
Fixed Non-negative Orthogonal Classifier: Inducing Zero-mean Neural Collapse with Feature Dimension Separation - Supplementary Material -

A RELATED WORK

We provided an elaborate treatment of related work.

A.1 ANALYSES AND EXTENSIONS OF NEURAL COLLAPSE

Extensive research has examined the neural collapse from various perspectives, including loss function (Lu & Steinerberger, 2020; Wojtowytsch et al., 2020), the use of mse loss (Han et al., 2021), feature normalization (Yaras et al., 2022), label smoothing (Zhou et al., 2022), fine-grained structure (Yang et al., 2023), generalization (Hui et al., 2022), out-of-distribution with normalization (Haas et al., 2022), imbalanced learning (Dang et al., 2023; Zhong et al., 2023), transfer learning (Galanti et al., 2021; Kornblith et al., 2021), federated learning (Li et al., 2023) or complex of them (Xu et al., 2023). Additionally, some studies have considered extensions of neural collapse such as intermediate layer collapse (Rangamani et al., 2023), generalized neural collapse (Liu et al., 2023), a novel metric (Xu & Liu, 2023). Naturally, geometric analyses have also arisen (Zhu et al., 2021; Thrampoulidis et al., 2022; Tirer et al., 2023).

A.2 ADVANTAGES OF ORTHOGONALITY

Orthogonality has also shown advantages in various environments, such as regularization in multi-task learning (He et al., 2020), vision applications (Cao et al., 2021), features normalization (Lu et al., 2023), optimization (Chiley et al., 2019; Hu et al., 2020; Wu et al., 2020), regularization-better training stability- robustness-generalization (Liu et al., 2021; Achour et al., 2022; Xu et al., 2022; Wang et al., 2020; Li et al., 2019; Ranasinghe et al., 2021; Trockman & Kolter, 2021), explainability (Zhang et al., 2022a; Yang et al., 2020b), continual learning (Pernici et al., 2021; Hersche et al., 2022), reduction on computation load (Yang et al., 2020a), quantization (Ma et al., 2023), graph neural networks (Bodnar et al., 2022; Yang et al., 2022a), transformer (Huang et al., 2022; Kong et al., 2022), representation learning (Medini et al., 2020; Tiao et al., 2023), continual learning (Chaudhry et al., 2020; Ramasesh et al., 2021; Saha et al., 2021; Farajtabar et al., 2020), disentanglement (Sarhan et al., 2020; Liu et al., 2020; Cha & Thiyaalingam, 2023).

B ZERO-MEAN NEURAL COLLAPSE

Class Mean. In the common neural collapse, the trained class mean is as Eq. 1.

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^{N_k} h_{i,k} \quad (1)$$

where K is the number of classes and N_k is the number of input samples included in the k -th class. $h_{i,k} \in \mathbb{R}^D$ means the feature vector of i -th input sample in the k -th class.

The zero class mean μ_k^0 is calculated in the same way with the trained class mean Eq. 2.

$$\mu_k^0 = \mu_k \quad (2)$$

Global Mean. In the common neural collapse, the trained global mean μ_G is as Eq. 3.

$$\mu_G = \frac{1}{K} \sum_{k=1}^K \mu_k \quad (3)$$

where μ_k is the k -th class mean.

This global mean will not be trained and fixed to zero-mean vector like Eq. 4.

$$\mu_G^0 = \mathbf{0}_D \quad (4)$$

where μ_G^0 denotes the zero global mean and $\mathbf{0}_D \in \mathbb{R}^D$ is the zeros vector.

Total Covariance. In the common neural collapse, the trained total covariance $\Sigma_T \in \mathbb{R}^{D \times D}$ is as Eq. 5.

$$\Sigma_T = \frac{1}{K} \sum_{k=1}^K \left(\frac{1}{N_k} \sum_{i=1}^{N_k} (\mathbf{h}_{i,k} - \boldsymbol{\mu}_G) (\mathbf{h}_{i,k} - \boldsymbol{\mu}_G)^\top \right) \quad (5)$$

This total covariance will be changed like Eq. 6.

$$\Sigma_T^0 = \frac{1}{K} \sum_{k=1}^K \left(\frac{1}{N_k} \sum_{i=1}^{N_k} \mathbf{h}_{i,k} \mathbf{h}_{i,k}^\top \right) \quad (\because \boldsymbol{\mu}_G^0 = \mathbf{0}_D) \quad (6)$$

where $\Sigma_T^0 \in \mathbb{R}^{D \times D}$ denotes the total covariance in the zero-mean NC.

Between-class Covariance. In the common neural collapse, the between-class covariance $\Sigma_B \in \mathbb{R}^{D \times D}$ is as Eq. 7.

$$\Sigma_B = \frac{1}{K} \sum_{k=1}^K ((\boldsymbol{\mu}_k - \boldsymbol{\mu}_G) (\boldsymbol{\mu}_k - \boldsymbol{\mu}_G)^\top) \quad (7)$$

This between-class covariance will be changed like Eq. 8.

$$\Sigma_B^0 = \frac{1}{K} \sum_{k=1}^K \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top \quad (\because \boldsymbol{\mu}_G^0 = \mathbf{0}_D) \quad (8)$$

where $\Sigma_B^0 \in \mathbb{R}^{D \times D}$ denotes the between-class covariance in the zero-mean NC.

Within-class Covariance. In the common neural collapse, the within-class covariance $\Sigma_W \in \mathbb{R}^{D \times D}$ is as Eq. 9.

$$\Sigma_W = \frac{1}{K} \sum_{k=1}^K \left(\frac{1}{N_k} \sum_{i=1}^{N_k} (\mathbf{h}_{i,k} - \boldsymbol{\mu}_k) (\mathbf{h}_{i,k} - \boldsymbol{\mu}_k)^\top \right) \quad (9)$$

The zero within-class covariance Σ_W^0 is calculated in the same way with the within-class covariance (Eq. 10).

$$\Sigma_W^0 = \Sigma_W \quad (10)$$

where $\Sigma_W^0 \in \mathbb{R}^{D \times D}$ denotes the within-class covariance in the zero-mean NC.

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

(ZNC1) Variability collapse: $\Sigma_W^0 \rightarrow \mathbf{0}_D$

(NC2) Convergence to simplex ETF:

$$||\boldsymbol{\mu}_k - \boldsymbol{\mu}_G||_2 - ||\boldsymbol{\mu}_{k'} - \boldsymbol{\mu}_G||_2 \rightarrow 0 \quad \forall_{k,k'}$$

$$\langle \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\mu}}_{k'} \rangle \rightarrow \frac{K}{K-1} \delta_{k,k'} - \frac{K}{K-1} \quad \forall_{k,k'}$$

where $\tilde{\boldsymbol{\mu}}_k = (\boldsymbol{\mu}_k - \boldsymbol{\mu}_G) / ||\boldsymbol{\mu}_k - \boldsymbol{\mu}_G||_2$ denotes the normalized k -th class mean and $\delta_{k,k'}$ is the Kronecker delta symbol.

(ZNC2) Convergence to orthogonal matrix:

$$||\boldsymbol{\mu}_k||_2 - ||\boldsymbol{\mu}_{k'}||_2 \rightarrow 0 \quad \forall_{k,k'}$$

$$\langle \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\mu}}_{k'} \rangle \rightarrow 0 \quad \forall_{k,k'}$$

where $\tilde{\boldsymbol{\mu}}_k = \boldsymbol{\mu}_k / ||\boldsymbol{\mu}_k||_2$ denotes the normalized k -th class mean.

Table 1: Test Accuracy (%) at zero error and last epoch.

Arch \ Dataset	MNIST		Fashion MNIST		SVHN		CIFAR10		CIFAR100	
	Zero	Last	Zero	Last	Zero	Last	Zero	Last	Zero	Last
VGG	99.40	99.56	92.92	93.31	93.82	94.53	87.85	88.65	63.03	63.85
ResNet	99.32	99.71	93.29	93.64	94.64	95.70	88.72	89.44	66.19	66.21
DenseNet	99.65	99.70	94.18	94.35	95.87	95.93	91.14	91.19	77.19	76.56
VGG [†]	99.49	99.56	93.25	93.56	93.91	94.51	87.54	88.44	63.14	63.95
ResNet [†]	99.52	99.64	93.96	93.77	92.45	95.21	87.49	89.45	66.03	65.11
DenseNet [†]	99.64	99.61	93.92	94.21	95.21	95.41	88.20	88.43	70.66	70.95
VGG+FNO [†]	99.49	99.49	92.98	93.31	93.89	94.16	88.14	89.10	62.59	64.17
ResNet+FNO [†]	99.40	99.67	93.90	94.28	92.30	95.41	87.61	89.59	67.15	66.20
DenseNet+FNO [†]	99.62	99.60	93.91	94.08	95.11	95.22	88.64	88.76	72.23	72.54

(NC3) Convergence to self-duality:

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

where $\dot{\mathbf{M}} = [\boldsymbol{\mu}_k - \boldsymbol{\mu}_G, 1 \leq k \leq K] \in \mathbb{R}^{D \times K}$ denotes the matrix obtained by concatenating the class means into the columns of a matrix.

(ZNC3) Convergence to self-duality:

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

where $\dot{\mathbf{M}} = [\boldsymbol{\mu}_k, 1 \leq k \leq K] \in \mathbb{R}^{D \times K}$.

(NC4) Simplification to NCC:

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

(ZNC4) Simplification to NCC:

$$\arg \max_{k'} \langle \mathbf{w}_{k'}, \mathbf{h} \rangle \rightarrow \arg \min_{k'} \|\mathbf{h} - \boldsymbol{\mu}_{k'}\|_2$$

B.1 ZERO-MEAN NEURAL COLLAPSE IN IMAGE CLASSIFICATION BENCHMARKS

Datasets. To reproduce the analyses environment for zero-mean neural collapse, we follow (Papayan et al., 2020). We utilize the MNIST (Deng, 2012), Fashion MNIST (Xiao et al., 2017), SVHN (Netzer et al., 2011), CIFAR10, and CIFAR100 datasets (Krizhevsky et al., 2009). To eliminate imbalance factor in the datasets, we subsamples 5,000 samples per class for MNIST and 4,600 samples per class for SVHN. The remaining datasets are already balanced: 6,000 examples for Fashion MNIST, 5,000 examples for CIFAR10, and 500 examples for CIFAR100. There was no data augmentation except normalization.

Architectures and Implementation Details. We train three types of convolutional networks: VGG (Simonyan & Zisserman, 2014), ResNet (He et al., 2015), and DenseNet (Huang et al., 2017). Following (Papayan et al., 2020), we minimize the cross-entropy loss using SGD with momentum 0.9 and the weight decay is set to $5e^{-4}$. The batch size and the number of epochs are set to 256 and 350, respectively. The initial learning is set differently for datasets and architectures and it is decayed by 10 at 1/3 and 2/3. We itemize architectures for each dataset with the initial learning rate:

- **MNIST:** VGG11 (0.06786), ResNet18 (0.013296), and DenseNet40 (0.094015)
- **Fashion MNIST:** VGG11 (0.009597), ResNet18 (0.13025), and DenseNet250 (0.009597)
- **SVHN:** VGG11 (0.094015), ResNet18 (0.009597), and DenseNet40 (0.06786)
- **CIFAR10:** VGG13 (0.048982), ResNet18 (0.06786), and DenseNet40 (0.094015)
- **CIFAR100:** VGG13 (0.180451), ResNet50 (0.13025), and DenseNet250 (0.13025)

Results and Analyses. We visualize the observations of various experiments. Best descriptions in captions of each figure.

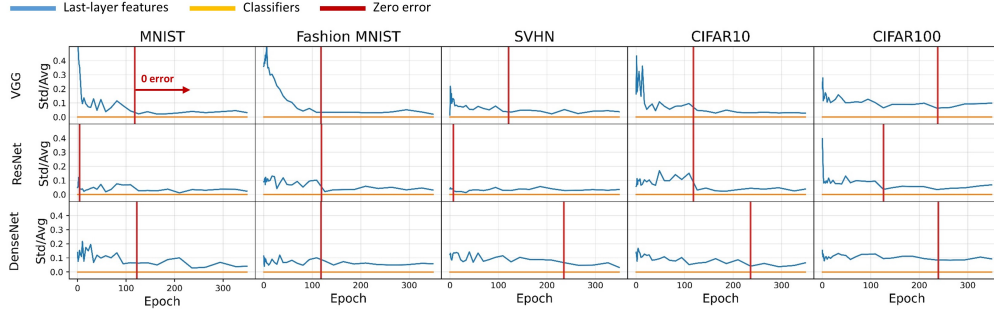


Figure 1: Class mean vectors become equinorm. The y-axis in each cell displays the coefficient of variation of averaged normalization of class means and class weight vectors. The blue lines represent $\text{Std}_k(\|\mu_k\|) / \text{Avg}_k(\|\mu_k\|)$ where μ_k , $1 \leq k \leq K$ denotes the class means of the last-layer features of the training input samples. The orange lines show $\text{Std}_k(\|w_k\|) / \text{Avg}_k(\|w_k\|)$ where w_k is the class weight vector of the k -th class. As training advances, the coefficients of variation for both class means and class weight vectors decrease.

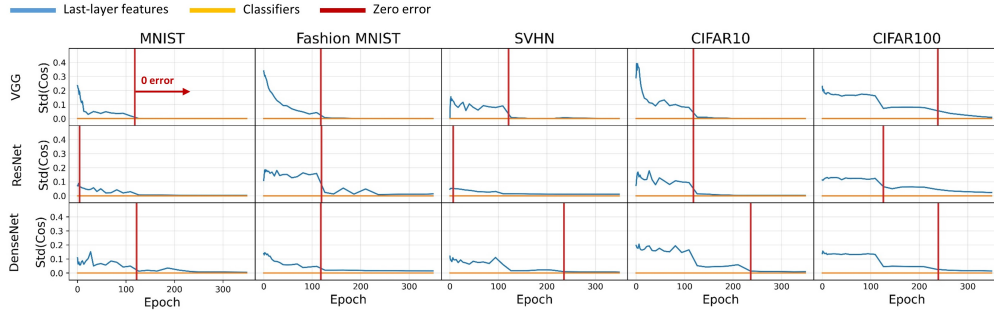


Figure 2: Class means converge to equiangularity. The y-axis in each cell displays the standard deviation of the cosine similarity between pairs of centered class means and class weight vectors across all pairs of classes k and k' , $\forall k \neq k'$. w_k , μ_k , and μ_G are as in Figure 1. The blue lines represent $\cos_\mu(k, k') = (\mu_k - \mu_G)(\mu_{k'} - \mu_G)^T / (\|\mu_k - \mu_G\|_2 \|\mu_{k'} - \mu_G\|_2) \forall k \neq k'$. The orange lines show $\cos_w(k, k') = w_k w_{k'}^T / (\|w_k\|_2 \|w_{k'}\|_2) \forall k \neq k'$. As training advances, the standard deviations of $\cos_\mu(k, k')$ and $\cos_w(k, k')$ converge to zero, signifying equiangularity.

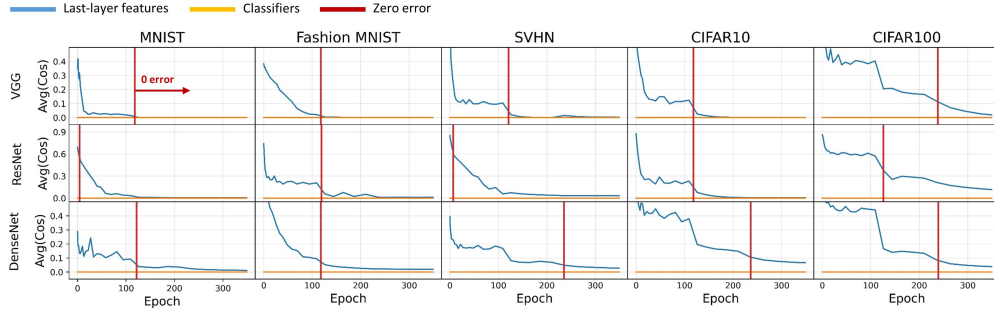


Figure 3: Class means converge to maximal-angle equiangularity. The y-axis in each cell displays the average of cosine similarity of between pairs of centered class means and between pairs of class weight vectors across all pairs of classes k and k' , $\forall k \neq k'$. The blue lines represent $\text{Avg}_{k,k'} \cos_\mu(k, k')$, $\forall k \neq k'$. The orange lines show $\text{Avg}_{k,k'} \cos_w(k, k')$, $\forall k \neq k'$. This result represents the maximum separation achievable among globally centered, equiangular vectors.

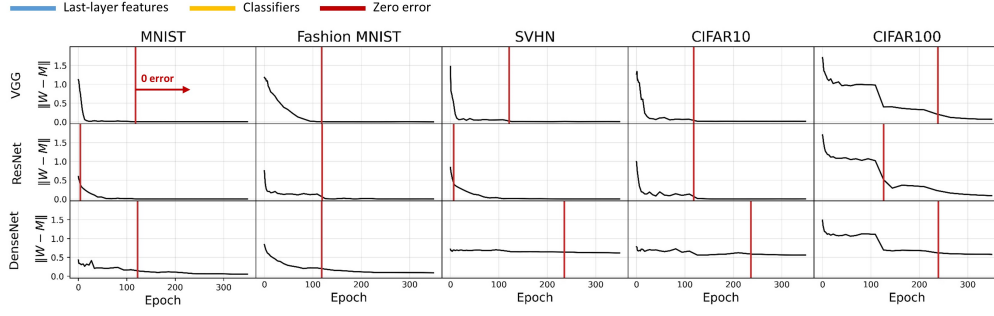


Figure 4: Class means approach to class weight vectors. The y-axis in each cell displays the distance between the class weight vectors and the normalized, centered class means. The lines represent the quantities of $\|\mathbf{W} - \hat{\mathbf{M}}\|_F^2$. \mathbf{W} is the class weight matrix. $\hat{\mathbf{M}} = \mathbf{M}/\|\mathbf{M}\|_F$ where $\mathbf{M} = (\boldsymbol{\mu}_k)_{1 \leq k \leq K} \in \mathbb{R}^{D \times K}$ is the matrix whose columns consist of the centered class means. This distance decreases as training advances. This result indicates that the centered class means are proportionally related to the class weight vectors like a self-dual manner.

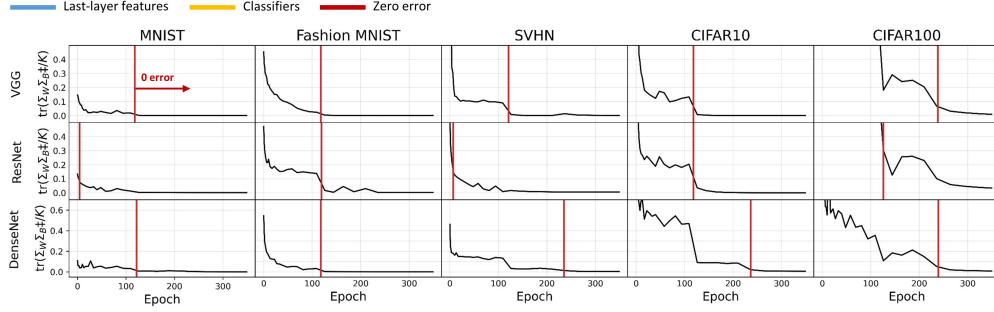


Figure 5: Within-class variation collapses. The y-axis in each cell displays the magnitude of the between-class covariance compared with the within-class covariance of the last-layer features. The lines represent $\text{tr}(\Sigma_W \Sigma_B^\dagger / K)$, where $\text{tr}(\cdot)$ denotes the trace operator, $[\cdot]^\dagger$ indicates Moore-Penrose pseudoinverse, and K is the number of classes. Σ_W and Σ_B are as in Table 1 of the main paper. This magnitude decreases as training advances. This result indicates that collapse of within-class variations occurs.

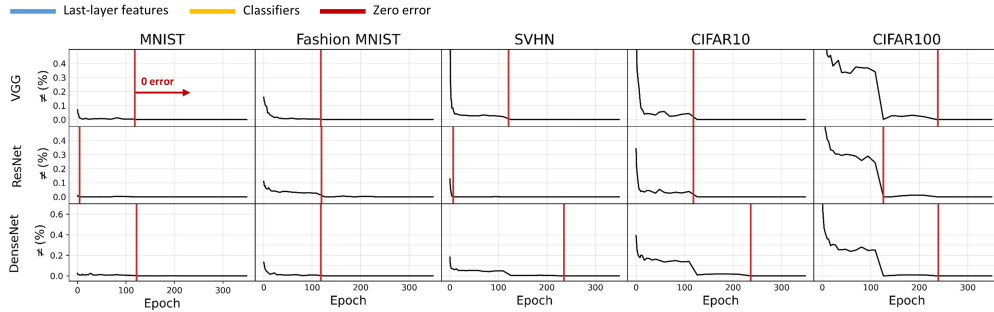


Figure 6: Classifier works in a similar way to NCC. The y-axis in each cell displays the percentage of input samples in the training set where there is a mismatch between the classifier's output and the result that would have been obtained by selecting $\arg \min_k \|\mathbf{h} - \boldsymbol{\mu}_k\|_2$ where \mathbf{h} is a last-layer feature and $\boldsymbol{\mu}_k, \forall 1 \leq k \leq K$ are the class means of last-layer features. The proportion of disagreement decreases as training advances. This result indicates that the classifier's decision is gradually simplified to the nearest class mean decision rule.

C PROOF FOR THEOREM 1

Note that the constraint constrained optimization problem of OLPM is reduced to an entangled constraint. We consider the k -th ($1 \leq k \leq K$ and $K \geq 2$) problem as:

$$\begin{aligned} \min_H \quad & \frac{1}{N} \sum_{i=1}^{n_k} \mathcal{L}_{ce}(\mathbf{h}_{k,i}, \mathbf{Q}^*), \\ \text{s.t.} \quad & \|\mathbf{h}_{k,i}\|^2 - 2 \sum_{j \neq k}^K \mathbf{h}_{k,i}^\top \mathbf{w}_j \leq E_{HW}, \quad 1 \leq i \leq n_k, \end{aligned} \quad (11)$$

where \mathbf{Q}^* is the fixed non-negative orthogonal classifier. The problem above is convex as the objective is a sum of affine functions and log-sum-exp functions with convex constraints. We have the Lagrange function as:

$$\tilde{L} = \frac{1}{N} \sum_{i=1}^{n_k} -\log \frac{\exp(\mathbf{h}_{k,i}^\top \mathbf{q}_k^*)}{\sum_{j=1}^K \exp(\mathbf{h}_{k,i}^\top \mathbf{q}_j^*)} + \sum_{i=1}^{n_k} \lambda_i \left(\|\mathbf{h}_{k,i}\|^2 - 2 \sum_{j \neq k}^K \mathbf{h}_{k,i}^\top \mathbf{q}_j - E_{HW} \right), \quad (12)$$

where λ_i is the Langrange multiplier. We have its gradient with respect to $\mathbf{h}_{k,i}$ as:

$$\frac{\partial \tilde{L}}{\partial \mathbf{h}_{k,i}} = -\frac{(1-p_k) \mathbf{q}_k^*}{N} + \frac{1}{N} \sum_{j \neq k}^K p_j \mathbf{q}_j^* + 2\lambda_i \left(\mathbf{h}_{k,i} - \sum_{j \neq k}^K \mathbf{q}_j^* \right), \quad 1 \leq i \leq n_k. \quad (13)$$

First we consider the case when $\lambda_i = 0$. $\partial \tilde{L} / \partial \mathbf{h}_{k,i} = 0$ gives the following equation:

$$\begin{aligned} \frac{(1-p_k) \mathbf{q}_k^*}{N} &= \frac{1}{N} \sum_{j \neq k}^K p_j \mathbf{q}_j^* \\ \sum_{j \neq k}^K p_j \mathbf{q}_k^* &= \sum_{j \neq k}^K p_j \mathbf{q}_j^*. \end{aligned} \quad (14)$$

Multiplying \mathbf{q}_k^* by both sides of the equation, we should have:

$$\sum_{j \neq k}^K p_j = 0 \quad (\because \mathbf{q}_k^{*\top} \mathbf{q}_{k'}^* = \delta_{k,k'}, \quad \forall k, k' \in [1, K]), \quad (15)$$

which contradicts with $p_j > 0, \forall 1 \leq i \leq K$ when the ℓ_2 norm of $\mathbf{h}_{k,i}$ is constrained and \mathbf{Q}^* has a fixed ℓ_2 norm. So we have $\lambda_i > 0$ and according to the KKT condition:

$$\|\mathbf{h}_{k,i}\|^2 - 2 \sum_{j \neq k}^K \mathbf{h}_{k,i}^\top \mathbf{q}_j = E_{HW}. \quad (16)$$

Then we have the equation:

$$\frac{\partial \tilde{L}}{\partial \mathbf{h}_{k,i}} = \frac{1}{N} \sum_{j \neq k}^K p_j (\mathbf{q}_j^* - \mathbf{q}_k^*) + 2\lambda_i \left(\mathbf{h}_{k,i}^* - \sum_{j \neq k}^K \mathbf{q}_j^* \right) = 0, \quad (17)$$

where $\mathbf{h}_{k,i}^*$ is the optimal solution of $\mathbf{h}_{k,i}$. Multiplying $\mathbf{q}_{j'}^*$ ($j' \neq k$) by both sides of Eq. 17, we get:

$$\frac{p_{j'}}{N} + 2\lambda_i (\mathbf{h}_{k,i}^{*\top} \mathbf{q}_{j'}^* - 1) = 0. \quad (18)$$

Since $p_{j'} > 0$, we have $\mathbf{h}_{k,i}^{*\top} \mathbf{q}_{j'}^* < 1$. Then for any pair $j, j' \neq k$, we have:

$$\frac{p_j}{p_{j'}} = \frac{\exp(\mathbf{h}_{k,i}^{*\top} \mathbf{q}_j^*)}{\exp(\mathbf{h}_{k,i}^{*\top} \mathbf{q}_{j'}^*)} = \frac{\exp(\mathbf{h}_{k,i}^{*\top} \mathbf{q}_j^* - 1)}{\exp(\mathbf{h}_{k,i}^{*\top} \mathbf{q}_{j'}^* - 1)} = \frac{\mathbf{h}_{k,i}^{*\top} \mathbf{q}_j^* - 1}{\mathbf{h}_{k,i}^{*\top} \mathbf{q}_{j'}^* - 1}. \quad (19)$$

Considering that the function $f(x) = x/\exp(x)$ is monotonically increasing when $x < 0$, we have:

$$\frac{\mathbf{h}_{k,i}^* \mathbf{q}_j^* - 1}{\exp(\mathbf{h}_{k,i}^* \mathbf{q}_j^* - 1)} = \frac{\mathbf{h}_{k,i}^* \mathbf{q}_{j'}^* - 1}{\exp(\mathbf{h}_{k,i}^* \mathbf{q}_{j'}^* - 1)} \quad (20)$$

$$\mathbf{h}_{k,i}^* \mathbf{q}_j^* = \mathbf{h}_{k,i}^* \mathbf{q}_{j'}^* = C, p_j = p_{j'} = p, \forall j, j' \neq k,$$

where C and p are constants. From Eq. 18, we have:

$$p = -2N\lambda_i(C - 1), \quad (21)$$

and

$$1 - p_k = (K - 1)p = -2N\lambda_i(K - 1)(C - 1). \quad (22)$$

From Eq. 17, we have:

$$\begin{aligned} \mathbf{h}_{k,i}^* &= \frac{1}{2N\lambda_i} \left((1 - p_k) \mathbf{q}_k^* - \sum_{j \neq k}^K p_j \mathbf{q}_j^* \right) + \sum_{j \neq k}^K \mathbf{q}_j^* \\ &= \frac{1}{2N\lambda_i} \left(-2N\lambda_i(K - 1)(C - 1) \mathbf{q}_k^* + 2N\lambda_i(C - 1) \sum_{j \neq k}^K \mathbf{q}_j^* \right) + \sum_{j \neq k}^K \mathbf{q}_j^* \\ &= -(K - 1)(C - 1) \mathbf{q}_k^* + C \sum_{j \neq k}^K \mathbf{q}_j^*, \end{aligned} \quad (23)$$

From the theorem in (Ji et al., 2022), the margin of a single feature $\mathbf{h}_{k,i}^*$ is defined:

$$\mathcal{M}_{k,i} := \mathbf{h}_{k,i}^* \mathbf{q}_k^* - \max_{j \neq k} \mathbf{h}_{k,i}^* \mathbf{q}_j^*. \quad (24)$$

Multiplying \mathbf{q}_k^* by both sides of Eq. 23, we should have:

$$\mathbf{h}_{k,i}^* \mathbf{q}_k^* = -(K - 1)(C - 1) \leq K - 1 \quad (\because 0 \leq C < 1), \quad (25)$$

and

$$\mathbf{h}_{k,i}^* \mathbf{q}_j^* = \mathbf{h}_{k,i}^* \mathbf{q}_{j'}^* = C \geq 0, \quad \forall j, j' \neq k \quad (26)$$

When the equality holds, we have:

$$\begin{aligned} \mathbf{h}_{k,i}^* \mathbf{q}_k^* &= K - 1 \\ \mathbf{h}_{k,i}^* \mathbf{q}_j^* &= 0, \quad \forall j \neq k, \end{aligned} \quad (27)$$

which is equivalent to Eq. 8 of the main paper and concludes the proof.

D ALGORITHMS FOR SECTION 6.1 AND 6.2

In this section, we provide algorithms that are more detailed and may help with implementation for *masked and weighted softmax with FNO classifier in a task T_t* in section 6.1 and for *arc-mixup with FNO classifier and feature masking in a mini-batch \mathbb{B}* in section 6.2. These algorithms represent the whole process in training classification models in continual learning and imbalanced learning, respectively.

Algorithm 1 Masked and Weighted Softmax with FNO classifier in a Task T_t

Input: $(\mathbf{X}_t, \mathbf{Y}_t), \mathbf{Q}$

Output: $\mathbf{P} \in \mathbb{R}^{N_t \times K}$

- 1: $\mathbf{H} \leftarrow \text{RELU}(f_\theta(\mathbf{X}_t))$ # get features from \mathbf{X}_t as Eq. 1 of the main paper.
 - 2: $\mathbb{K} = \{c_i \mid c_i \text{ of } \mathbf{X}_t\}$ # initialize a set of class labels in the task T_t .
 - 3: $M^{(-\infty)} = (m_{i,j})_{1 \leq i \leq N_t, 1 \leq j \leq K}$, where $m_{i,j} = -\infty$ if $j \notin \mathbb{K}$ otherwise 1 # initialize $M^{(-\infty)}$
 - 4: $\mathbf{P} = \text{W-SOFTMAX}(M^{(-\infty)} \odot \text{MATMUL}(\mathbf{Q}, \mathbf{H}))$ # get the confidence of \mathbf{H}
-

Algorithm 2 Arc-mixup with FNO classifier and feature masking in a mini-batch \mathbb{B}

Input: $(\mathbf{X}, \mathbf{Y}) \in \mathbb{B}, \mathbf{Q}$

Output: $\mathbf{P} \in \mathbb{R}^{|\mathbb{B}| \times K}$

- 1: $(\hat{\mathbf{X}}, \hat{\mathbf{Q}}) \leftarrow \text{ArcMixup}(\mathbf{X}, \mathbf{Q})$ # mixup input samples and class vectors as Eq. 11 of the main paper
 - 2: $\mathbb{K} = \{c_i \mid c_i \in \mathbb{B}, \forall 1 \leq i \leq |\mathbb{B}|\}$ # initialize a set of class labels in the mini-batch
 - 3: $\hat{\mathbb{J}} = \bigcup_{i=1}^{|\mathbb{B}|} \mathbb{J}_i$ # initialize an index set including all index sets of \mathbf{Q}
 - 4: $M^{(0)} = (m_{i,j})_{1 \leq i \leq |\mathbb{B}|, 1 \leq j \leq D}$, where $m_{i,j} = \mathbf{1}_{j \in \hat{\mathbb{J}}}$ # initialize a zero mask $M^{(0)}$
 - 5: $\hat{\mathbf{H}} \leftarrow \text{RELU}(f_\theta(\hat{\mathbf{X}}))$ # get features from $\hat{\mathbf{X}}$ as Eq. 1 of the main paper.
 - 6: $\mathbf{P} = \text{MATMUL}(\hat{\mathbf{Q}}, \text{LAYERNORM}(M^{(0)} \odot \hat{\mathbf{H}}))$ # get the confidence of $\hat{\mathbf{H}}$
-

E ADDITIONAL IMPLEMENTATION DETAILS FOR EXPERIMENTS IN SECTION 7

E.1 CONTINUAL LEARNING IN SPLIT DATASETS

Dataset Explanation. Following (Hsu et al., 2018; Van de Ven & Tolias, 2019; Buzzega et al., 2020; Kim et al., 2023), the split datasets of the MNIST, CIFAR10, CIFAR100, and Tiny-ImageNet are described as:

- **S-MNIST, S-CIFAR10:** $T_1(0-1), T_2(2-3), T_3(4-5), T_4(6-7), T_5(8-9)$
- **S-CIFAR100:** $T_1(0-9), T_2(10-19), T_3(20-29), T_4(30-39), T_5(40-49), T_6(50-59), T_7(60-69), T_8(70-79), T_9(80-89), T_{10}(90-99)$
- **S-Tiny-ImageNet:** $T_1(0-9), T_2(10-19), T_3(20-29), T_4(30-39), T_5(40-49), T_6(50-59), T_7(60-69), T_8(70-79), T_9(80-89), T_{10}(90-99), T_{11}(100-109), T_{12}(110-119), T_{13}(120-129), T_{14}(130-139), T_{15}(140-149), T_{16}(150-159), T_{17}(160-169), T_{18}(170-179), T_{19}(180-189), T_{20}(190-199)$

where $T_t(i-j)$ indicates that the t -th task has class labels from i to j and all classes are sequentially divided to each task. In CIFAR10, for instance, we can convert the class index in each task as below:

- $T_1(0-1)$: {airplane, automobile}
- $T_2(2-3)$: {bird, cat}
- $T_3(4-5)$: {deer, dog}
- $T_4(6-7)$: {frog, horse}
- $T_5(8-9)$: {ship, truck}

Table 2: Hyperparameter settings on all experiments in Continual Learning Benchmarks. (Table credit: (Kim et al., 2023))

Method	Buffer	S-MNIST			S-CIFAR10			S-CIFAR100			S-Tiny-ImageNet		
SGD	-	$lr: 0.03$			$lr: 0.1$			$lr: 0.03$			$lr: 0.03$		
ER	200	$lr: 0.01$			$lr: 0.1$			$lr: 0.1$			$lr: 0.1$		
	500	$lr: 0.1$			$lr: 0.1$			$lr: 0.1$			$lr: 0.03$		
	5120	$lr: 0.1$			$lr: 0.1$			$lr: 0.1$			$lr: 0.1$		
DER++	200	$lr: 0.03$	$\alpha: 0.2$	$\beta: 1.0$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 1.0$
	500	$lr: 0.03$	$\alpha: 1.0$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.2$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.2$	$\beta: 0.5$
	5120	$lr: 0.1$	$\alpha: 0.2$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 1.0$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 0.5$	$lr: 0.03$	$\alpha: 0.1$	$\beta: 0.5$

Hyperparameter Settings. We have the same settings to (Buzzega et al., 2020; Kim et al., 2023). To summarize them, we borrow the table of hyperparameter settings in (Kim et al., 2023) (Please refer to Table 2).

Task-Incremental Learning (Task-IL) and Class-Incremental Learning (Class-IL). As described in (Hsu et al., 2018; Van de Ven & Tolias, 2019), the sequence of tasks can be modeled in two ways: Task-IL and Class-IL. Both divide datasets and train models for tasks in the same way. The only one difference lies in the evaluation system. To be more specific, Task-IL provides task information the task information during prediction, so the models classify test samples into the target task’s classes. In contrast, Class-IL does not provide task information, so the models must predict the test samples as one of the total classes, regardless of the target task. We trained and validated our models in continual learning benchmarks using these training strategies.

E.2 IMBALANCED LEARNING IN LONG-TAILED DATASETS

Dataset Explanation. Following (Zhong et al., 2021) and (Zhou et al., 2020), the long-tailed datasets of CIFAR10, CIFAR100, ImageNet (Russakovsky et al., 2015), and Places365 (Zhou et al., 2017) are described as:

- **CIFAR10-LT** comprises ten imbalanced classes, subsampled at an exponentially decreasing rate from the initial class of CIFAR10 as mentioned in (Zhong et al., 2021).
- **CIFAR100-LT** includes one hundred imbalanced classes, constructed in the same way as CIFAR10-LT.
- **ImageNet-LT** is a long-tailed dataset for large-scale object classification, derived from the ImageNet. The sampling is based on Pareto distribution with a power value $\alpha = 5$ and the classes have varying cardinality from 5 to 1,280. Therefore, it contains 115.8K images sorted into 1,000 classes.
- **Places-LT** is an extended version of the large-scale scene classification dataset Places. The classes differ in their cardinality, ranging from 5 to 4,980, and therefore, it contains 184.5K images from 365 classes.

Architectures. We utilize a ResNet32 (Zhong et al., 2021) consisting of three residual blocks, with each output dimensions of 16, 32, and 64, respectively, for the CIFAR10-LT dataset. For the CIFAR100-LT dataset, each output dimension is twice that of the CIFAR10-LT dataset. Unlike the ResNet architecture for ImageNet, the kernel size, stride, and padding are set to 3, 1, and 1, respectively for the first convolutional layer. ResNet50 and 152 are same to He et al. (2015).

Hyperparameters Settings. We employed ResNet32 for CIFAR10/100-LT and trained the model with 128 mini-batches, utilizing SGD with momentum of 0.9 and weight decay of $2e-4$, for 200 epochs. The learning rate was warmed up from 0.02 to the initial learning rate in a linear fashion, divided by 0.1 at epochs 160 and 180. For the other datasets, we utilized ResNet50 and 152, trained the models using SGD with momentum of 0.9 and weight decay of $5e-4$, and updated the learning rate using a cosine annealing scheduler. Additionally, we differently set mixup alpha of mixup and arc-mixup as the datasets: $\alpha = 1.0$ for both in CIFAR10/100-LT, $\alpha = 0.2$ for mixup and $\alpha = 5.0$ for arc-mixup in ImageNet, and $\alpha = 0.2$ for both in Places-LT.

F ADDITIONAL EXPERIMENTAL RESULTS.

F.1 CONTINUAL LEARNING IN SPLIT DATASETS

Table 3: Classification results for standard continual learning benchmarks including the performance of commonly used methods in these benchmarks. Best in bold for each buffer setting. (Final Average Accuracy \uparrow (%): $mean_{std}$)

Buffer	Method	S-MNIST		S-CIFAR-10		S-Tiny-ImageNet	
		Class-IL	Task-IL	Class-IL	Task-IL	Class-IL	Task-IL
-	JOINT	95.57 _{0.24}	99.51 _{0.07}	92.20 _{0.15}	98.31 _{0.12}	59.99 _{0.19}	82.04 _{0.10}
	SGD	19.60 _{0.04}	94.94 _{2.18}	19.62 _{0.05}	61.02 _{3.33}	7.92 _{0.26}	18.31 _{0.68}
	oEWC (Schwarz et al., 2018)	20.46 _{1.01}	98.39 _{0.48}	19.49 _{0.12}	68.29 _{3.92}	7.58 _{0.10}	19.20 _{0.31}
	SI (Zenke et al., 2017)	19.27 _{0.30}	96.00 _{2.04}	19.48 _{0.17}	68.05 _{5.91}	6.58 _{0.31}	36.32 _{0.13}
	LwF (Li & Hoiem, 2017)	19.62 _{0.01}	94.11 _{3.01}	19.61 _{0.05}	63.29 _{2.35}	8.46 _{0.22}	15.85 _{0.58}
	PNN (Rusu et al., 2016)		99.23 _{0.20}		95.13 _{0.72}		67.84 _{0.29}
	ER (Riemer et al., 2018)	80.43 _{1.89}	97.86 _{0.35}	44.79 _{1.86}	91.19 _{0.94}	8.49 _{0.16}	38.17 _{2.00}
	GEM (Lopez-Paz & Ranzato, 2017)	80.11 _{1.54}	97.78 _{0.25}	25.54 _{0.76}	90.44 _{0.94}		
	A-GEM (Chaudhry et al., 2018)	45.72 _{4.26}	98.61 _{0.24}	20.04 _{0.34}	83.88 _{1.49}	8.07 _{0.08}	22.77 _{0.03}
	iCaRL (Rebuffi et al., 2017)	70.51 _{0.53}	98.28 _{0.09}	49.02 _{3.20}	88.99 _{2.13}	7.53 _{0.79}	28.19 _{1.47}
200	FDR (Benjamin et al., 2018)	79.43 _{3.26}	97.66 _{0.18}	30.91 _{2.74}	91.01 _{0.68}	8.70 _{0.19}	40.36 _{0.68}
	GSS (Aljundi et al., 2019)	38.90 _{2.49}	95.02 _{1.85}	39.07 _{5.59}	88.80 _{2.89}		
	HAL (Chaudhry et al., 2021)	84.70 _{0.87}	97.96 _{0.21}	32.36 _{2.70}	82.51 _{3.20}		
	DER (Buzzega et al., 2020)	84.55 _{1.64}	98.80 _{0.15}	61.93 _{1.79}	91.40 _{0.92}	11.87 _{0.78}	40.22 _{0.67}
	DER++ (Buzzega et al., 2020)	85.61 _{1.40}	98.76 _{0.28}	64.88 _{1.17}	91.92 _{0.60}	10.96 _{1.17}	40.87 _{1.16}
	ER [†]	78.27 _{1.37}	97.73 _{0.26}	49.38 _{2.15}	91.54 _{0.81}	8.58 _{0.19}	38.39 _{0.72}
	FNOER [†]	78.59 _{1.27}	98.03 _{0.22}	50.56 _{1.52}	91.79 _{0.64}	9.21 _{0.28}	40.46 _{1.15}
	ERM [†]	82.98 _{1.03}	98.13 _{0.16}	61.75 _{6.07}	91.39 _{2.13}	15.47 _{0.67}	44.11 _{0.50}
	FNOERM [†]	84.26 _{1.16}	98.45 _{0.19}	63.84 _{1.47}	92.03 _{0.52}	17.31 _{0.74}	44.76 _{0.90}
	DER++ [†]	85.64 _{1.02}	98.84 _{0.11}	63.67 _{1.01}	91.61 _{0.73}	11.59 _{1.07}	41.00 _{0.88}
500	FNODER++ [†]	83.29 _{1.64}	97.67 _{0.33}	64.15 _{1.50}	92.52 _{0.55}	13.75 _{1.28}	45.52 _{1.29}
	DERMR++ [†]	84.45 _{0.88}	99.03 _{0.09}	66.35 _{1.52}	93.17 _{0.54}	13.21 _{0.56}	49.75 _{0.99}
	FNODERM [†]	86.27 _{0.88}	99.11 _{0.08}	67.53 _{1.25}	93.98 _{0.39}	18.44 _{0.94}	53.06 _{0.67}
	ER (Riemer et al., 2018)	86.12 _{1.89}	99.04 _{0.18}	57.74 _{0.27}	93.61 _{0.27}	9.99 _{0.29}	48.64 _{0.46}
	GEM (Lopez-Paz & Ranzato, 2017)	85.99 _{1.35}	98.71 _{0.20}	26.20 _{1.26}	92.16 _{0.69}		
	A-GEM (Chaudhry et al., 2018)	46.66 _{5.85}	98.93 _{0.21}	22.67 _{0.57}	89.48 _{1.45}	8.06 _{0.04}	25.33 _{0.49}
	iCaRL (Rebuffi et al., 2017)	70.10 _{1.08}	98.32 _{0.07}	47.55 _{3.95}	88.22 _{2.62}	9.38 _{1.53}	31.55 _{3.27}
	FDR (Benjamin et al., 2018)	85.87 _{4.04}	97.54 _{1.90}	28.71 _{3.23}	93.29 _{0.59}	10.54 _{0.21}	49.88 _{0.71}
	GSS (Aljundi et al., 2019)	49.76 _{4.73}	97.71 _{0.53}	49.73 _{4.78}	91.02 _{1.57}		
	HAL (Chaudhry et al., 2021)	87.21 _{0.49}	98.03 _{0.22}	41.79 _{4.46}	84.54 _{2.36}		
5120	DER (Buzzega et al., 2020)	90.54 _{1.18}	98.84 _{0.13}	70.51 _{1.67}	93.40 _{0.39}	17.75 _{1.14}	51.78 _{0.88}
	DER++ (Buzzega et al., 2020)	91.00 _{1.49}	98.94 _{0.27}	72.70 _{1.36}	93.88 _{0.50}	19.38 _{1.41}	51.91 _{0.68}
	ER [†]	85.99 _{1.52}	99.14 _{0.07}	62.38 _{1.40}	94.12 _{0.31}	10.12 _{0.22}	48.06 _{0.80}
	FNOER [†]	85.92 _{1.76}	99.10 _{0.15}	63.41 _{1.36}	93.99 _{0.45}	11.07 _{0.41}	45.77 _{0.46}
	ERM [†]	89.35 _{0.59}	99.20 _{0.16}	70.64 _{1.28}	94.22 _{0.41}	20.43 _{0.38}	53.21 _{0.84}
	FNOERM [†]	89.42 _{0.72}	99.16 _{0.17}	71.43 _{0.95}	94.38 _{0.43}	22.41 _{0.57}	52.60 _{0.58}
	DER++ [†]	91.01 _{0.46}	98.95 _{0.07}	73.15 _{0.80}	94.07 _{0.39}	19.82 _{0.87}	52.24 _{0.94}
	FNODER++ [†]	90.56 _{0.64}	98.06 _{0.16}	72.70 _{1.21}	94.37 _{0.67}	19.42 _{0.72}	53.35 _{0.95}
	DERMR++ [†]	83.10 _{1.22}	99.08 _{0.09}	71.85 _{3.76}	94.28 _{1.49}	17.71 _{0.58}	59.86 _{1.08}
	FNODERM [†]	86.75 _{0.75}	99.00 _{0.10}	74.77 _{0.66}	95.56 _{0.16}	22.45 _{0.36}	59.87 _{1.91}
5120	ER (Riemer et al., 2018)	93.40 _{1.29}	99.33 _{0.22}	82.47 _{0.52}	96.98 _{0.17}	27.40 _{0.31}	67.29 _{0.23}
	GEM (Lopez-Paz & Ranzato, 2017)	95.11 _{0.87}	99.44 _{0.12}	25.26 _{3.46}	95.55 _{0.02}		
	A-GEM (Chaudhry et al., 2018)	54.24 _{6.49}	98.93 _{0.20}	21.99 _{2.29}	90.10 _{2.09}	7.96 _{0.13}	26.22 _{0.65}
	iCaRL (Rebuffi et al., 2017)	70.60 _{1.03}	98.32 _{0.11}	55.07 _{1.55}	92.23 _{0.84}	14.08 _{1.92}	40.83 _{3.11}
	FDR (Benjamin et al., 2018)	87.47 _{3.15}	97.79 _{1.33}	19.70 _{0.07}	94.32 _{0.97}	28.97 _{0.41}	68.01 _{0.42}
	GSS (Aljundi et al., 2019)	89.39 _{0.75}	98.33 _{0.17}	67.27 _{4.27}	94.19 _{1.15}		
	HAL (Chaudhry et al., 2021)	89.52 _{0.96}	98.35 _{0.17}	59.12 _{4.41}	88.51 _{3.32}		
	DER (Buzzega et al., 2020)	94.90 _{0.57}	99.29 _{0.11}	83.81 _{0.33}	95.43 _{0.33}	36.73 _{0.64}	69.50 _{0.26}
	DER++ (Buzzega et al., 2020)	95.30 _{1.20}	99.47 _{0.07}	85.24 _{0.49}	96.12 _{0.21}	39.02 _{0.97}	69.84 _{0.63}
	ER [†]	93.42 _{1.08}	99.41 _{0.15}	84.31 _{0.38}	97.02 _{0.26}	27.30 _{0.51}	67.69 _{0.33}
5120	FNOER [†]	92.95 _{1.57}	99.41 _{0.10}	84.33 _{1.47}	96.90 _{0.49}	28.38 _{0.46}	65.05 _{0.81}
	ERM [†]	93.51 _{0.60}	99.38 _{0.12}	82.63 _{1.34}	96.45 _{0.27}	35.73 _{0.41}	67.50 _{0.53}
	FNOERM [†]	93.98 _{0.39}	99.47 _{0.09}	82.88 _{1.35}	96.79 _{0.38}	36.90 _{0.41}	66.86 _{0.32}
	DER++ [†]	95.09 _{0.56}	99.50 _{0.08}	85.56 _{0.38}	96.30 _{0.22}	39.66 _{0.89}	69.95 _{0.32}
	FNODER++ [†]	94.99 _{0.74}	99.40 _{0.08}	85.83 _{0.36}	96.68 _{0.14}	31.72 _{1.31}	65.76 _{0.78}
	DERMR++ [†]	93.75 _{0.23}	99.62 _{0.05}	84.71 _{0.65}	96.78 _{0.16}	34.72 _{0.46}	72.40 _{0.25}
	FNODERM [†]	94.26 _{0.24}	99.59 _{0.05}	85.65 _{0.38}	97.20 _{0.13}	38.95 _{0.71}	72.70 _{0.27}

Evaluation Metrics. (Kumari et al., 2022) Final Average Accuracy (A_T) and Forgetting (A_F) where $a_{j,t}$ denotes the test accuracy on the t -task after the model has trained all task up to j .

$$A_T = \frac{1}{T} \sum_{t=1}^T a_{T,t} \quad F_T = \frac{1}{T-1} \sum_{t=1}^T \max_{j \in \{1, \dots, T-1\}} (a_{j,t} - a_{T,t})$$

Table 4: Classification results for standard CL benchmarks. The experiment on S-Tiny-ImageNet was conducted using 5 trials with random seeds while other experiments were conducted using 10 trials with random seeds. Best in bold for each buffer setting. (Final Average Forgetting \downarrow (%): $mean_{std}$)

Buffer	Method	S-MNIST		S-CIFAR-10	
		Class-IL	Task-IL	Class-IL	Task-IL
-	SGD	99.10 _{0.55}	5.15 _{2.74}	96.39 _{0.12}	46.24 _{2.12}
-	oEWC (Schwarz et al., 2018)	97.79 _{1.24}	0.44 _{0.16}	91.64 _{3.07}	29.33 _{3.84}
	SI (Zenke et al., 2017)	98.89 _{0.86}	5.15 _{2.74}	95.78 _{0.64}	38.76 _{0.89}
	LwF (Li & Hoiem, 2017)	99.30 _{0.11}	5.15 _{2.74}	96.69 _{0.25}	32.56 _{0.56}
	PNN (Rusu et al., 2016)		0.00 _{0.00}		0.00 _{0.00}
200	ER (Riemer et al., 2018)	21.36 _{2.46}	0.84 _{0.41}	61.24 _{2.62}	7.08 _{0.64}
	GEM (Lopez-Paz & Ranzato, 2017)	22.32 _{2.04}	1.19 _{0.38}	82.61 _{1.60}	9.27 _{2.07}
	A-GEM (Chaudhry et al., 2018)	66.15 _{6.84}	0.96 _{0.28}	95.73 _{0.20}	16.39 _{0.86}
	iCaRL (Rebuffi et al., 2017)	11.73 _{0.73}	0.28 _{0.08}	28.72 _{0.49}	2.63 _{3.48}
	FDR (Benjamin et al., 2018)	21.15 _{4.18}	0.52 _{0.18}	86.40 _{2.67}	7.36 _{0.03}
	GSS (Aljundi et al., 2019)	74.10 _{3.03}	4.30 _{2.31}	75.25 _{4.07}	8.56 _{1.78}
	HAL (Chaudhry et al., 2021)	14.54 _{1.49}	0.53 _{0.19}	69.11 _{4.21}	12.26 _{0.02}
	DER (Buzzega et al., 2020)	17.66 _{2.10}	0.57 _{0.18}	40.76 _{0.42}	6.57 _{0.20}
	DER++ (Buzzega et al., 2020)	16.27 _{1.73}	0.66 _{0.28}	32.59 _{2.32}	5.16 _{0.21}
	ER [†]	24.09 _{1.61}	0.89 _{0.24}	58.97 _{2.70}	6.49 _{0.90}
	FNOER [†]	24.28 _{1.58}	0.84 _{0.16}	57.78 _{1.94}	6.41 _{0.83}
	ERM [†]	10.22 _{1.20}	0.69 _{0.12}	26.30 _{6.18}	6.18 _{1.56}
	FNOERM [†]	9.88 _{1.79}	0.50 _{0.16}	21.93 _{3.75}	5.91 _{0.69}
	DER++ [†]	16.34 _{1.23}	0.53 _{0.12}	34.80 _{1.68}	6.45 _{1.07}
	FNODER++ [†]	19.61 _{2.07}	2.13 _{0.43}	35.58 _{2.18}	5.40 _{0.81}
	DERMR++ [†]	4.29 _{0.75}	0.32 _{0.09}	23.10 _{1.46}	4.82 _{0.71}
	FNODERM [†]	5.59 _{0.76}	0.25 _{0.05}	27.00 _{2.20}	3.83 _{0.68}
	ER (Riemer et al., 2018)	15.97 _{2.46}	0.39 _{0.20}	45.35 _{0.07}	3.54 _{0.35}
	GEM (Lopez-Paz & Ranzato, 2017)	15.57 _{1.77}	0.54 _{0.15}	74.31 _{4.62}	9.12 _{0.21}
	A-GEM (Chaudhry et al., 2018)	65.84 _{7.24}	0.64 _{0.20}	94.01 _{1.16}	14.26 _{4.18}
500	iCaRL (Rebuffi et al., 2017)	11.84 _{0.73}	0.30 _{0.09}	25.71 _{1.10}	2.66 _{2.47}
	FDR (Benjamin et al., 2018)	13.90 _{5.19}	1.35 _{2.40}	85.62 _{0.36}	4.80 _{0.00}
	GSS (Aljundi et al., 2019)	60.35 _{6.03}	0.89 _{0.40}	62.88 _{2.67}	7.73 _{3.99}
	HAL (Chaudhry et al., 2021)	9.97 _{1.62}	0.35 _{0.21}	62.21 _{4.34}	5.41 _{1.10}
	DER (Buzzega et al., 2020)	9.58 _{1.52}	0.45 _{0.13}	26.74 _{0.15}	4.56 _{0.45}
	DER++ (Buzzega et al., 2020)	8.85 _{1.86}	0.35 _{0.15}	22.38 _{4.41}	4.66 _{1.15}
	ER [†]	16.28 _{1.99}	0.38 _{0.11}	42.83 _{1.89}	3.51 _{0.40}
	FNOER [†]	16.47 _{2.30}	0.48 _{0.16}	41.52 _{1.94}	3.79 _{0.65}
	ERM [†]	7.18 _{1.37}	0.39 _{0.17}	17.31 _{1.58}	3.09 _{0.36}
	FNOERM [†]	6.79 _{1.80}	0.46 _{0.20}	15.60 _{2.61}	3.20 _{0.48}
	DER++ [†]	9.00 _{0.50}	0.41 _{0.05}	23.12 _{1.88}	3.56 _{0.38}
	FNODER++ [†]	9.86 _{0.84}	0.88 _{0.17}	24.12 _{2.00}	3.22 _{0.82}
	DERMR++ [†]	2.07 _{0.89}	0.23 _{0.09}	16.95 _{3.77}	3.03 _{0.76}
	FNODERM [†]	3.81 _{0.85}	0.85 _{0.28}	17.34 _{1.35}	1.97 _{0.37}
	ER (Riemer et al., 2018)	6.08 _{1.84}	0.25 _{0.23}	13.99 _{0.12}	0.27 _{0.06}
	GEM (Lopez-Paz & Ranzato, 2017)	4.30 _{1.16}	0.16 _{0.09}	75.27 _{4.41}	6.91 _{2.33}
	A-GEM (Chaudhry et al., 2018)	55.10 _{10.79}	0.63 _{0.21}	84.49 _{3.08}	11.36 _{1.68}
	iCaRL (Rebuffi et al., 2017)	11.64 _{0.72}	0.26 _{0.06}	24.94 _{0.14}	1.59 _{0.57}
	FDR (Benjamin et al., 2018)	11.58 _{3.97}	0.95 _{1.61}	96.64 _{0.19}	1.93 _{0.48}
	GSS (Aljundi et al., 2019)	7.90 _{1.21}	0.18 _{0.11}	58.11 _{9.12}	7.71 _{2.31}
5120	HAL (Chaudhry et al., 2021)	6.55 _{1.63}	0.13 _{0.07}	27.19 _{7.53}	5.21 _{0.50}
	DER (Buzzega et al., 2020)	4.53 _{0.83}	0.32 _{0.08}	10.12 _{0.80}	2.59 _{0.08}
	DER++ (Buzzega et al., 2020)	4.19 _{1.63}	0.23 _{0.06}	7.27 _{0.84}	1.18 _{0.19}
	ER [†]	6.26 _{1.56}	0.21 _{0.10}	14.44 _{0.65}	0.46 _{0.21}
	FNOER [†]	6.92 _{2.32}	0.20 _{0.11}	14.34 _{2.02}	0.70 _{0.63}
	ERM [†]	3.12 _{1.10}	0.24 _{0.10}	7.57 _{1.31}	0.81 _{0.28}
	FNOERM [†]	2.44 _{0.60}	0.16 _{0.13}	6.16 _{1.03}	0.56 _{0.20}
	DER++ [†]	4.59 _{0.79}	0.24 _{0.10}	7.50 _{0.80}	1.05 _{0.34}
	FNODER++ [†]	4.72 _{0.97}	0.32 _{0.06}	7.68 _{0.71}	0.79 _{0.18}
	DERMR++ [†]	0.90 _{0.18}	0.09 _{0.06}	5.02 _{0.69}	0.71 _{0.27}
	FNODERM [†]	1.47 _{0.25}	0.11 _{0.04}	5.56 _{0.37}	0.43 _{0.13}

F.2 IMBALANCED LEARNING IN LONG-TAILED DATASETS

We present the top-1 test accuracy for three class divisions: *Head-Many* (more than 100 images), *Medium* (20 to 100 images), and *Tail-Few* (less than 20 images). The imbalance factor ω in

CIFAR10/100-LT datasets means the ratio of the number of a head class n_{max} to the number of a tail class n_{min} , i.e., $\omega = n_{max}/n_{min}$.

Table 5: Results of Imbalanced Learning with other methods on CIFAR10/100-LT. Each column of Method indicates the augmentation method, the type of classifier, and the type of loss function, respectively. Best in bold. (Accuracy (%): $mean_{std}$)

Method			Reference	CIFAR10-LT			CIFAR100-LT		
Augmentation	Classifier	Loss		100	50	10	100	50	10
-	FC	CE	(Yang et al., 2022b)	72.10 _{0.30}	77.60 _{0.30}	87.40 _{0.30}	-	-	-
-	ETF	CE	(Yang et al., 2022b)	72.90 _{0.30}	79.50 _{0.20}	87.20 _{0.10}	-	-	-
-	ETF	DR	(Yang et al., 2022b)	73.00 _{0.20}	78.40 _{0.30}	86.90 _{0.20}	-	-	-
mixup	FC	CE	(Yang et al., 2022b)	73.90 _{0.30}	79.30 _{0.20}	87.80 _{0.10}	43.00	48.10	-
B-mixup	FC	CE	(Zhang et al., 2022b)	78.70	-	89.60	-	-	-
mixup	ETF	CE	(Yang et al., 2022b)	67.00 _{0.40}	77.20 _{0.30}	87.00 _{0.20}	-	-	-
mixup	ETF	DR	(Yang et al., 2022b)	76.50 _{0.30}	81.00 _{0.20}	87.70 _{0.20}	45.30	50.40	-
-	FC	CE	(reproduced.) [†]	71.86 _{0.65}	77.58 _{0.48}	87.42 _{0.30}	41.65 _{0.41}	46.89 _{0.38}	60.48 _{0.34}
mixup	FC	CE	(reproduced.) [†]	74.24 _{0.44}	80.00 _{0.54}	89.08 _{0.32}	43.80 _{0.42}	49.57 _{0.37}	63.90 _{0.33}
arc-mixup	FNO	CE	(reproduced.) [†]	82.59 _{0.26}	85.13 _{0.25}	89.50 _{0.14}	49.26 _{2.82}	54.44 _{2.32}	63.14 _{3.82}

Table 6: Detailed Results of Imbalanced Learning on CIFAR10/100-LT. Best in bold. (Aug: the augmentation method, Clf: the type of classifier, \mathcal{L} : the type of loss function) (Accuracy (%): $mean_{std}$)

	Method			CIFAR10-LT				CIFAR100-LT			
	Aug	Clf	\mathcal{L}	Many	Median	Few	All	Many	Median	Few	All
imb 100	-	FC	CE [†]	91.99 _{3.02}	71.72 _{1.67}	51.92 _{5.34}	71.86 _{0.65}	68.41 _{0.55}	40.41 _{0.75}	9.94 _{0.55}	41.65 _{0.41}
	mixup	FC	CE [†]	94.72 _{0.43}	75.70 _{1.03}	51.82 _{1.61}	74.24 _{0.44}	73.25 _{0.52}	42.56 _{0.81}	8.73 _{0.49}	43.80 _{0.42}
	arc-mixup	FNO	CE [†]	84.60 _{0.78}	80.19 _{0.71}	83.79 _{0.86}	82.59 _{0.26}	63.46 _{3.55}	53.04 _{2.93}	27.05 _{2.08}	49.26 _{2.82}
imb 50	-	FC	CE [†]	93.30 _{1.26}	76.74 _{1.00}	62.98 _{1.91}	77.58 _{0.48}	69.95 _{0.61}	47.05 _{0.41}	18.08 _{0.63}	46.89 _{0.38}
	mixup	FC	CE [†]	94.98 _{0.30}	79.98 _{0.70}	65.04 _{1.58}	80.00 _{0.54}	74.53 _{0.42}	50.26 _{0.75}	17.75 _{0.66}	49.57 _{0.37}
	arc-mixup	FNO	CE [†]	86.00 _{0.69}	82.80 _{0.56}	87.38 _{0.81}	85.13 _{0.25}	64.38 _{2.97}	57.81 _{1.77}	38.01 _{2.49}	54.44 _{2.32}
imb 10	-	FC	CE [†]	94.00 _{0.40}	84.77 _{0.45}	84.37 _{0.64}	87.42 _{0.30}	72.53 _{0.37}	61.12 _{0.53}	44.73 _{0.76}	60.48 _{0.34}
	mixup	FC	CE [†]	95.37 _{0.29}	86.81 _{0.47}	85.80 _{0.89}	89.08 _{0.32}	76.89 _{0.43}	64.65 _{0.45}	46.88 _{0.68}	63.90 _{0.33}
	arc-mixup	FNO	CE [†]	89.94 _{0.55}	86.59 _{0.53}	92.95 _{0.38}	89.50 _{0.14}	66.26 _{5.32}	64.70 _{3.06}	57.39 _{3.07}	63.14 _{3.82}

Table 7: Detailed Results of Imbalanced Learning on Places-LT. All experiments were conducted using 3 trials with random seeds. Best in bold. (Aug: the augmentation method, Clf: the type of classifier, \mathcal{L} : the type of loss function) (Accuracy (%): $mean_{std}$)

Method			ResNet152				ResNet152 (FT)			
Aug	Clf	\mathcal{L}	Many	Median	Few	All	Many	Median	Few	All
-	FC	CE [†]	40.63 _{0.21}	17.69 _{0.64}	2.04 _{0.44}	22.62 _{0.29}	40.77 _{0.08}	19.99 _{0.67}	4.41 _{0.55}	24.16 _{0.41}
mixup	FC	CE [†]	42.10 _{0.31}	15.82 _{0.60}	0.86 _{0.18}	22.10 _{0.27}	43.14 _{0.63}	20.10 _{1.61}	3.32 _{1.33}	24.83 _{1.19}
arc-mixup	FNO	CE [†]	35.82 _{0.32}	31.74 _{0.64}	16.89 _{0.23}	30.07 _{0.40}	40.31 _{0.18}	35.31 _{0.19}	20.89 _{0.33}	34.06 _{0.10}

REFERENCES

- El Mehdi Achour, François Malgouyres, and Franck Mamalet. Existence, stability and scalability of orthogonal convolutional neural networks. *The Journal of Machine Learning Research*, 23(1): 15743–15798, 2022.
- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *Advances in neural information processing systems*, 32, 2019.
- Ari S Benjamin, David Rolnick, and Konrad Kording. Measuring and regularizing networks in function space. *arXiv preprint arXiv:1805.08289*, 2018.
- Cristian Bodnar, Francesco Di Giovanni, Benjamin Chamberlain, Pietro Liò, and Michael Bronstein. Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing in gnns. *Advances in Neural Information Processing Systems*, 35:18527–18541, 2022.

-
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- Zhiwen Cao, Zongcheng Chu, Dongfang Liu, and Yingjie Chen. A vector-based representation to enhance head pose estimation. In *Proceedings of the IEEE/CVF Winter Conference on applications of computer vision*, pp. 1188–1197, 2021.
- Jaehoon Cha and Jeyan Thiyaalingam. Orthogonality-enforced latent space in autoencoders: An approach to learning disentangled representations. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 3913–3948. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/cha23b.html>.
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- Arslan Chaudhry, Naeemullah Khan, Puneet Dokania, and Philip Torr. Continual learning in low-rank orthogonal subspaces. *Advances in Neural Information Processing Systems*, 33:9900–9911, 2020.
- Arslan Chaudhry, Albert Gordo, Puneet Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In *Proceedings of the AAAI conference on artificial intelligence*, pp. 6993–7001, 2021.
- Vitaliy Chiley, Ilya Sharapov, Atli Kossou, Urs Koster, Ryan Reece, Sofia Samaniego de la Fuente, Vishal Subbiah, and Michael James. Online normalization for training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Hien Dang, Tan Nguyen, Tho Tran, Hung Tran, and Nhat Ho. Neural collapse in deep linear network: From balanced to imbalanced data. *arXiv preprint arXiv:2301.00437*, 2023.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 3762–3773. PMLR, 2020.
- Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning. *arXiv preprint arXiv:2112.15121*, 2021.
- Jarrold Haas, William Yolland, and Bernhard T Rabus. Linking neural collapse and l2 normalization with improved out-of-distribution detection in deep neural networks. *Transactions on Machine Learning Research*, 2022.
- XY Han, Vardan Papyan, and David L Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path. *arXiv preprint arXiv:2106.02073*, 2021.
- Guiqing He, Yincheng Huo, Mingyao He, Haixi Zhang, and Jianping Fan. A novel orthogonality loss for deep hierarchical multi-task learning. *IEEE Access*, 8:67735–67744, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- Michael Hersche, Geethan Karunaratne, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas Rahimi. Constrained few-shot class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9057–9067, 2022.
- Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. *arXiv preprint arXiv:1810.12488*, 2018.
- Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. *arXiv preprint arXiv:2001.05992*, 2020.

-
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*, pp. 2261–2269. IEEE Computer Society, 2017. ISBN 978-1-5386-0457-1. URL <http://dblp.uni-trier.de/db/conf/cvpr/cvpr2017.html#HuangLMW17>.
- Huaibo Huang, Xiaoqiang Zhou, and Ran He. Orthogonal transformer: An efficient vision transformer backbone with token orthogonalization. *Advances in Neural Information Processing Systems*, 35: 14596–14607, 2022.
- Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. *arXiv preprint arXiv:2202.08384*, 2022.
- Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J. Su. An unconstrained layer-peeled perspective on neural collapse, 2022.
- Hoyong Kim, Minchan Kwon, and Kangil Kim. Revisiting softmax masking for stability in continual learning, 2023.
- Lingkai Kong, Yuqing Wang, and Molei Tao. Momentum stiefel optimizer, with applications to suitably-orthogonal attention, and optimal transport. *arXiv preprint arXiv:2205.14173*, 2022.
- Simon Kornblith, Ting Chen, Honglak Lee, and Mohammad Norouzi. Why do better loss functions lead to less transferable features? *Advances in Neural Information Processing Systems*, 34:28648–28662, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Lilly Kumari, Shengjie Wang, Tianyi Zhou, and Jeff A Bilmes. Retrospective adversarial replay for continual learning. *Advances in Neural Information Processing Systems*, 35:28530–28544, 2022.
- Shuai Li, Kui Jia, Yuxin Wen, Tongliang Liu, and Dacheng Tao. Orthogonal deep neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1352–1368, 2019.
- Zexi Li, Xinyi Shang, Rui He, Tao Lin, and Chao Wu. No fear of classifier biases: Neural collapse inspired federated learning with synthetic and fixed classifier. *arXiv preprint arXiv:2303.10058*, 2023.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Bingchen Liu, Yizhe Zhu, Zuohui Fu, Gerard De Melo, and Ahmed Elgammal. Oogan: Disentangling gan with one-hot sampling and orthogonal regularization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 4836–4843, 2020.
- Sheng Liu, Xiao Li, Yuexiang Zhai, Chong You, Zhihui Zhu, Carlos Fernandez-Granda, and Qing Qu. Convolutional normalization: Improving deep convolutional network robustness and training. *Advances in neural information processing systems*, 34:28919–28928, 2021.
- Weiyang Liu, Longhui Yu, Adrian Weller, and Bernhard Schölkopf. Generalizing and decoupling neural collapse via hyperspherical uniformity gap. *arXiv preprint arXiv:2303.06484*, 2023.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- Jianfeng Lu and Stefan Steinerberger. Neural collapse with cross-entropy loss. *arXiv preprint arXiv:2012.08465*, 2020.
- Yao Lu, Stephen Gould, and Thalaiyasingam Ajanthan. Bidirectionally self-normalizing neural networks. *Neural Networks*, 167:283–291, 2023.
- Yuexiao Ma, Taisong Jin, Xiawu Zheng, Yan Wang, Huixia Li, Yongjian Wu, Guannan Jiang, Wei Zhang, and Rongrong Ji. Ompq: Orthogonal mixed precision quantization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 9029–9037, 2023.

-
- Tharun Medini, Beidi Chen, and Anshumali Shrivastava. Solar: Sparse orthogonal learned and random embeddings. *arXiv preprint arXiv:2008.13225*, 2020.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Vardan Pappayan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40): 24652–24663, 2020. doi: 10.1073/pnas.2015509117. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2015509117>.
- Federico Pernici, Matteo Bruni, Claudio Baccchi, Francesco Turchini, and Alberto Del Bimbo. Class-incremental learning with pre-allocated fixed classifiers. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 6259–6266. IEEE, 2021.
- Vinay Venkatesh Ramasesh, Aitor Lewkowycz, and Ethan Dyer. Effect of scale on catastrophic forgetting in neural networks. In *International Conference on Learning Representations*, 2021.
- Kanchana Ranasinghe, Muzammal Naseer, Munawar Hayat, Salman Khan, and Fahad Shahbaz Khan. Orthogonal projection loss. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 12333–12343, 2021.
- Akshay Rangamani, Marius Lindegaard, Tomer Galanti, and Tomaso A Poggio. Feature learning in deep classifiers through intermediate neural collapse. In *International Conference on Machine Learning*, pp. 28729–28745. PMLR, 2023.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauero. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. *arXiv preprint arXiv:2103.09762*, 2021.
- Mhd Hasan Sarhan, Nassir Navab, Abouzar Eslami, and Shadi Albarqouni. Fairness by learning orthogonal disentangled representations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, pp. 746–761. Springer, 2020.
- Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International conference on machine learning*, pp. 4528–4537. PMLR, 2018.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. URL <http://arxiv.org/abs/1409.1556>.
- Christos Thrampoulidis, Ganesh Ramachandra Kini, Vala Vakilian, and Tina Behnia. Imbalance trouble: Revisiting neural-collapse geometry. *Advances in Neural Information Processing Systems*, 35:27225–27238, 2022.
- Louis C Tiao, Vincent Dutordoir, and Victor Picheny. Spherical inducing features for orthogonally-decoupled gaussian processes. *arXiv preprint arXiv:2304.14034*, 2023.

-
- Tom Tirer, Haoxiang Huang, and Jonathan Niles-Weed. Perturbation analysis of neural collapse. In *International Conference on Machine Learning*, pp. 34301–34329. PMLR, 2023.
- Asher Trockman and J Zico Kolter. Orthogonalizing convolutional layers with the cayley transform. *arXiv preprint arXiv:2104.07167*, 2021.
- Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X Yu. Orthogonal convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11505–11515, 2020.
- Stephan Wojtowytsch et al. On the emergence of simplex symmetry in the final and penultimate layers of neural network classifiers. *arXiv preprint arXiv:2012.05420*, 2020.
- Xia Wu, Xueyuan Xu, Jianhong Liu, Hailing Wang, Bin Hu, and Feiping Nie. Supervised feature selection with orthogonal regression and feature weighting. *IEEE Transactions on Neural Networks and Learning Systems*, 32(5):1831–1838, 2020.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. URL <http://arxiv.org/abs/1708.07747>. cite arxiv:1708.07747Comment: Dataset is freely available at <https://github.com/zalandoresearch/fashion-mnist> Benchmark is available at <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/>.
- Jing Xu and Haoxiong Liu. Quantifying the variability collapse of neural networks. *arXiv preprint arXiv:2306.03440*, 2023.
- Mengjia Xu, Akshay Rangamani, Qianli Liao, Tomer Galanti, and Tomaso Poggio. Dynamics in deep classifiers trained with the square loss: Normalization, low rank, neural collapse, and generalization bounds. *Research*, 6:0024, 2023.
- Xiaojun Xu, Linyi Li, and Bo Li. Lot: Layer-wise orthogonal training on improving l2 certified robustness. *Advances in Neural Information Processing Systems*, 35:18904–18915, 2022.
- Huanrui Yang, Minxue Tang, Wei Wen, Feng Yan, Daniel Hu, Ang Li, Hai Li, and Yiran Chen. Learning low-rank deep neural networks via singular vector orthogonality regularization and singular value sparsification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 678–679, 2020a.
- Liang Yang, Lina Kang, Qiuliang Zhang, Mengzhe Li, Dongxiao He, Zhen Wang, Chuan Wang, Xiaochun Cao, Yuanfang Guo, et al. Open: Orthogonal propagation with ego-network modeling. *Advances in Neural Information Processing Systems*, 35:9249–9261, 2022a.
- Yibo Yang, Liang Xie, Shixiang Chen, Xiangtai Li, Zhouchen Lin, and Dacheng Tao. Do we really need a learnable classifier at the end of deep neural network? *arXiv preprint arXiv:2203.09081*, 2022b.
- Yongyi Yang, Jacob Steinhardt, and Wei Hu. Are neurons actually collapsed? on the fine-grained structure in neural representations. 2023.
- Zebin Yang, Aijun Zhang, and Agus Sudjianto. Enhancing explainability of neural networks through architecture constraints. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6): 2610–2621, 2020b.
- Can Yaras, Peng Wang, Zhihui Zhu, Laura Balzano, and Qing Qu. Neural collapse with normalized features: A geometric analysis over the riemannian manifold. *Advances in neural information processing systems*, 35:11547–11560, 2022.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International conference on machine learning*, pp. 3987–3995. PMLR, 2017.

-
- Borui Zhang, Wenzhao Zheng, Jie Zhou, and Jiwen Lu. Bort: Towards explainable neural networks with bounded orthogonal constraint. *arXiv preprint arXiv:2212.09062*, 2022a.
- Shaoyu Zhang, Chen Chen, Xiujuan Zhang, and Silong Peng. Label-occurrence-balanced mixup for long-tailed recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3224–3228. IEEE, 2022b.
- Zhisheng Zhong, Jiequan Cui, Shu Liu, and Jiaya Jia. Improving calibration for long-tailed recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- Zhisheng Zhong, Jiequan Cui, Yibo Yang, Xiaoyang Wu, Xiaojuan Qi, Xiangyu Zhang, and Jiaya Jia. Understanding imbalanced semantic segmentation through neural collapse. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19550–19560, 2023.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.
- Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. BBN: Bilateral-branch network with cumulative learning for long-tailed visual recognition. pp. 1–8, 2020.
- Jinxin Zhou, Chong You, Xiao Li, Kangning Liu, Sheng Liu, Qing Qu, and Zhihui Zhu. Are all losses created equal: A neural collapse perspective. *Advances in Neural Information Processing Systems*, 35:31697–31710, 2022.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. 2021.