow 0 \quad \forall c, c' \\ \langle \tilde{\mu}_c, \tilde{\mu}_{c'} \rangle &\rightarrow \frac{C}{C-1} \delta_{c,c'} - \frac{1}{C-1} \quad \forall c, c'. \end{aligned}$$

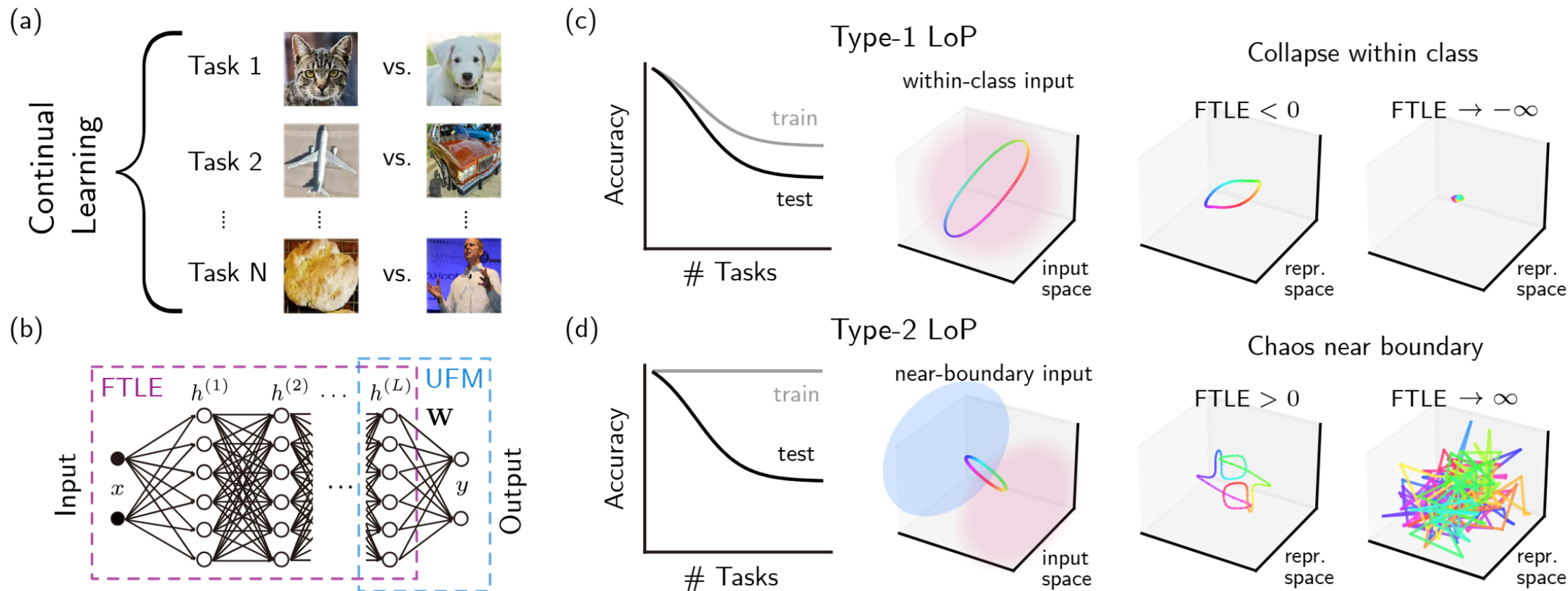
(NC3) Convergence to self-duality:

$$\left\| \frac{W^\top}{\|W\|_F} - \frac{\dot{M}}{\|\dot{M}\|_F} \right\|_F \rightarrow 0.$$

(NC4) Simplification to NCC:

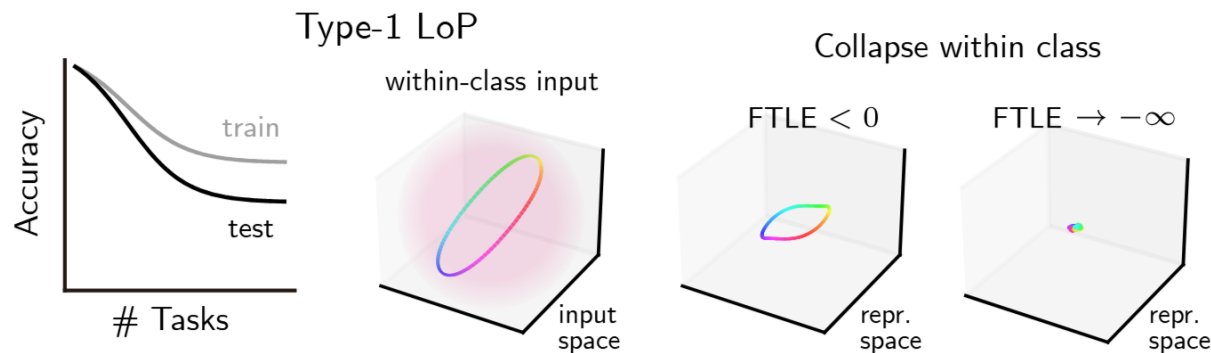
$$\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'} \rightarrow \arg \min_{c'} \|\mathbf{h} - \mu_{c'}\|_2,$$

Vardan Papyan et al. 2020



- **Propose** an FTLE + UFM framework to analyze representation-space dynamics in continual learning.
- **Identify and validate two different types of LoP, both theoretically and experimentally.**
- **Develop** a generalized Mixup method that effectively mitigates both types of LoP.

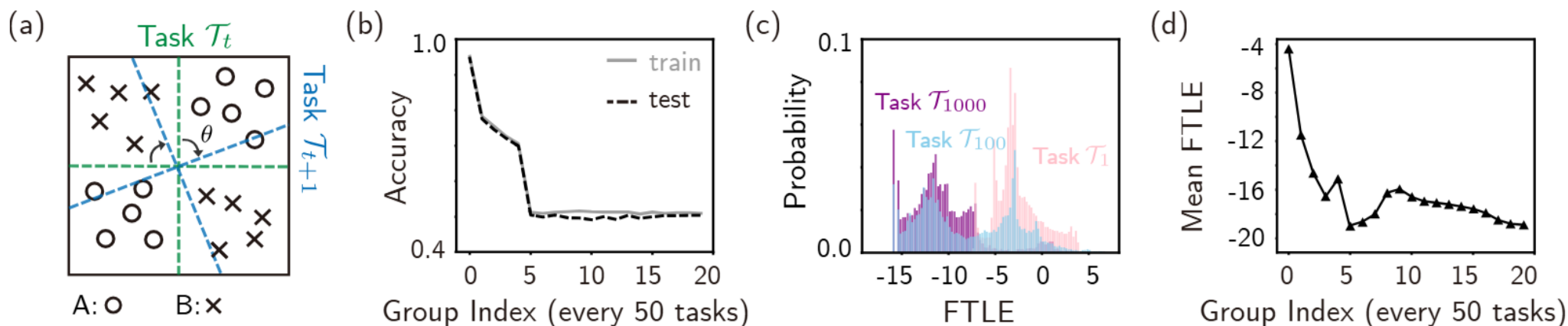
# Type-1 LoP: the Collapse of Representation Space



## Type-1 LoP:

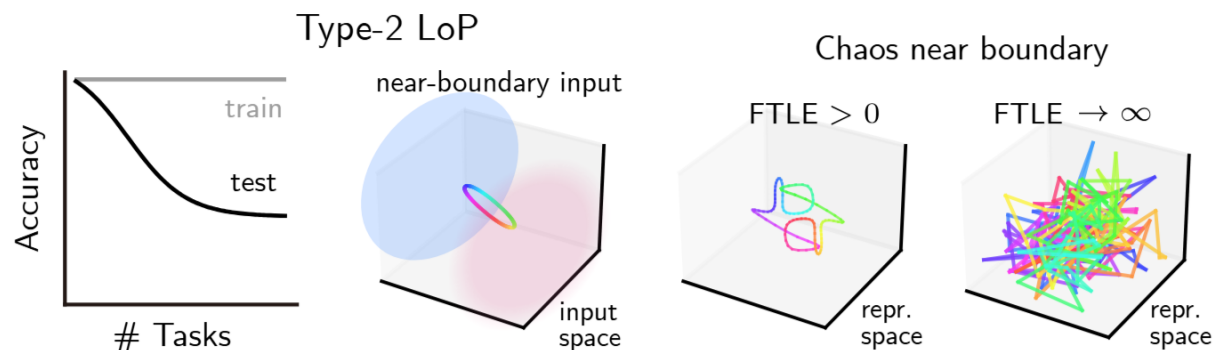
- Learning progressively causes **within-class regions to collapse**, and these collapsed areas accumulate during continual learning.
- Type-1 LoP occurs when representations of a new task approach these previously collapsed regions, characterized by **highly negative FTLEs**.
- Both training and test accuracies drop sharply, indicating a **loss of learning capacity**.

$$\log \beta = \lambda_{\text{trained}}^{k,k} - \lambda_0^{k,k} < 0 \quad \lambda_{\text{trained},T}^{k,k} \rightarrow -\infty \quad \text{as} \quad \prod_{t=1}^T \beta_t \rightarrow 0$$



According to our theory, **Type-1 LoP** is more likely to occur in **low-dimensional representation spaces**. Thus we verified this using a **toy neural network** with a low-dimensional representation layer.

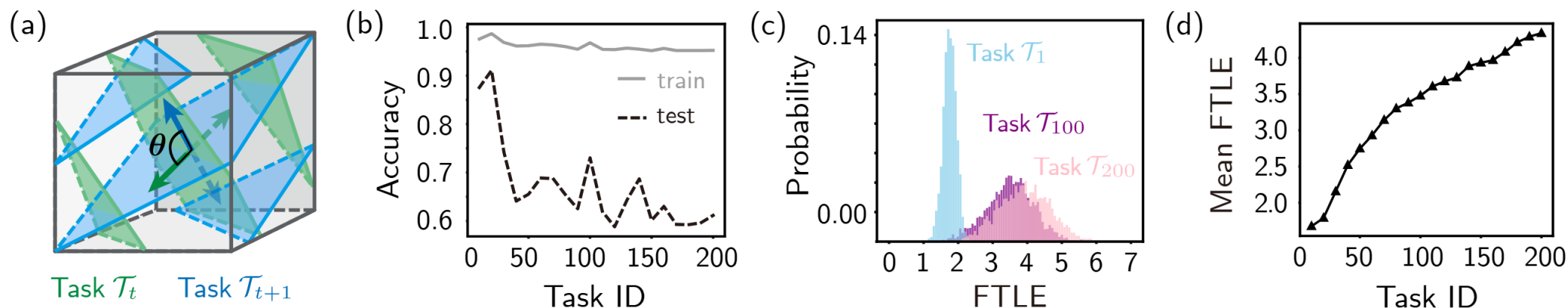
# Type-2 LoP : Over-stretched Boundaries and Chaotic Behaviors



## Type-2 LoP:

- Learning progressively causes **inter-class regions to expand**, and this expansion accumulates during continual learning.
- Type-2 LoP occurs when representations of a new task approach these **overly stretched regions** in representation space, characterized by **highly positive FTLEs**.
- Training accuracy remains high due to the **chaotic and over-expressive representation space**, while test accuracy degrades, indicating a **loss of generalization capacity**.

$$\log \alpha = \lambda_{\text{trained}}^{k,k'} - \lambda_0^{k,k'} > 0 \quad \lambda_{\text{trained},T}^{k,k'} \rightarrow \infty \quad \text{as} \quad \prod_{t=1}^T \alpha_t \rightarrow \infty$$

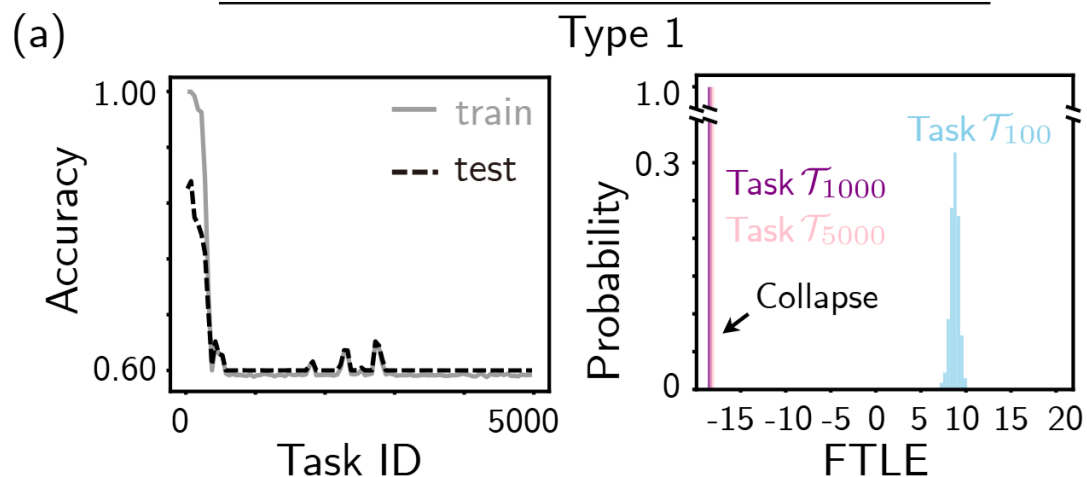


According to our theory, **Type-2 LoP** is more likely to occur in **high-dimensional representation spaces**. Thus we verified this using a **toy neural network** with a high-dimensional representation layer.

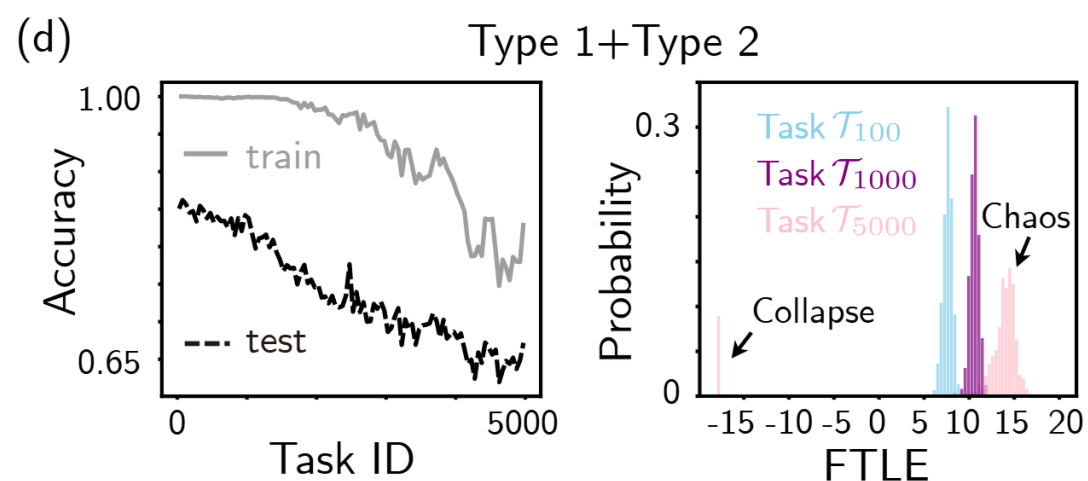
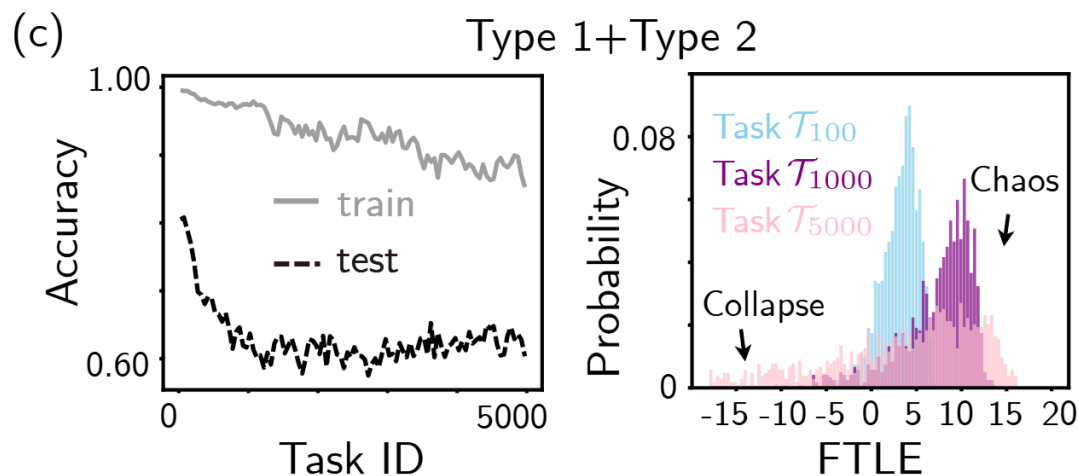
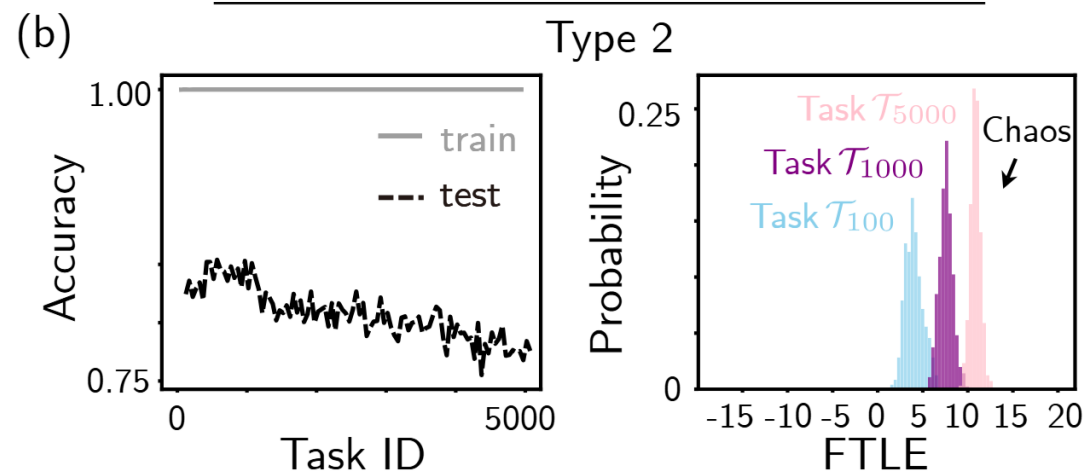
# Verifying LoP in Real Datasets



## Low Dim Space

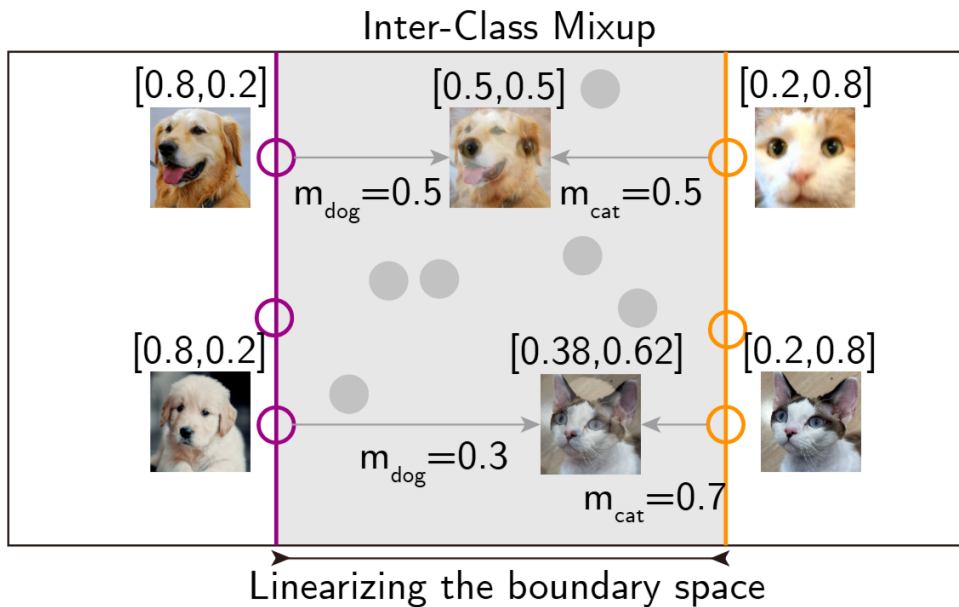


## High Dim Space



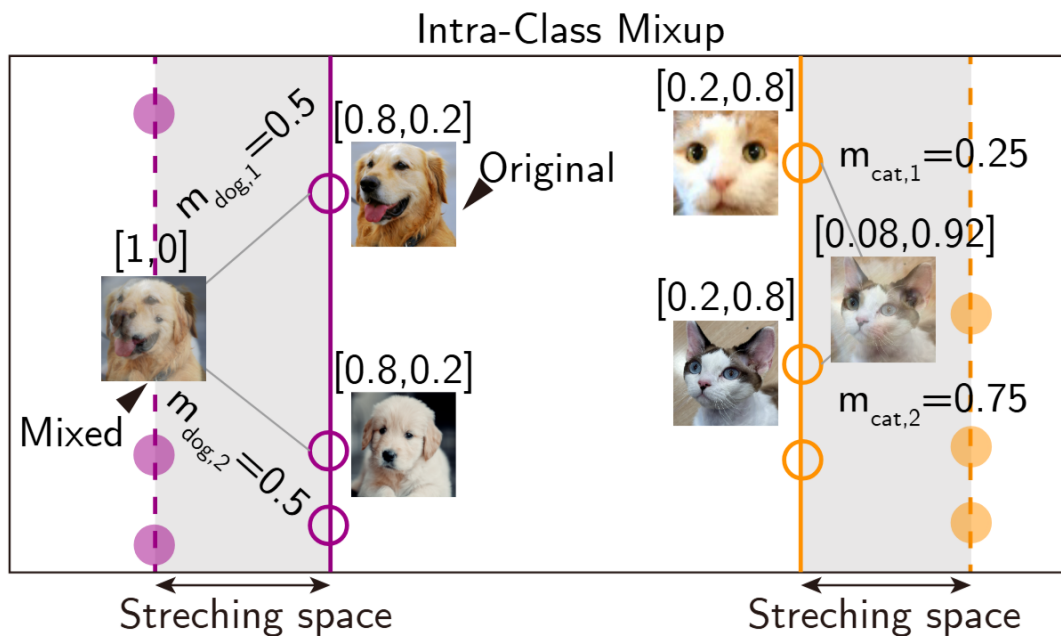
We validate our theoretical prediction on real datasets, confirming the existence of **two different types of LoP** and their **mixed forms**.

# Generalized Mixup: A Simple Method to Mitigate LoP



$$x^m = mx^K + (1 - m)x^{K'}, \quad y^m = my^K + (1 - m)y^{K'}$$

**Inter-class Mixup**, akin to the classical Mixup method, **mitigates Type-2 LoP** by alleviating over-expansion in representation space



$$x^m = mx_i^K + (1 - m)x_j^K$$

$$y_K^m = y^K + \frac{M}{2} - M|0.5 - m|, \quad y_{K'}^m = y^{K'} - \frac{M}{2} + M|0.5 - m|$$

**Intra-class Mixup**, extending classical Mixup by combining samples **within the same class** using different **labels**, mitigates **Type-1 LoP** by preventing **within-class space collapse**

Table 1: SmallConv Acc. on Continual ImageNet

Task ( $\times 1000$ )	0-1	1-2	2-3	3-4	4-5
No Intv.	0.817	0.805	0.562	0.500	0.500
Retrained	0.853	0.845	0.845	0.840	0.840
L2 init	0.804	0.796	0.786	0.785	0.788
Layernorm	0.753	0.760	0.759	0.751	0.751
CBP	0.834	0.847	0.846	0.847	0.857
G-mixup[ours]	<b>0.866</b>	<b>0.881</b>	<b>0.885</b>	<b>0.880</b>	<b>0.879</b>

Table 2: ConvNet Acc. on Continual ImageNet

Task ( $\times 1000$ )	0-1	1-2	2-3	3-4	4-5
No Intv.	0.794	0.778	0.604	0.537	0.500
Retrained	0.857	0.851	0.850	0.849	0.846
L2 init	0.814	0.805	0.800	0.803	0.807
Layernorm	0.782	0.768	0.752	0.749	0.755
CBP	0.848	0.867	0.864	0.863	0.878
G-mixup[ours]	<b>0.875</b>	<b>0.896</b>	<b>0.899</b>	<b>0.894</b>	<b>0.896</b>

Table 3: 0.25x Resnet-18 Acc. on CIFAR100

Task ( $\times 4$ )	0-1	1-2	2-3	3-4	4-5
No Intv.	0.927	0.812	0.754	0.708	0.661
Retrained	0.926	0.800	0.745	0.702	0.666
L2 init	0.925	0.788	0.724	0.682	0.649
Layernorm	0.851	0.755	0.723	0.674	0.641
CBP	0.923	0.812	0.760	0.713	0.678
G-mixup[ours]	<b>0.932</b>	<b>0.825</b>	<b>0.772</b>	<b>0.732</b>	<b>0.697</b>

Table 4: Resnet-18 Acc. on CIFAR100

Task ( $\times 4$ )	0-1	1-2	2-3	3-4	4-5
No Intv.	0.907	0.847	0.812	0.778	0.743
Retrained	0.901	0.842	0.815	0.788	0.767
L2 init	0.922	0.829	0.785	0.752	0.723
Layernorm	0.847	0.782	0.763	0.728	0.697
CBP	0.907	0.847	0.816	0.788	<b>0.768</b>
G-mixup[ours]	<b>0.928</b>	<b>0.864</b>	<b>0.832</b>	<b>0.800</b>	<b>0.768</b>