

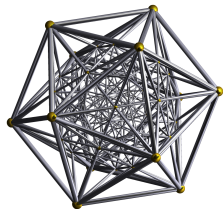
ACDL Summer Course 2023

## Lecture 2: Low-Dimensional Structures in Deep Representation Learning I

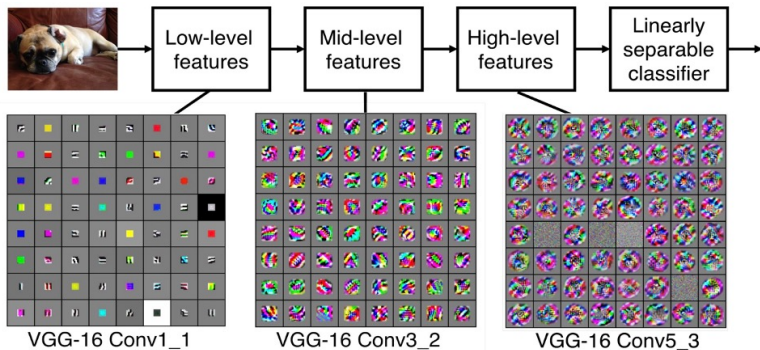
**Qing Qu**

EECS, University of Michigan

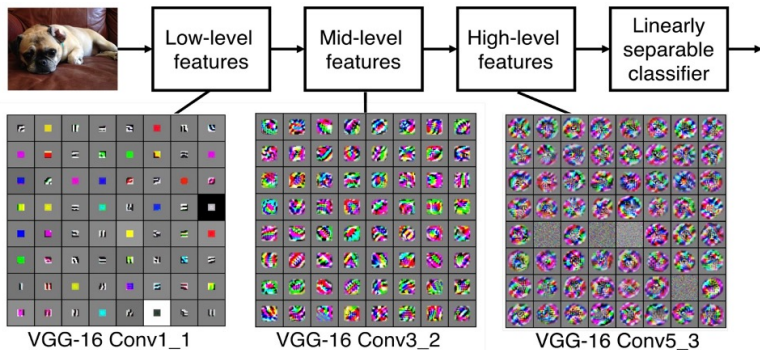
June 10th, 2023



# What Representations are DNNs Designed to Learn?

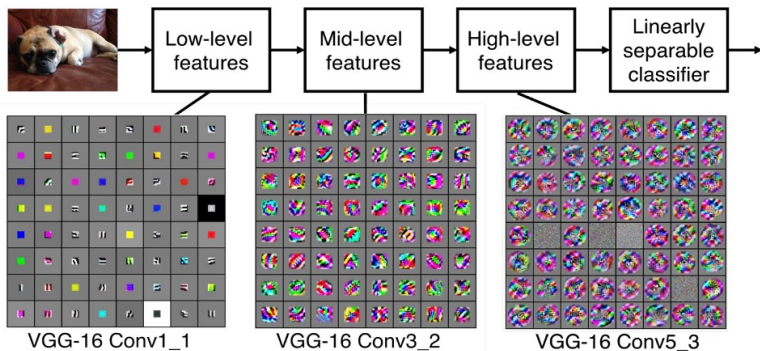


# What Representations are DNNs Designed to Learn?



- **Wishful Design:** DNNs learn rich representations across different layers.

# What Representations are DNNs Designed to Learn?



- **Wishful Design:** DNNs learn rich representations across different layers.
- **Reality:** Is it really the case in the practice of modern DNNs?

# Outline

- 1 Low-Dimensional Representation: Neural Collapse (NC)
- 2 Understanding NC from Optimization
- 3 Prevalence of NC under Different Training Scenarios
- 4 Conclusion

# Multi-Class Image Classification Problem

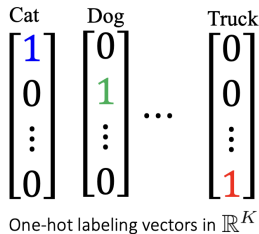
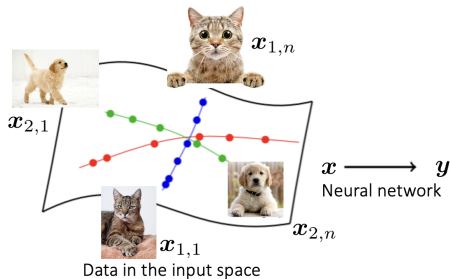
- **Goal:** Learn a deep network predictor from a labelled training dataset  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, n\}$ .

---

<sup>1</sup>If not, we can use data augmentation to make them balanced

# Multi-Class Image Classification Problem

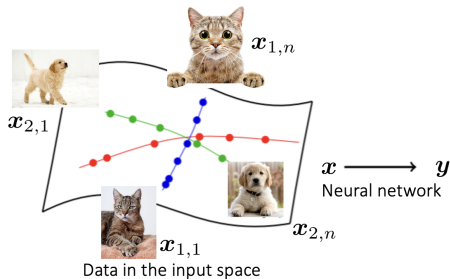
- **Goal:** Learn a deep network predictor from a labelled training dataset  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, n\}$ .
- **Training Labels:**  $k = 1, \dots, K$ 
  - $K = 10$  classes (MNIST, CIFAR10, etc)
  - $K = 1000$  classes (ImageNet)



<sup>1</sup>If not, we can use data augmentation to make them balanced

# Multi-Class Image Classification Problem

- **Goal:** Learn a deep network predictor from a labelled training dataset  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}); i = 1, \dots, n\}$ .
- **Training Labels:**  $k = 1, \dots, K$ 
  - $K = 10$  classes (MNIST, CIFAR10, etc)
  - $K = 1000$  classes (ImageNet)



$$\begin{array}{c} \text{Cat} \\ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{array}
 \quad
 \begin{array}{c} \text{Dog} \\ \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \end{array}
 \quad
 \dots
 \quad
 \begin{array}{c} \text{Truck} \\ \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \end{array}$$

One-hot labeling vectors in  $\mathbb{R}^K$

- For simplicity, we assume **balanced** dataset where each class has  $n$  training samples.<sup>1</sup>

<sup>1</sup>If not, we can use data augmentation to make them balanced



# Deep Neural Network Classifiers

- **A vanilla deep network:**

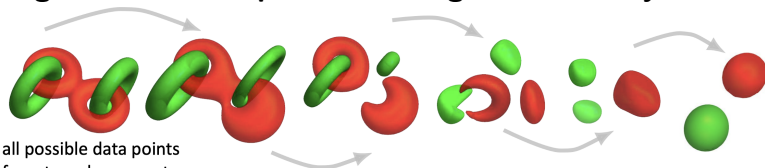
$$f_{\Theta}(\mathbf{x}) = \underbrace{\mathbf{W}_L}_{\text{linear classifier } \mathbf{W}} \underbrace{\sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_{L-1})}_{\text{feature } \phi_{\theta}(\mathbf{x})=:h} + \mathbf{b}_L$$

# Deep Neural Network Classifiers

- **A vanilla deep network:**

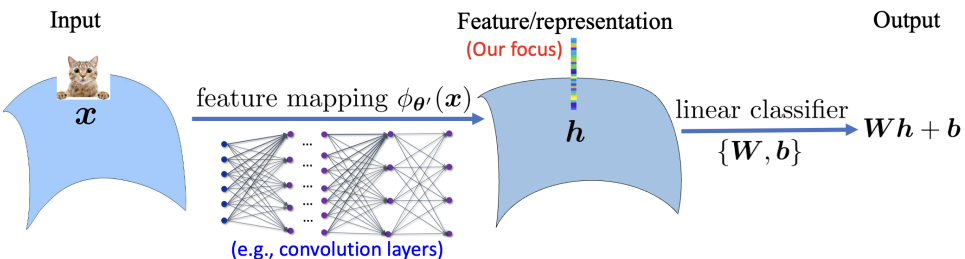
$$f_{\Theta}(x) = \underbrace{W_L}_{\text{linear classifier } W} \underbrace{\sigma(W_{L-1} \cdots \sigma(W_1 x + b_1) + b_{L-1})}_{\text{feature } \phi_{\theta}(x) =: h} + b_L$$

- **Progressive linear separation through nonlinear layers:**



all possible data points  
from two classes; not a  
single input!

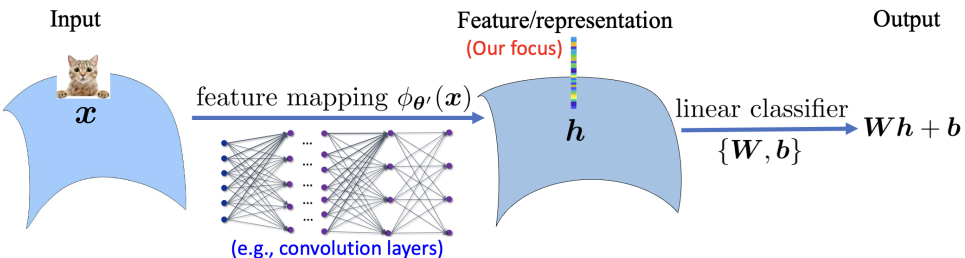
# Deep Neural Network Classifiers



- Training a deep neural network:

$$\min_{\theta, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \underbrace{\mathcal{L}_{\text{CE}}(\mathbf{W}\phi_{\theta}(\mathbf{x}_{k,i}) + \mathbf{b}, \mathbf{y}_k)}_{\text{cross-entropy (CE) loss}} + \lambda \underbrace{\|(\theta, \mathbf{W}, \mathbf{b})\|_F^2}_{\text{weight decay}}$$

# Deep Neural Network Classifiers



Output:  $f(\mathbf{x}; \theta) = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \xrightarrow{\text{Softmax function}} \begin{bmatrix} 0.6 \\ 0.3 \\ 0.1 \end{bmatrix}$

Cat  
Dog  
Panda

$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

CE(Cat):  $= -q(\text{Cat}) \cdot \log p(\text{Cat})$   
 $= -1 \cdot \log 0.6$   
 $= 0.51\dots$

Prediction (probability)      Target

# Neural Collapse in Multi-Class Classification

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

 Vardan Papayan,  X. Y. Han, and David L. Donoho

[+ See all authors and affiliations](#)

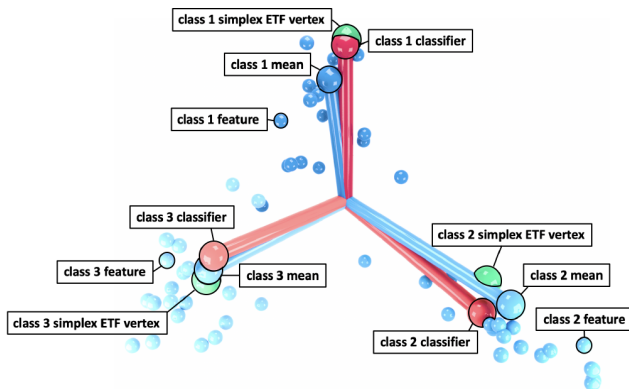
PNAS October 6, 2020 117 (40) 24652-24663; first published September 21, 2020;

<https://doi.org/10.1073/pnas.2015509117>

Contributed by David L. Donoho, August 18, 2020 (sent for review July 22, 2020; reviewed by Helmut Boelsckei and Stéphane Mallat)

- Reveals common outcome of learned features and classifiers across a variety of architectures and dataset
- Precise mathematical structure within the features and classifier

# Neural Collapse in Multi-Class Classification



**Credit:** Han et al. Neural Collapse Under MSE Loss: Proximity to and Dynamics on the Central Path. ICLR, 2022.

# Neural Collapse: Symmetry and Structures

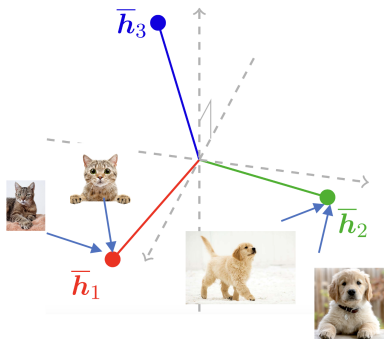
- **NC1: Within-Class Variability Collapse:** features of each class collapse to class-mean with **zero** variability:

$$k\text{-th class, } i\text{-th sample} : \mathbf{h}_{k,i} \rightarrow \bar{\mathbf{h}}_k,$$

# Neural Collapse: Symmetry and Structures

- NC1: Within-Class Variability Collapse:** features of each class collapse to class-mean with **zero** variability:

$k$ -th class,  $i$ -th sample :  $\mathbf{h}_{k,i} \rightarrow \bar{\mathbf{h}}_k$ ,

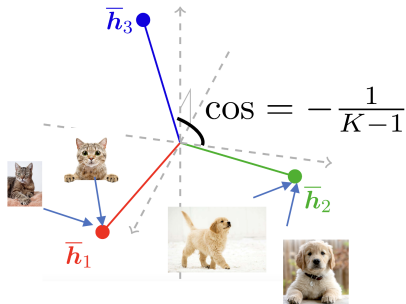




# Neural Collapse: Symmetry and Structures

- NC2: Convergence to Simplex Equiangular Tight Frame (ETF):** the class means are linearly separable, and maximally distant

$$\frac{\langle \bar{\mathbf{h}}_k, \bar{\mathbf{h}}_{k'} \rangle}{\|\bar{\mathbf{h}}_k\| \|\bar{\mathbf{h}}_{k'}\|} \rightarrow \begin{cases} 1, & k = k' \\ -\frac{1}{K-1}, & k \neq k' \end{cases}$$

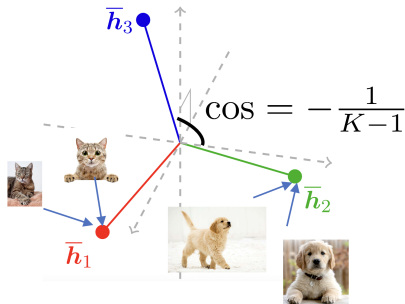


# Neural Collapse: Symmetry and Structures

- NC2: Convergence to Simplex Equiangular Tight Frame (ETF):** the class means are linearly separable, and maximally distant

$$\overline{\mathbf{H}}^\top \overline{\mathbf{H}} \sim \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top,$$

$$\overline{\mathbf{H}} = [\overline{\mathbf{h}}_1 \quad \cdots \quad \overline{\mathbf{h}}_K]$$

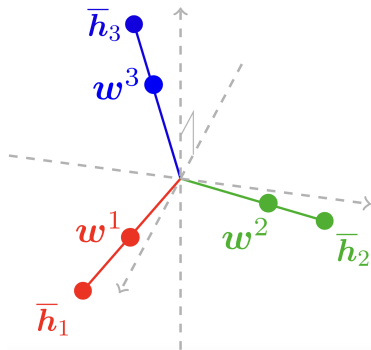


# Neural Collapse: Symmetry and Structures

- **NC3: Convergence to Self-Duality**: the last-layer classifiers are perfectly matched with the class-means of features

$$\frac{w^k}{\|w^k\|} \rightarrow \frac{\bar{h}_k}{\|\bar{h}_k\|},$$

where  $w^k$  represents the  $k$ -th row of  $W$ .



# Understanding the Prevalence of Neural Collapse

**Question.** Given the prevalence of Neural Collapse across datasets and network architectures, why would such a phenomenon happen in training overparameterized networks?

# Outline

- ① Low-Dimensional Representation: Neural Collapse (NC)
- ② Understanding NC from Optimization
- ③ Prevalence of NC under Different Training Scenarios
- ④ Conclusion

# Dealing with a Highly Nonconvex Problem

The training problem is highly **nonconvex** [Li et al.'18]:

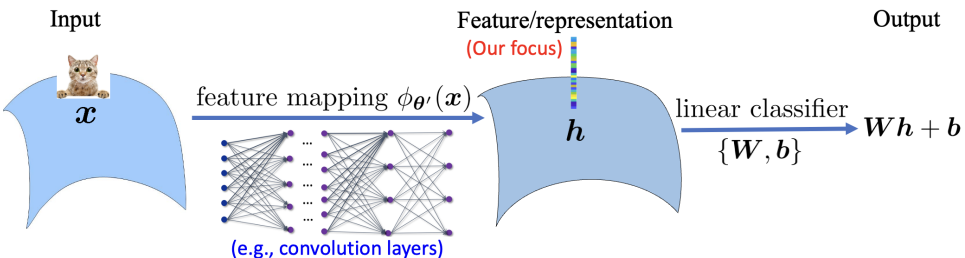
$$\min_{\theta', \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \phi_{\theta'}(\mathbf{x}_{k,i}) + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\theta', \mathbf{W}, \mathbf{b})\|_F^2,$$

due to the fact that the network

$$f_{\Theta}(\mathbf{x}) = \underbrace{\mathbf{W}_L}_{\text{linear classifier } W} \underbrace{\sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_{L-1})}_{\text{feature } \phi_{\Theta}(\mathbf{x})=:h} + \mathbf{b}_L$$

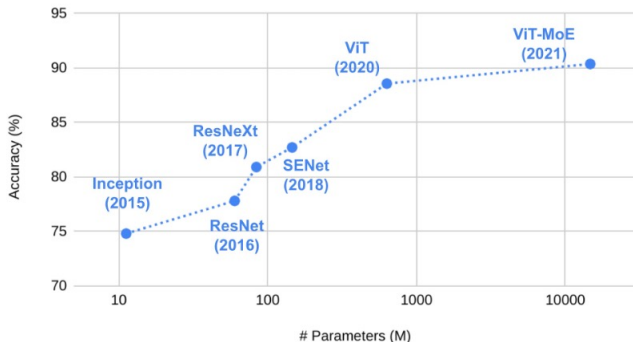
- **Nonlinear interaction across layers.**
- **Nonlinear activation functions.**

## Simplification: Unconstrained Feature Model



**Assumption.** We treat  $\mathbf{H} = [\mathbf{h}_{1,1} \ \cdots \ \mathbf{h}_{K,n}]$  as a **free** optimization variable, ignoring the constraint  $\mathbf{h}\phi_{\theta}(\mathbf{x})$ .

# The Trend of Large Models...



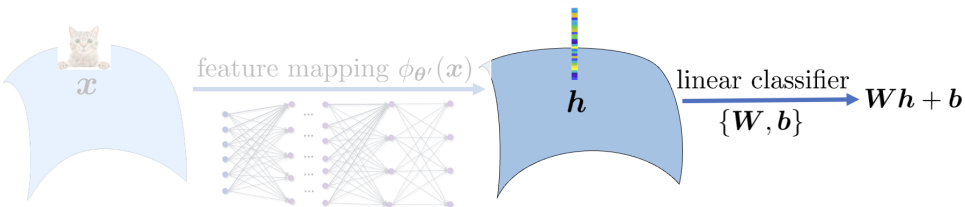
**Figure:** Accuracy vs. model size for image classification on ImageNet dataset

~23 million (# Parameters in ResNet-50)  $\gg$  ~1 million (# Samples in ImageNet)

**In principle, deep network can fit *any* training labels!**  
*(i.e., not only clean, but also corrupted labels)*

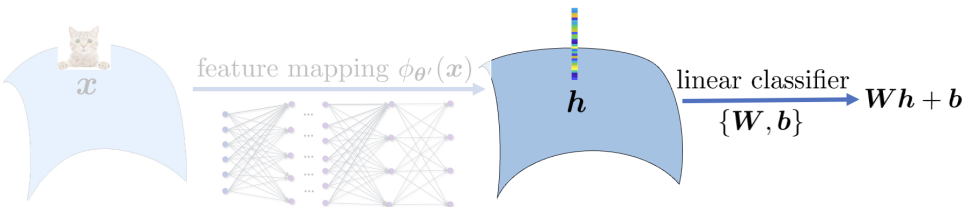


## Simplification: Unconstrained Feature Model



**Assumption.** We treat  $\mathbf{H} = [\mathbf{h}_{1,1} \ \cdots \ \mathbf{h}_{K,n}]$  as a **free** optimization variable, ignoring the constraint  $\mathbf{h}\phi_{\theta}(\mathbf{x})$ .

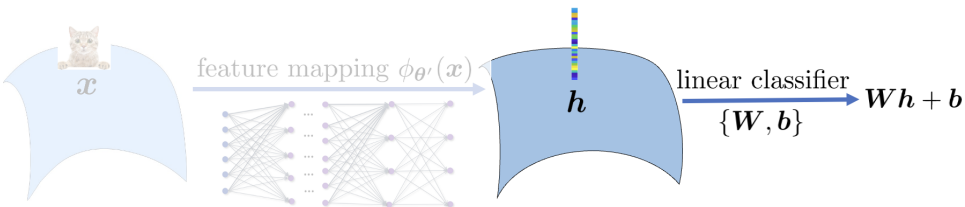
## Simplification: Unconstrained Feature Model



**Assumption.** We treat  $\mathbf{H} = [\mathbf{h}_{1,1} \ \cdots \ \mathbf{h}_{K,n}]$  as a **free** optimization variable, ignoring the constraint  $\mathbf{h}\phi_{\theta}(\mathbf{x})$ .

- **Validity:** modern network are highly overparameterized, that they are **universal approximators** [Shaham'18];

## Simplification: Unconstrained Feature Model

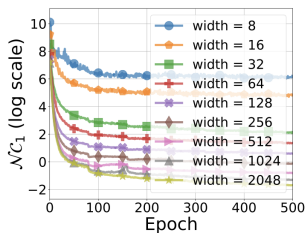


**Assumption.** We treat  $\mathbf{H} = [\mathbf{h}_{1,1} \ \cdots \ \mathbf{h}_{K,n}]$  as a **free** optimization variable, ignoring the constraint  $\mathbf{h}\phi_{\theta}(\mathbf{x})$ .

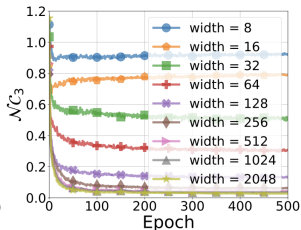
- **Validity:** modern networks are highly overparameterized, that they are **universal approximators** [Shaham'18];
- **State-of-the-Art:** also called **Layer-Peeled Model** [Fang'21], existing work [E'20, Lu'20, Mixon'20, Fang'21] only studied global optimality conditions;

# Experiments: NC Occurs on Random Labels/Inputs

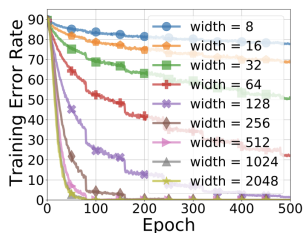
CIFAR-10 with **random** labels, MLP with **varying network widths**



Within-Class Variability (NC1)



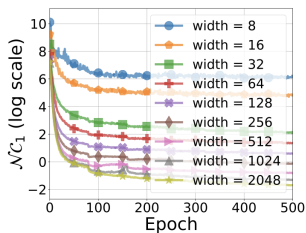
Self-Duality Collapse (NC2)



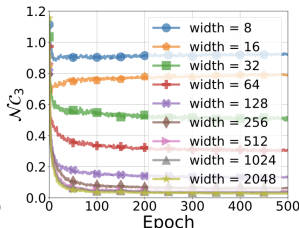
Training Error

# Experiments: NC Occurs on Random Labels/Inputs

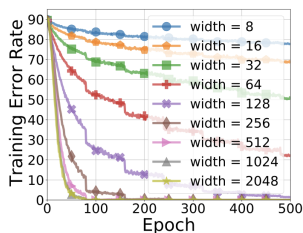
CIFAR-10 with **random** labels, MLP with **varying network widths**



Within-Class Variability (NC1)



Self-Duality Collapse (NC2)



Training Error

- **Validity of unconstrained features model:** Learn NC last-layer features and classifiers for any inputs
- The network memorizes training data in a very special way: NC
- We observe similar results on **random inputs (random pixels)**

# Geometric Analysis of Global Landscape

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

Theorem (Global Optimality & Benign Global Landscape, Zhu et al.'21)

Let feature dimension  $d$  is larger than the class number  $K$ , i.e.,  $d > K$ . Consider the above nonconvex optimization problem w.r.t.  $(\mathbf{W}, \mathbf{H})$ . Then

- **Global optimality:** Any global solution  $(\{\mathbf{H}^*, \mathbf{W}^*, \mathbf{b}^*\})$  obeys Neural Collapse, with  $\mathbf{b}^* = 0$  and

$$\underbrace{\mathbf{h}_{k,i}^* = \bar{\mathbf{h}}_k^*}_{\text{NC1}}, \quad \underbrace{\frac{\langle \bar{\mathbf{h}}_k^*, \bar{\mathbf{h}}_{k'}^* \rangle}{\|\bar{\mathbf{h}}_k^*\| \|\bar{\mathbf{h}}_{k'}^*\|} = \begin{cases} 1, & k = k' \\ -\frac{1}{K-1}, & k \neq k' \end{cases}}_{\text{NC2}}, \quad \underbrace{\frac{\mathbf{w}^{k*}}{\|\mathbf{w}^{k*}\|} = \frac{\bar{\mathbf{h}}_k^*}{\|\bar{\mathbf{h}}_k^*\|}}_{\text{NC3}}$$

# Geometric Analysis of Global Landscape

[Lu et al.'20] study the following one-example-per class model

$$\min_{\{\mathbf{h}_k\}} \frac{1}{K} \sum_{k=1}^K \mathcal{L}_{\text{CE}}(\mathbf{h}_k, \mathbf{y}_k), \text{ s.t. } \|\mathbf{h}_k\|_2 = 1$$

[E et al.'20, Fang et al.'21, Gral et al.'21, etc.] study constrained formulation

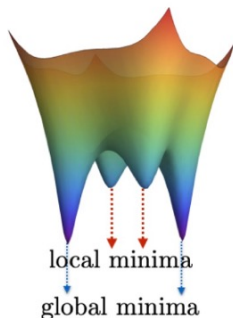
$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k), \text{ s.t. } \|\mathbf{W}\|_F \leq 1, \|\mathbf{h}_{k,i}\|_2 \leq 1$$

These work show that any global solution has NC, but

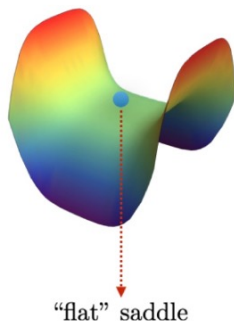
- What about **local minima/saddle points**?
- The constrained formulations are not aligned with practice

## Global Optimality Does Not Imply Efficient Optimization

“bad” local minima



“flat” saddle point



Our loss is still highly nonconvex:

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

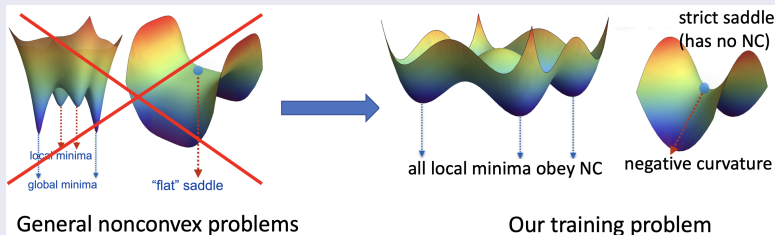


# Geometric Analysis of Global Landscape

## Theorem (Global Optimality & Benign Global Landscape, Zhu et al.'21)

Let feature dimension  $d$  is larger than the class number  $K$ , i.e.,  $d > K$ . Consider the above nonconvex optimization problem w.r.t.  $(\mathbf{W}, \mathbf{H})$ . Then

- **Global optimality:** Any global solution  $(\{\mathbf{H}^*, \mathbf{W}^*, \mathbf{b}^*\})$  obeys Neural Collapse.
- **Benign global landscape:** The objective function (i) has no spurious local minima, and (ii) any non-global critical point is a strict saddle with negative curvature.



# Geometric Analysis of Global Landscape

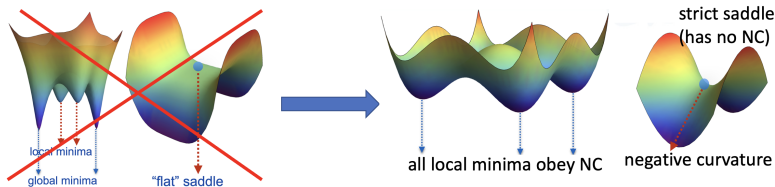
## Theorem (Global Optimality & Benign Global Landscape, Zhu et al.'21)

*Let feature dimension  $d$  is larger than the class number  $K$ , i.e.,  $d > K$ . Consider the above nonconvex optimization problem w.r.t.  $(\mathbf{W}, \mathbf{H})$ . Then*

- **Global optimality:** *Any global solution  $(\{\mathbf{H}^*, \mathbf{W}^*, \mathbf{b}^*\})$  obeys Neural Collapse.*
- **Benign global landscape:** *The objective function (i) has no spurious local minima, and (ii) any non-global critical point is a strict saddle with negative curvature.*

**Message.** Iterative algorithms such as (stochastic) gradient descent will always learn Neural Collapse features and classifiers.

# Implications of Our Results

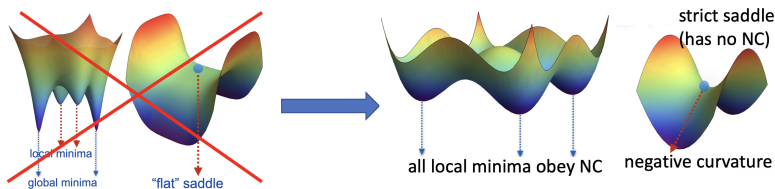


General nonconvex problems

Our training problem

- **A feature learning perspective.**
  - **Top down:** unconstrained feature model, representation learning, but no input information.
  - **Bottom up:** shallow network, strong assumptions, far from practice.

# Implications of Our Results



General nonconvex problems

Our training problem

- **A feature learning perspective.**
  - **Top down:** unconstrained feature model, representation learning, but no input information.
  - **Bottom up:** shallow network, strong assumptions, far from practice.
- **Connections to empirical phenomena.**

# Implications of Our Results

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

variational form:  $\|\mathbf{Z}\|_* = \min_{\mathbf{Z}=\mathbf{W}\mathbf{H}} \frac{1}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2)$

- Closely relates to **low-rank matrix factorization** problems [Burer et al'03, Bhojanapalli et al'16, Ge et al'16, Zhu et al'18, Li et al'19, Chi et al'19]

# Implications of Our Results

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

$$\text{variational form: } \|\mathbf{Z}\|_* = \min_{\mathbf{Z}=\mathbf{WH}} \frac{1}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2)$$

- Closely relates to **low-rank matrix factorization** problems [Burer et al'03, Bhojanapalli et al'16, Ge et al'16, Zhu et al'18, Li et al'19, Chi et al'19]
- However, we have more **structured** observation

$$\mathbf{Y} = \begin{bmatrix} 1 & \cdots & 1 & & & & \\ & & & 1 & \cdots & 1 & \\ & & & & & & 1 & \cdots & 1 \end{bmatrix} = \mathbf{I}_K \otimes \mathbf{1}_n^\top$$

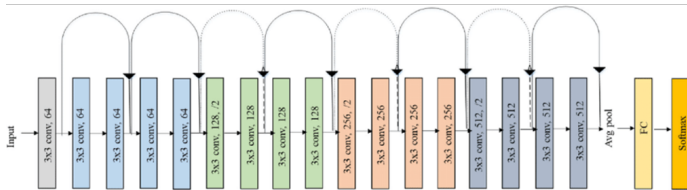
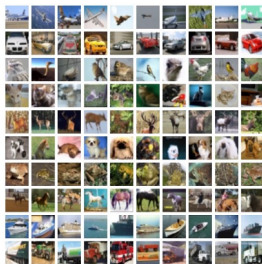
# Experiments on Practical Neural Networks

Conduct experiments with **practical networks** to verify our findings:

Use a Residual Neural Network (ResNet) on CIFAR-10 Dataset:

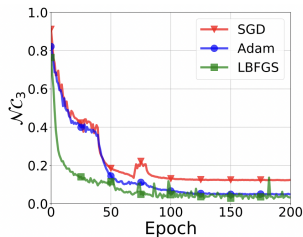
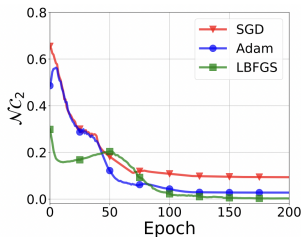
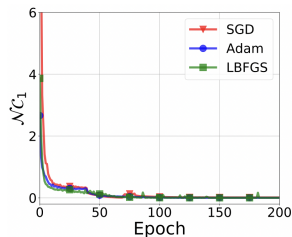
- $K = 10$  classes
- 50K training images
- 10K testing images

airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck



# Experiments: NC is Algorithm Independent

ResNet18 on CIFAR-10 with **different training algorithms**



Within-Class Variability ( $\mathcal{NC}_1$ )

Between-Class Separation ( $\mathcal{NC}_2$ )

Self-Duality Collapse ( $\mathcal{NC}_3$ )

- The smaller the quantities, the severer NC
- NC is prevalent across **different training algorithms**



# Exploit NC for Improving Training & Memory

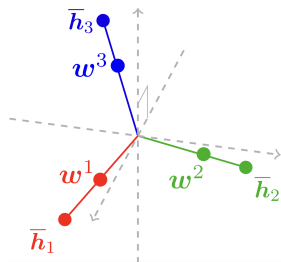
NC is prevalent, and classifier always converges to a Simplex ETF

- **Implication 1: No need to learn the classifier** [Hoffer et al. 2018]
  - Just fix it as a Simplex ETF
  - Save **8%, 12%, and 53%** parameters for ResNet50, DenseNet169, and ShuffleNet!

# Exploit NC for Improving Training & Memory

NC is prevalent, and classifier always converges to a Simplex ETF

- **Implication 1: No need to learn the classifier** [Hoffer et al. 2018]
  - Just fix it as a Simplex ETF
  - Save **8%, 12%, and 53%** parameters for ResNet50, DenseNet169, and ShuffleNet!
- **Implication 2: No need of large feature dimension  $d$** 
  - Just use feature dim.  $d = \# \text{class } K$  (e.g.,  $d = 10$  for CIFAR-10)
  - Further saves **21% and 4.5%** parameters for ResNet18 and ResNet50!



# Exploit NC for Improving Training & Memory

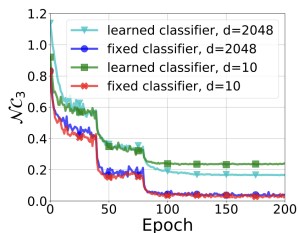
ResNet50 on CIFAR-10 with different settings

- **Learned** classifier (default) vs. **fixed** classifier as a simplex ETF
- Feature dim  $d = 2048$  (default) vs.  $d = 10$

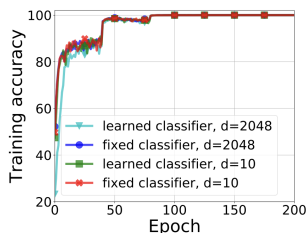
# Exploit NC for Improving Training & Memory

ResNet50 on CIFAR-10 with different settings

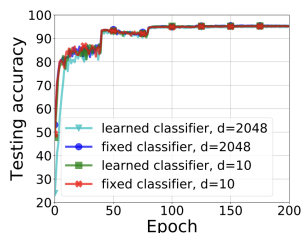
- **Learned** classifier (default) vs. **fixed** classifier as a simplex ETF
- Feature dim  $d = 2048$  (default) vs.  $d = 10$



Self-Duality Collapse (NC3)



Training Accuracy

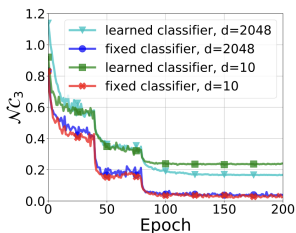


Testing Accuracy

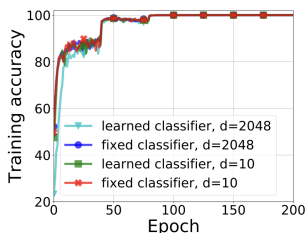
# Exploit NC for Improving Training & Memory

ResNet50 on CIFAR-10 with different settings

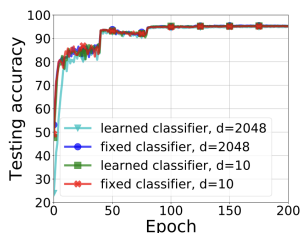
- **Learned** classifier (default) vs. **fixed** classifier as a simplex ETF
- Feature dim  $d = 2048$  (default) vs.  $d = 10$



Self-Duality Collapse (NC3)



Training Accuracy



Testing Accuracy

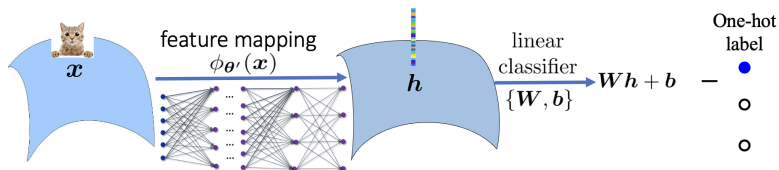
- Training with **small** dimensional features and **fixed** classifiers achieves on-par performance with **large** dimensional features and **learned** classifiers.

# Outline

- ① Low-Dimensional Representation: Neural Collapse (NC)
- ② Understanding NC from Optimization
- ③ Prevalence of NC under Different Training Scenarios
- ④ Conclusion

# Is Cross-entropy Loss Essential?

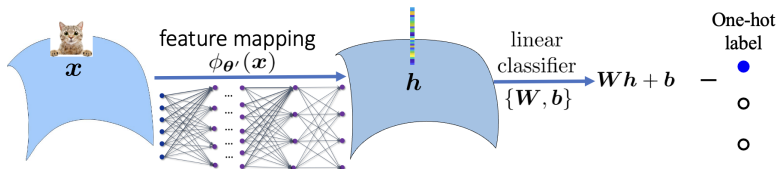
**Question.** Is cross-entropy loss essential to neural collapse?



<sup>2</sup>He et al., Bag of tricks for image classification with convolutional neural networks, CVPR'19.

# Is Cross-entropy Loss Essential?

**Question.** Is cross-entropy loss essential to neural collapse?



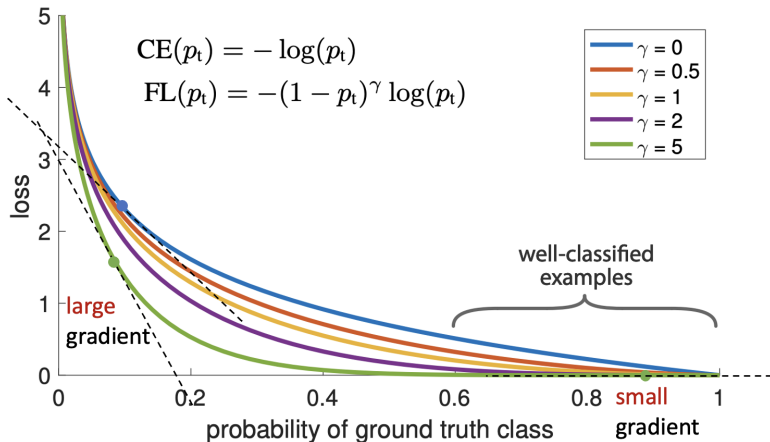
- We can measure the mismatch between the network output and the one-hot label in many ways.
- Various losses and tricks (e.g., label smoothing, focal loss) have been proposed to improve network training and performance<sup>2</sup>

<sup>2</sup>He et al., Bag of tricks for image classification with convolutional neural networks, CVPR'19.



## Example I: Focal Loss (FL)

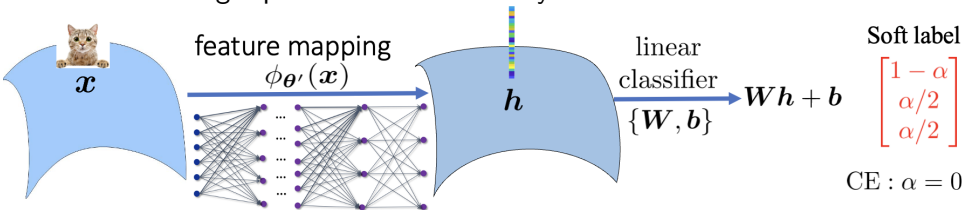
Focal loss puts more focus on hard, misclassified examples<sup>3</sup>



<sup>3</sup>Lin et al., Focal Loss for Dense Object Detection, CVPR'18.

## Example II: Label Smoothing (LS)

Label smoothing replaces the hard label by a soft label<sup>4</sup>

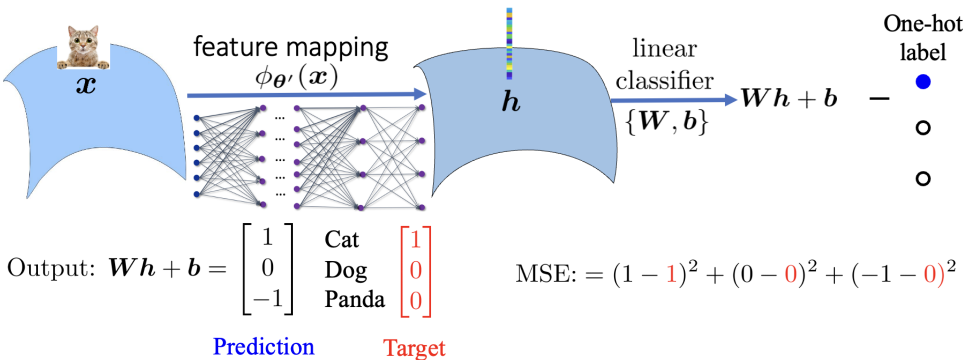


Output:  $Wh + b = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \xrightarrow{\text{Softmax function}} \begin{bmatrix} 0.6 \\ 0.3 \\ 0.1 \end{bmatrix}$

		Cat	$\begin{bmatrix} 1 - \alpha \\ \alpha/2 \\ \alpha/2 \end{bmatrix}$	LS = $-q(\text{Cat}) \cdot \log p(\text{Cat})$
		Dog		$-q(\text{Dog}) \cdot \log p(\text{Dog})$
		Panda		$-q(\text{Panda}) \cdot \log p(\text{Panda})$
	<b>Prediction</b>	<b>Target</b>		$= -(1 - \alpha) \log(0.6)$
				$- \frac{\alpha}{2} \log(0.3)$
				$- \frac{\alpha}{2} \log(0.1)$

<sup>4</sup>Szegedy et al., Rethinking the inception architecture for computer vision, CVPR'16.  
Muller, Kornblith, Hinton, When does label smoothing help?, NeurIPS'19.

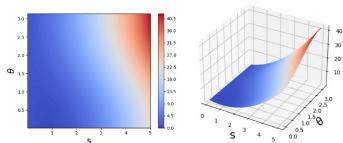
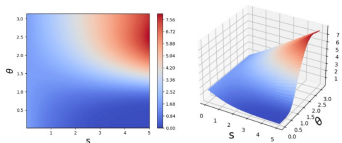
# Example III: Mean-squared Error (MSE) Loss



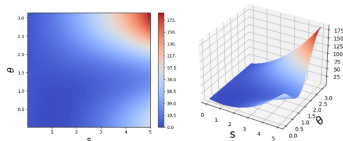
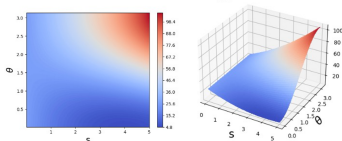
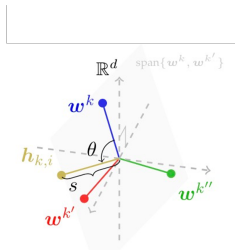
Compared with CE, **rescaled** MSE loss produces on par results for computer vision & NLP tasks.<sup>5</sup>

<sup>5</sup>Hui & Belkin, Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks, ICLR 2021.

## Example III: Mean-squared Error (MSE) Loss

(a) Vanilla MSE ( $\alpha = 1, M = 1$ )

(b) Cross Entropy

(c) Rescaled MSE ( $\alpha = 5, M = 1$ )(d) Rescaled MSE ( $\alpha = 1, M = 5$ )

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{2N} \|\Omega_{\alpha}^{\odot 1/2} \odot (\mathbf{W}\mathbf{H} + \mathbf{b}\mathbf{1}^{\top} - \mathbf{M}\mathbf{Y})\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2.$$

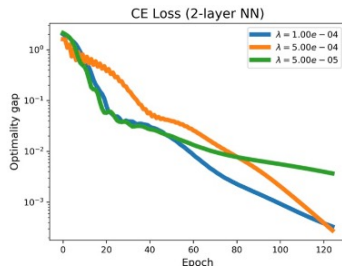
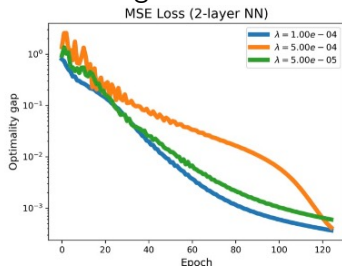
## Example III: Mean-squared Error (MSE) Loss

- Error bound condition for vanilla MSE loss:

$$\text{dist}((\mathbf{W}, \mathbf{H}, \mathbf{b}), \mathcal{X}) \leq \kappa \|\nabla F(\mathbf{W}, \mathbf{H}, \mathbf{b})\|_F$$

for any  $(\mathbf{W}, \mathbf{H}, \mathbf{b})$  with  $\text{dist}((\mathbf{W}, \mathbf{H}, \mathbf{b}), \mathcal{X}) \leq \delta$ .

- Local linear convergence of GD:



## Which Loss is the Best to Use?

Testing accuracy (%) for WideResNet18 on mini-ImageNet with different widths and training iterations

Loss	CE	FL	LS	MSE
Width = $\times 0.25$ Epochs = 200	71.95	70.20	70.40	69.15
Width = $\times 2$ Epochs = 800	79.30	79.32	80.20	79.62

- The performance is also affected by the choice of network architecture, training iterations, dataset, etc.

# Are All Losses Created Equal?—A NC Perspective

## Theorem (Informal, Zhou et al.'22)

*Under the unconstrained feature model, with feature dim.*

*$d \geq \#class K - 1$ , for all the one-hot labeling based losses (e.g., CE, FL, LS, MSE),*

- NC are the only global solutions for all losses.*
- All losses have benign global landscape w.r.t.  $(\mathbf{W}, \mathbf{H}, \mathbf{b})$*

# Are All Losses Created Equal?—A NC Perspective

## Theorem (Informal, Zhou et al.'22)

*Under the unconstrained feature model, with feature dim.*

*$d \geq \#class K - 1$ , for all the one-hot labeling based losses (e.g., CE, FL, LS, MSE),*

- *NC are the only global solutions for all losses.*
- *All losses have benign global landscape w.r.t.  $(\mathbf{W}, \mathbf{H}, \mathbf{b})$*

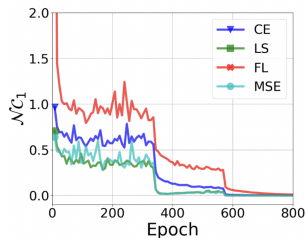
**Implication for practical networks** If network is *large enough and trained longer enough*

- All losses lead to largely identical features on **training data**—NC phenomena
- All losses lead to largely identical performance on **test data** (experiments in the following slides)

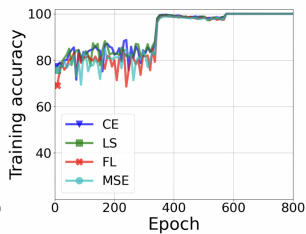


# Are All Losses Created Equal?—A NC Perspective

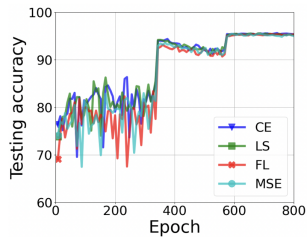
ResNet50 (with different network widths and training epochs) on CIFAR-10 with **different training losses**



Within-Class Variability (NC1)



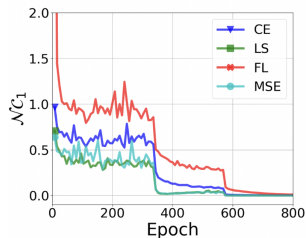
Train accuracy



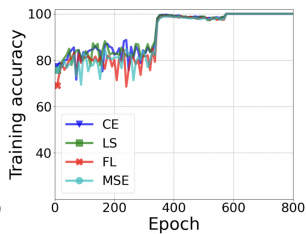
Testing accuracy

# Are All Losses Created Equal?—A NC Perspective

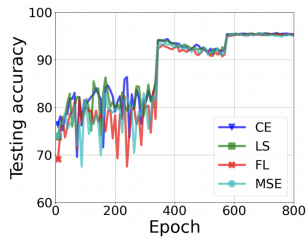
ResNet50 (with different network widths and training epochs) on CIFAR-10 with **different training losses**



Within-Class Variability (NC1)



Train accuracy

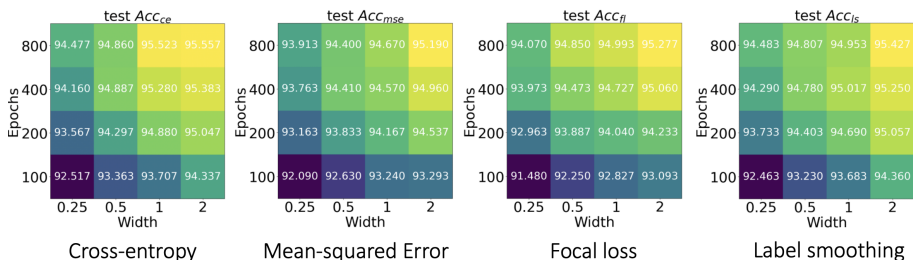


Testing accuracy

**Observation:** If network is *large enough and trained longer enough*, all losses lead to largely identical NC features on **training data**.

# All Losses Are Almost Created Equal

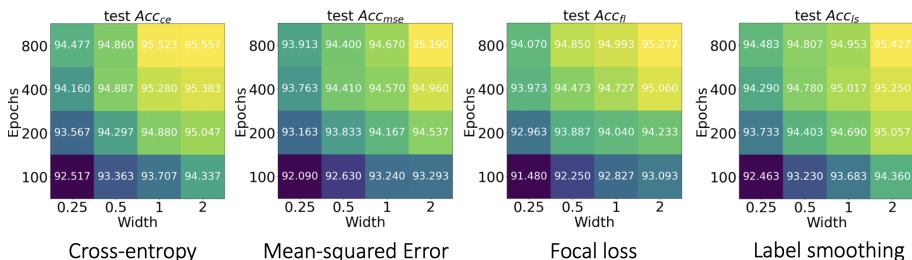
ResNet50 (with different network widths and training epochs) on CIFAR-10 with **different training losses**



- Right top corners not only have better performance, but also have **smaller** variance than left bottom corners

# All Losses Are Almost Created Equal

ResNet50 (with different network widths and training epochs) on CIFAR-10 with **different training losses**



- Right top corners not only have better performance, but also have **smaller** variance than left bottom corners

**Observation:** If network is *large enough and trained longer enough*, all losses lead to largely identical performance on **test data**.

# Neural Collapse with Feature Normalization

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k)$$

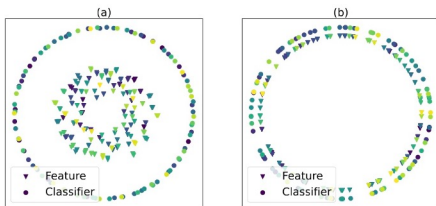
$$\text{s.t. } \|\mathbf{w}_k\|_2 = \tau, \|\mathbf{h}_{k,i}\|_2 = 1, \mathbf{h}_{k,i} = \phi_{\boldsymbol{\theta}}(\mathbf{x}_{k,i}), \quad \forall i \in [n], \forall k \in [K].$$

# Neural Collapse with Feature Normalization

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t. } \|\mathbf{w}_k\|_2 = \tau, \|\mathbf{h}_{k,i}\|_2 = 1, \mathbf{h}_{k,i} = \phi_{\theta}(\mathbf{x}_{k,i}), \forall i \in [n], \forall k \in [K].$$

- Improve the quality of learned features with larger class separation [Yu et al., 2020, Wang and Isola, 2020]
- Improve test performance in practice [Graf et al., 2021, Liu et al., 2021]



# Neural Collapse with Feature Normalization

- Under the unconstrained feature model, a similar global landscape result can be shown for:

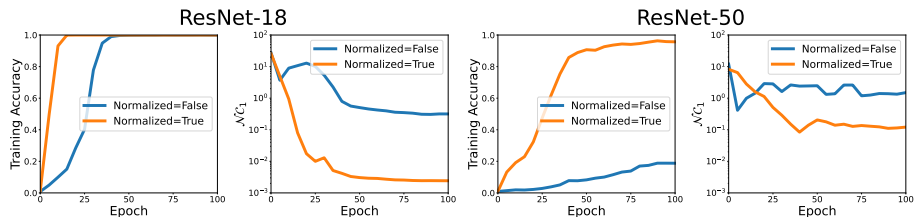
$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t. } \|\mathbf{w}_k\|_2 = \tau, \|\mathbf{h}_{k,i}\|_2 = 1, \forall i \in [n], \forall k \in [K].$$

- More advanced analysis based upon Riemannian optimization tools.

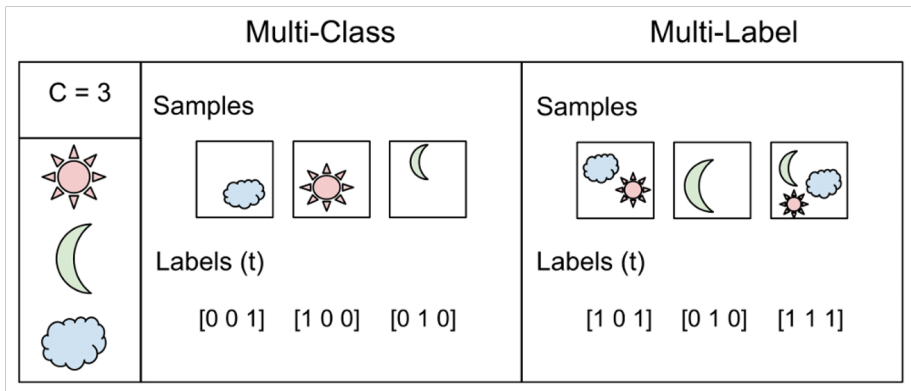
# Experimental Results with Feature Normalization

Faster training/feature collapse with ResNet on CIFAR100 with feature normalization



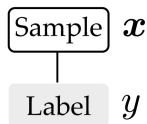


# Neural Collapse for Multi-Label Learning



# Multi-label Learning Setup

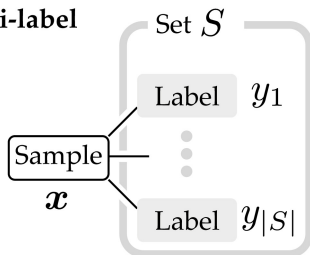
## Single-label



$$\mathcal{L}_{\text{CE}}(\psi_{\Theta}(\mathbf{x}), y)$$

Loss

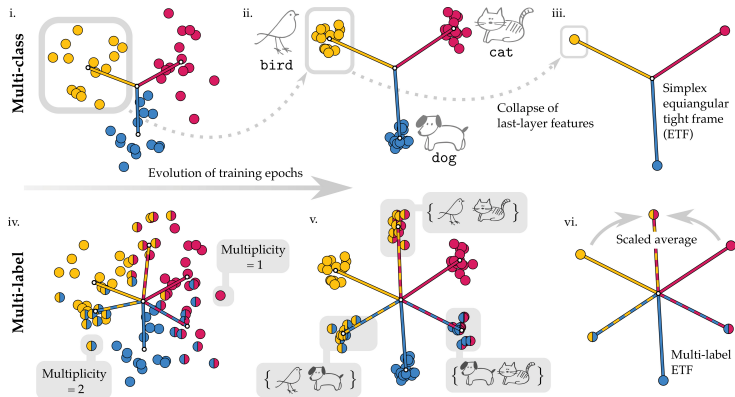
## Multi-label



$$\sum_{i=1}^{|S|} \mathcal{L}_{\text{CE}}(\psi_{\Theta}(\mathbf{x}), y_i)$$

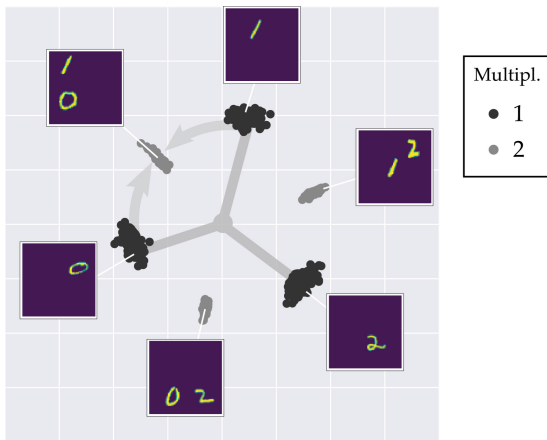
"Pick-all" Loss

# Last-Layer Geometry of Multi-label Learning



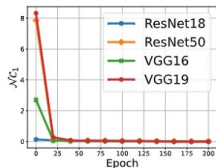
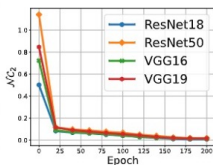
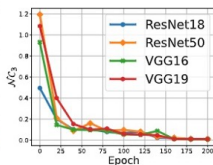
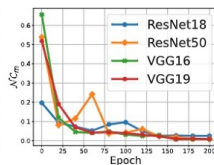
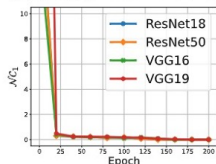
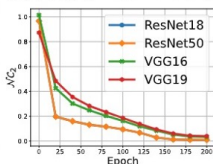
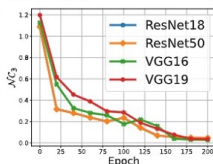
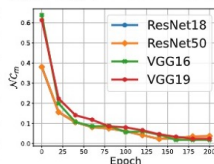
- Neural collapse in multi-label learning with 3 classes where the colors denote the class label;
- Respectively, left/mid/right panel shows representations during early/mid/late phase of training unconstrained feature model.

# Multilabel-MNIST Synthetic Example



- Experiments with simple MLP architectures.
- The ETF structure still holds for data imbalancedness.

# Neural Collapse for Multi-Label Learning

(a)  $\mathcal{NC}_1$  (MLab-MNIST)(b)  $\mathcal{NC}_2$  (MLab-MNIST)(c)  $\mathcal{NC}_3$  (MLab-MNIST)(d)  $\mathcal{NC}_m$  (MLab-MNIST)(e)  $\mathcal{NC}_1$  (MLab-Cifar10)(f)  $\mathcal{NC}_2$  (MLab-Cifar10)(g)  $\mathcal{NC}_3$  (MLab-Cifar10)(h)  $\mathcal{NC}_m$  (MLab-Cifar10)

# Outline

- 1 Low-Dimensional Representation: Neural Collapse (NC)
- 2 Understanding NC from Optimization
- 3 Prevalence of NC under Different Training Scenarios
- 4 Conclusion

## References

- 1 Z. Zhu\*, T. Ding\*, J. Zhou, X. Li, C. You, J. Sulam, and Q. Qu, A Geometric Analysis of Neural Collapse with Unconstrained Features, NeurIPS'2021 (spotlight, top 3%).
- 2 J. Zhou\*, X. Li\*, T. Ding, C. You, Q. Qu\*, Z. Zhu\*. On the Optimization Landscape of Neural Collapse under MSE Loss: Global Optimality with Unconstrained Features. ICML'2022.
- 3 C. Yaras\*, P. Wang\*, Z. Zhu, L. Balzano, Q. Qu, Neural Collapse with Normalized Features: A Geometric Analysis over the Riemannian Manifold. NeurIPS'2022.
- 4 J. Zhou, C. You, X. Li, K. Liu, S. Liu, Q. Qu, Z. Zhu. Are All Losses Created Equal? A Neural Collapse Perspective. NeurIPS'2022.
- 5 P. Wang\*, H. Liu\*, C. Yaras\*, L. Balzano, Q. Qu. Linear Convergence Analysis of Neural Collapse with Unconstrained Features. NeurIPS OPT Workshop, 2022.

# Conclusion and Coming Attractions

*Learning* common deep networks for low-dim structure

- **Low-dimensional features:** understand low-dim. features (sparse and neural collapse (NC)) learned in deep classifiers trained with one-hot labeling based losses in generic settings

**Thank You! Questions?**



# Call for Papers

- IEEE JSTSP Special Issue on Seeking Low-dimensionality in Deep Neural Networks (SLOWDNN) Manuscript Due: **Nov. 30, 2023.**
- Conference on Parsimony and Learning (CPAL) January 2024, Hongkong, Manuscript Due: **Aug. 28, 2023.**

