



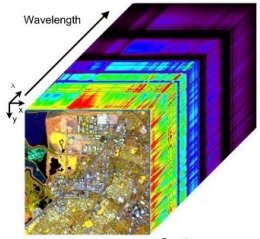
# From Shallow to Deep Representation Learning in Imaging and Beyond: Global Nonconvex Theory and Algorithms

Qing Qu

Dept. of EECS, University of Michigan

September 7, 2021

# Data Increasingly Massive & High-Dimensional...



hyperspectral imaging



autonomous driving

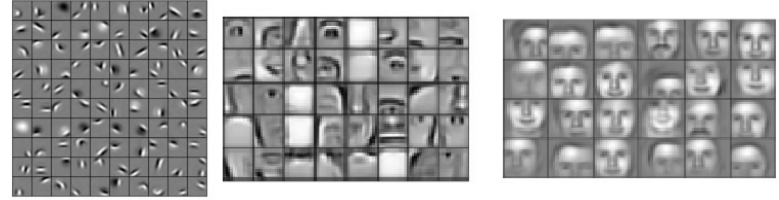


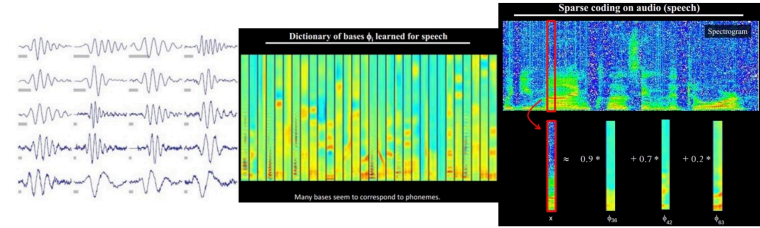
image representations



social network



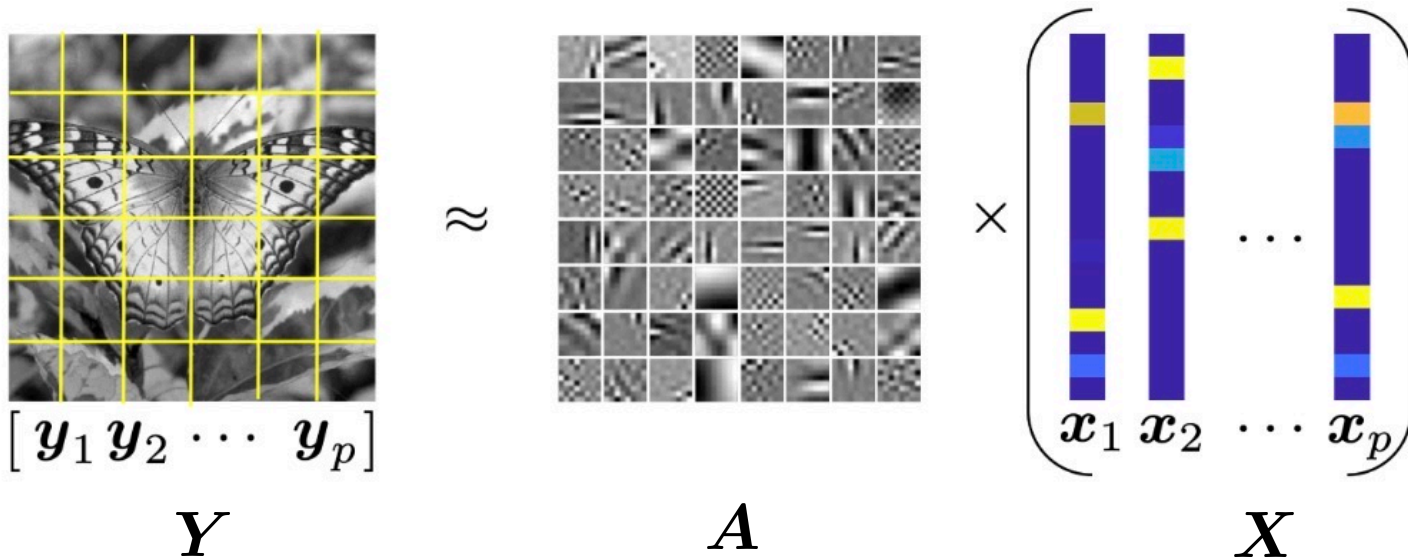
healthcare



audio representations

Data representation is *critical* for modern machine learning methods.

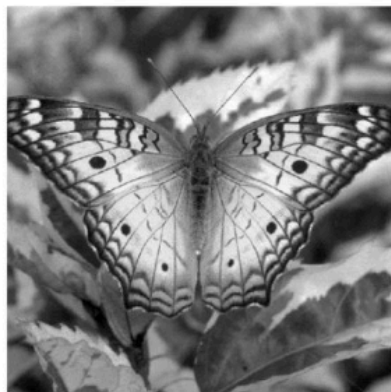
# Unsupervised Learning



- Learning sparsely-used dictionaries:

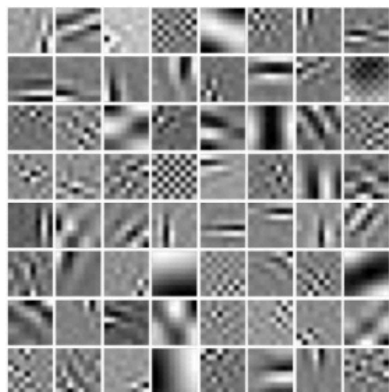
Given  $Y \in \mathbb{R}^{n \times p}$ , jointly find overcomplete dictionary  $A \in \mathbb{R}^{n \times m}$  and sparse  $X \in \mathbb{R}^{m \times p}$ .

# Unsupervised Learning

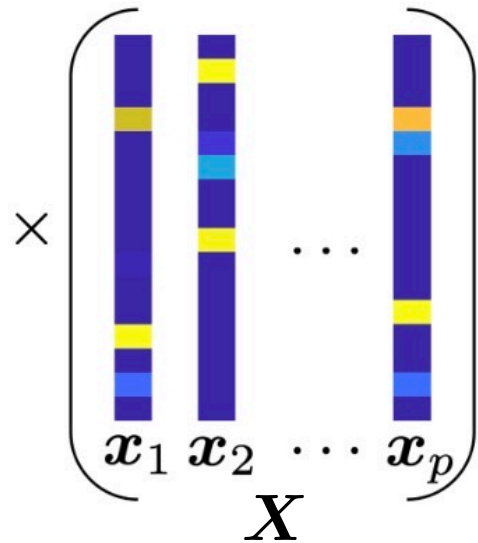


$Y$

$\approx$



$A$



- Learning sparsely-used dictionaries:

$$\min_{A \in \mathcal{M}, X} f(Y, A \cdot X) + \lambda \cdot g(X)$$

**data fidelity**                      **regularizer**

# Unsupervised Learning

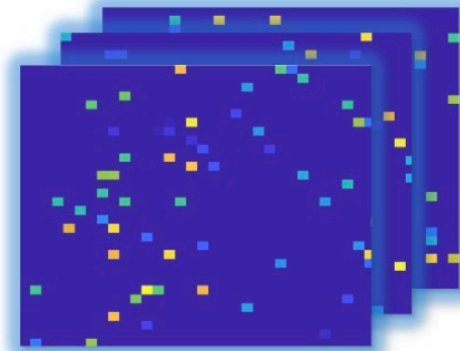


$y_i$

$$\approx \sum_k$$



$\otimes$



$$\approx \sum_k$$

$a_k$

$\otimes$

$x_{ki}$

- **Learning convolutional dictionaries:**

Given  $\{y_i\}_i$ , jointly learn convolutional dictionaries  $\{a_i\}_i$  and sparse coefficients  $\{x_{ki}\}_{i,k}$ .

➤ Translation invariant, can be viewed as one layer of ConvNets



Denoising

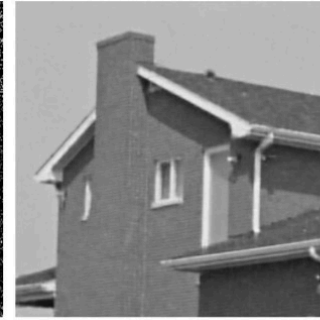
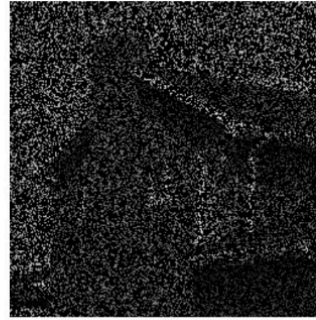


Image Restoration



Super Resolution

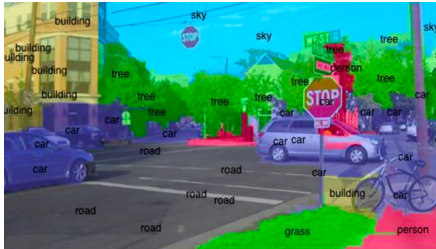


Image Half-toning

- Image courtesy of Julien Mairal et al.

# Supervised (Deep) Learning

Deep learning has attained superior performances for many tasks in practice:



Computer vision



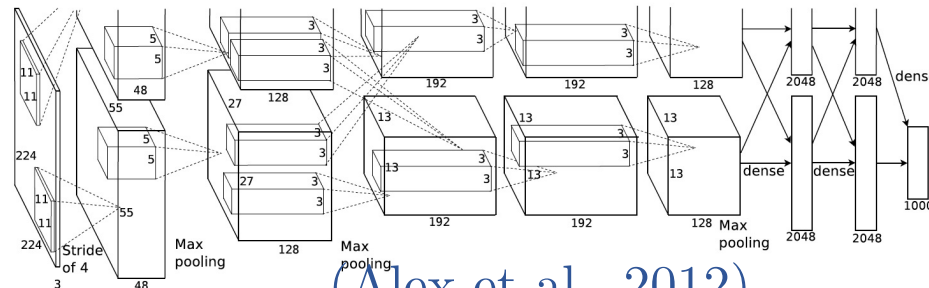
Natural language processing



Gameplay



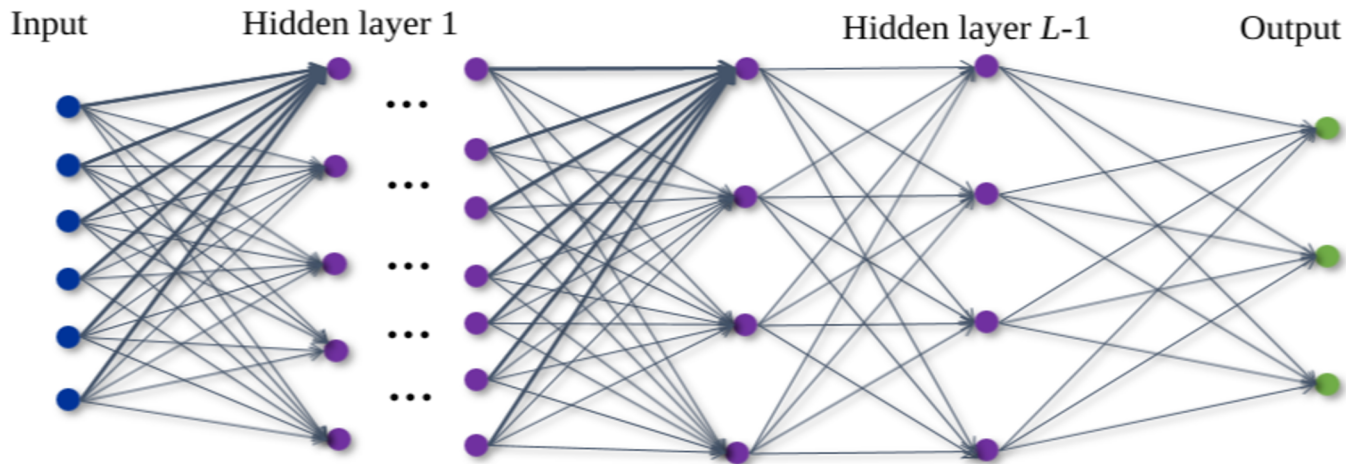
Protein modeling



(Alex et al., 2012)

→ "Cat"

# Training Deep Neural Networks



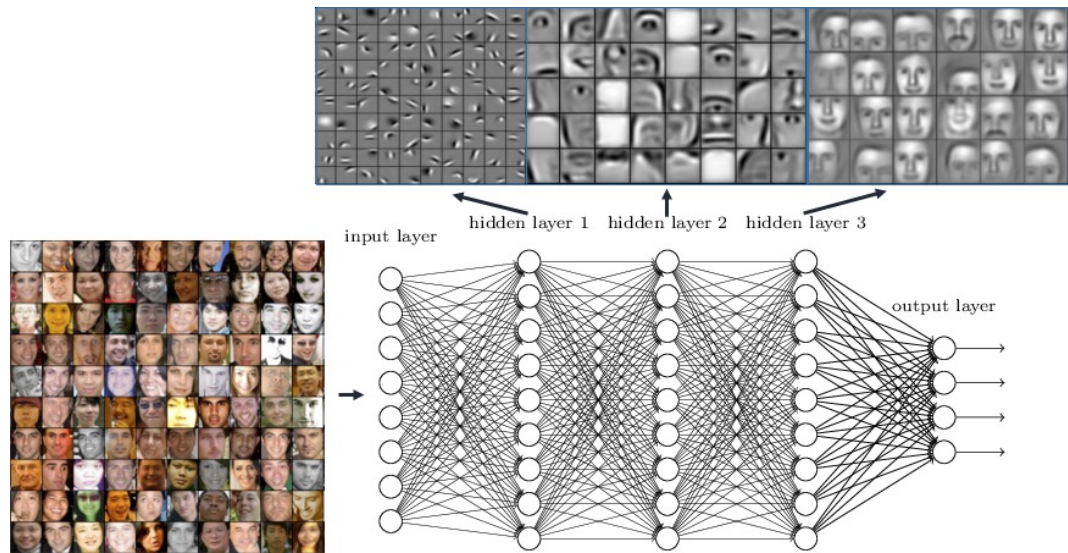
$$\psi_{\Theta}(x) = \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 x + \mathbf{b}_1) + \mathbf{b}_{L-1}) + \mathbf{b}_L$$

$$\Theta := \{\mathbf{W}_\ell, \mathbf{b}_\ell\}_{\ell=1}^L \quad \sigma(\cdot): \text{nonlinear activations}$$

↑                      ↑  
weights                bias



# Training Deep Neural Networks



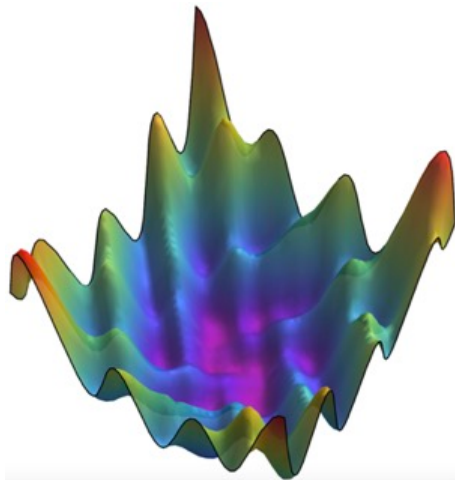
$$\min_{\Theta} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}_{\text{CE}}(\psi_{\Theta}(\mathbf{x}_{k,i}), \mathbf{y}_k) + \lambda \|\Theta\|_F^2$$

$\mathbf{x}_{k,i}$ :  $i$ -th input in the  $k$ -th class

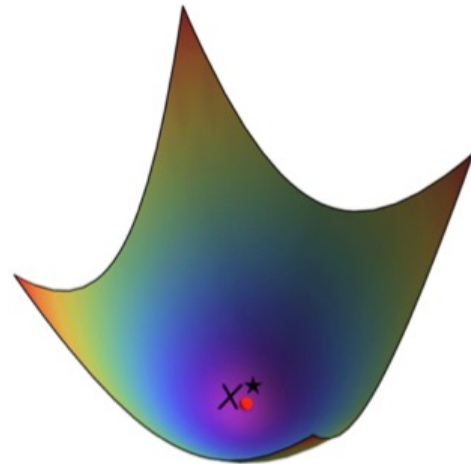
$\mathbf{y}_k$ : One-hot vector for the  $k$ -th class

# Nonconvex Problems in Representation Learning

$$\min_{\mathbf{x}} f(\mathbf{x}), \text{ s.t. } \mathbf{x} \in \mathbb{R}^n$$

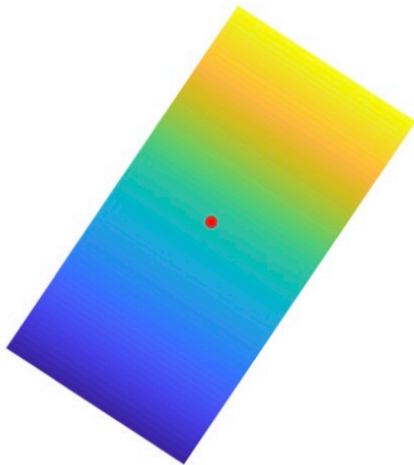


Nonconvex landscape

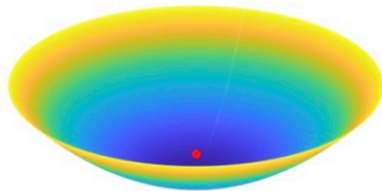


Convex landscape

# General Nonconvex Problems

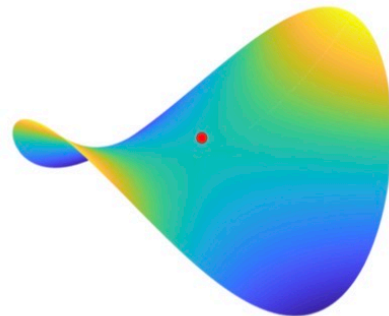


Noncritical Point ( $\nabla\varphi \neq \mathbf{0}$ )



Minimizer

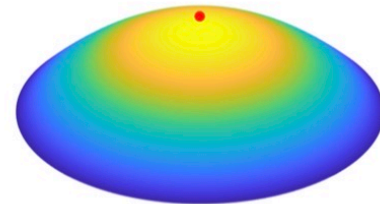
$$\nabla^2\varphi \succ \mathbf{0}$$



Saddle

$$\lambda_{\min}\nabla^2\varphi < 0$$

$$\lambda_{\max}\nabla^2\varphi > 0$$



Maximizer

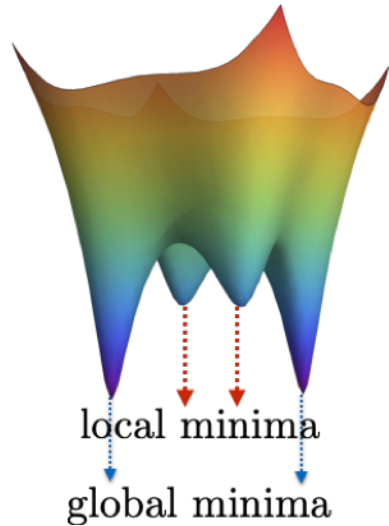
$$\nabla^2\varphi \prec \mathbf{0}$$

Critical Points ( $\nabla\varphi = \mathbf{0}$ )

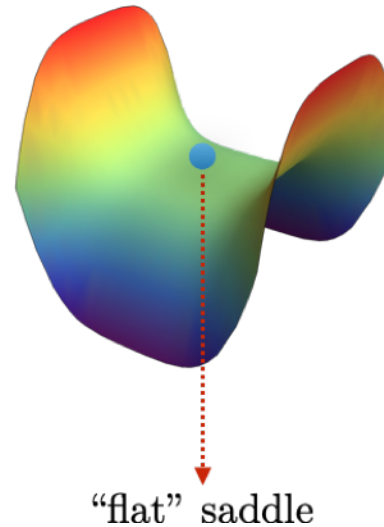
# General Nonconvex Problems

$$\min_{\mathbf{x}} f(\mathbf{x}), \text{ s.t. } \mathbf{x} \in \mathbb{R}^n$$

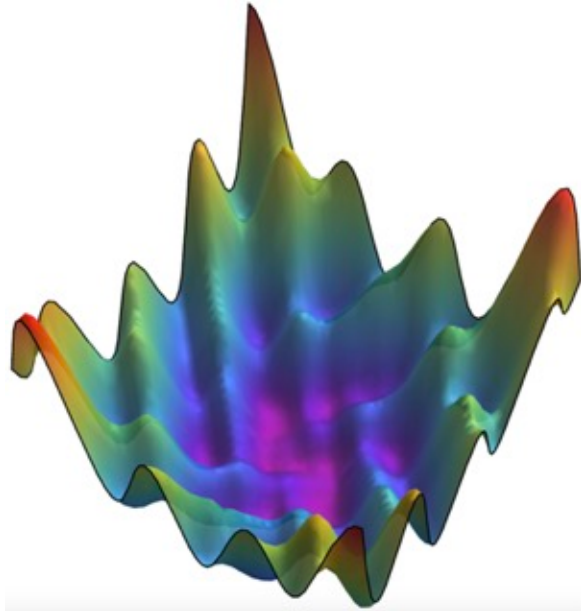
“bad” local minimizers



“flat” saddle points



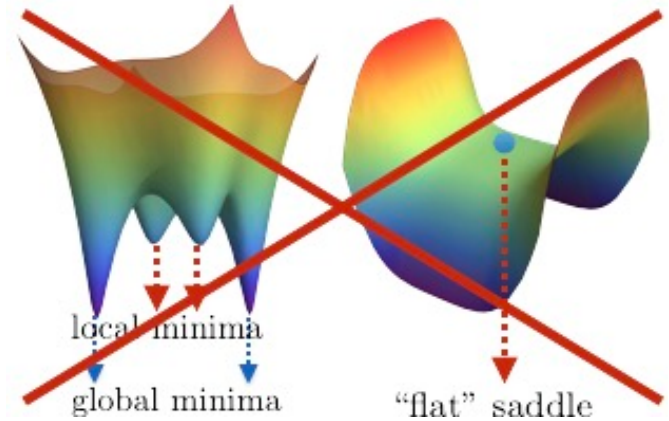
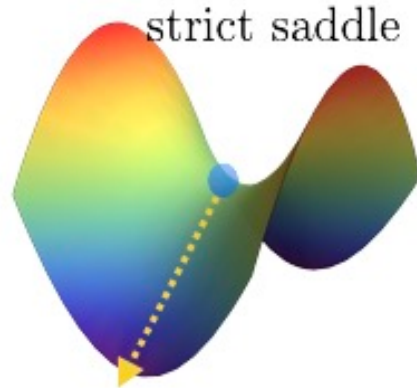
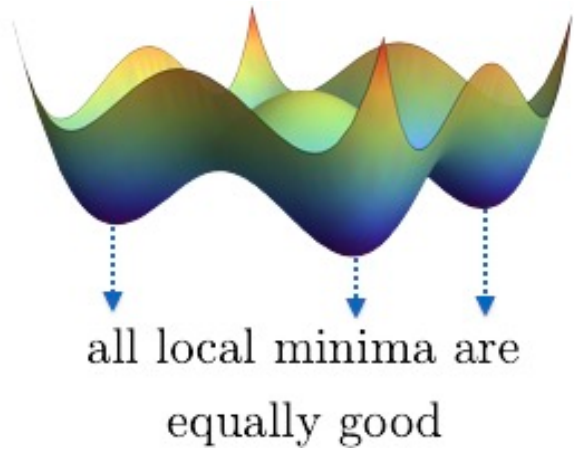
# General Nonconvex Problems



$$\min_{\mathbf{x}} f(\mathbf{x}), \text{ s.t. } \mathbf{x} \in \mathbb{R}^n$$

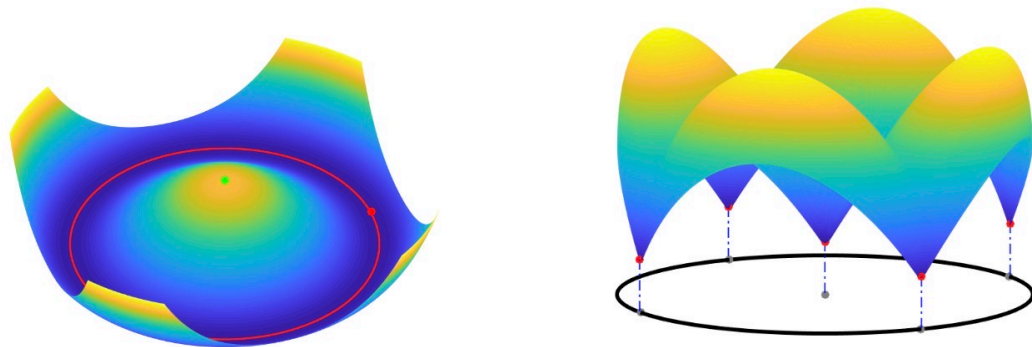
**In the worst case**, even finding a local minimizer is NP-hard (Murty et al. 1987)

# Optimizing Nonconvex Problems Globally



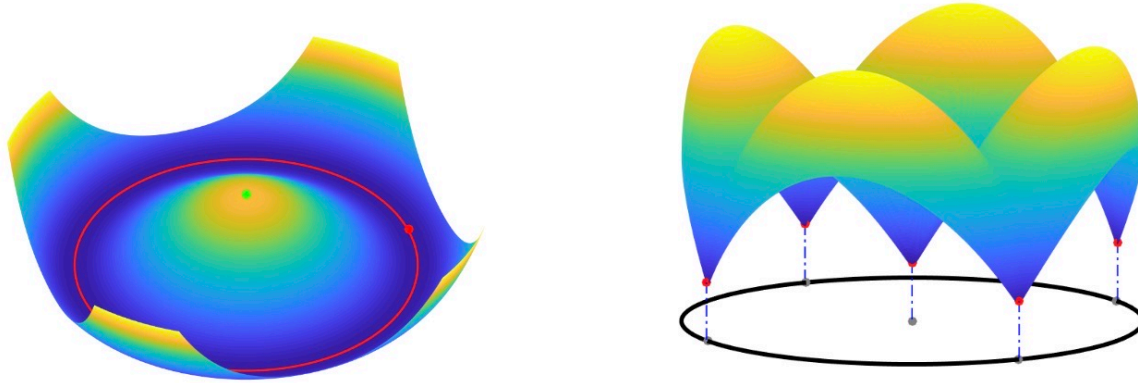
**Benign nonconvex landscapes enable efficient  
global optimization!**

# Nonconvex Problems with Benign Landscape



- Generalized Phase Retrieval [Sun'18]
- Low-rank Matrix Recovery [Ma'16, Jin'17, Chi'19]
- Sparse Dictionary Learning [Sun'16, Qu'20]
- (Orthogonal) Tensor Decomposition [Ge'15]
- Sparse Blind Deconvolution [Zhang'17, Li'18, Kuo'19]
- Deep Linear Network [Kawaguchi'16]
- ...

# Nonconvex Problems with Benign Landscape



- **Q. Qu** (\*), Z. Zhu (\*), X. Li, M. C. Tsakiris, J. Wright, R. Vidal, Finding the Sparsest Vectors in a Subspace: Theory, Algorithms, and Applications, In Submission, 2020.
- Y. Zhang, **Q. Qu**, J. Wright, Symmetry and Geometry in Nonconvex Optimization, In Submission, 2020.



# Outline of this Talk

- **Learning Shallow Representations:  
(Convolutional) Dictionary Learning**
- Learning Deep Representations:  
Deep Neural Collapse
- Conclusion & Discussion

# Landscape Analysis of Dictionary Learning

1. **Q. Qu**, Y. Zhai, X. Li, Y. Zhang, Z. Zhu, Analysis of optimization landscapes for overcomplete learning, *ICLR'20*, (oral, top 1.9%)
  2. Y. Lau (\*), **Q. Qu**(\*), H. Kuo, P. Zhou, Y. Zhang, J. Wright, Short-and-sparse Deconvolution – A Geometric Approach, *ICLR'20*
- Provide the first **global nonconvex landscape** analysis for *convolutional/overcomplete* dictionary learning problems.
  - Efficient nonconvex optimization methods to global solutions with applications in imaging.

# Convolutional Dictionary Learning (DL)

Given multiple measurements  $\{\mathbf{y}_i\}_i$  of circulant convolution

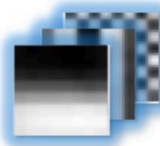
$$\mathbf{y}_i = \sum_{k=1}^K \mathbf{a}_k \circledast \mathbf{x}_{ki}, \quad (1 \leq i \leq p)$$

can we learn all  $\{\mathbf{a}_k\}_k$  and  $\{\mathbf{x}_{ki}\}_{k,i}$  simultaneously?

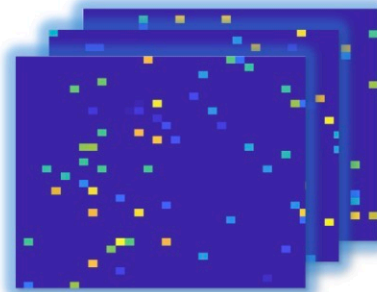


$\mathbf{y}_i$

$\approx \sum_k$



$\circledast$



$\approx \sum_k$

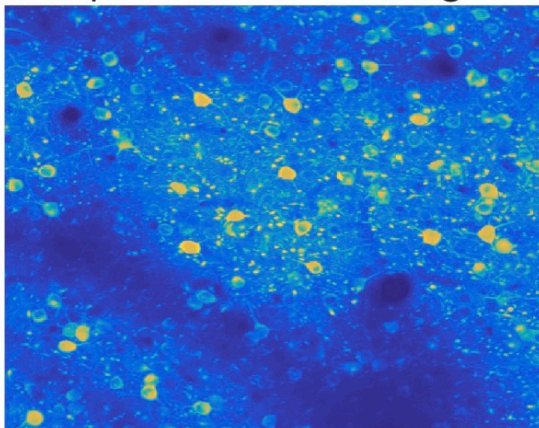
$\mathbf{a}_k$

$\circledast$

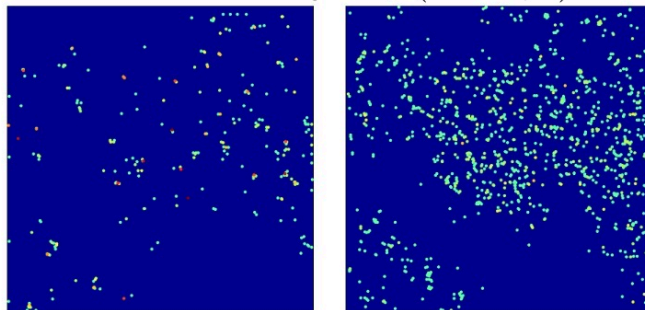
$\mathbf{x}_{ki}$

# Convolutional Dictionary Learning (DL)

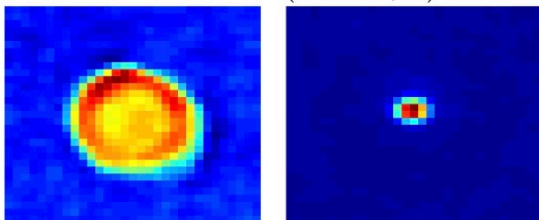
two-photon calcium image  $Y$



activation map  $X_k$  ( $k = 1, 2$ )

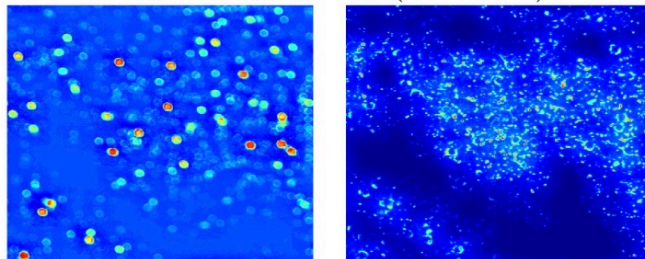


kernel  $A_k$  ( $k = 1, 2$ )



reconstruction

$$Y_k = A_k \circledast X_k \quad (k = 1, 2)$$



# Convolutional DL vs. Overcomplete DL

For each  $\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i$ , we can *equivalently* rewrite in the matrix form as:

$$\mathbf{C}_{\mathbf{y}_i} = \mathbf{C}_a \cdot \mathbf{C}_{\mathbf{x}_i}, \quad 1 \leq i \leq p$$

where a circulant matrix

$$\mathbf{C}_a = \begin{pmatrix} \begin{matrix} \color{blue} & \color{yellow} & \color{cyan} & \color{blue} & \color{purple} & \color{green} & \color{cyan} & \color{blue} \\ \color{cyan} & \color{blue} & \color{yellow} & \color{cyan} & \color{blue} & \color{purple} & \color{green} & \color{cyan} \\ \color{green} & \color{cyan} & \color{blue} & \color{yellow} & \color{cyan} & \color{blue} & \color{purple} & \color{green} \\ \color{purple} & \color{green} & \color{cyan} & \color{blue} & \color{yellow} & \color{cyan} & \color{blue} & \color{purple} \\ \color{blue} & \color{purple} & \color{green} & \color{cyan} & \color{blue} & \color{yellow} & \color{cyan} & \color{blue} \\ \color{yellow} & \color{cyan} & \color{blue} & \color{purple} & \color{green} & \color{cyan} & \color{blue} & \color{yellow} \\ \color{cyan} & \color{blue} & \color{purple} & \color{green} & \color{cyan} & \color{blue} & \color{yellow} & \color{cyan} \\ \color{yellow} & \color{cyan} & \color{blue} & \color{purple} & \color{green} & \color{cyan} & \color{blue} & \color{yellow} \end{matrix} \\ s_1[\mathbf{a}] & s_2[\mathbf{a}] & \cdots & s_n[\mathbf{a}] \end{pmatrix} \in \mathbb{R}^{n \times n}$$

# Convolutional DL vs. Overcomplete DL

For each sample

$$\mathbf{y}_i = \sum_{k=1}^K \mathbf{a}_k \circledast \mathbf{x}_{ki},$$

equivalently,

$$\mathbf{C}_{\mathbf{y}_i} = \underbrace{\begin{bmatrix} \mathbf{C}_{\mathbf{a}_1} & \cdots & \mathbf{C}_{\mathbf{a}_K} \end{bmatrix}}_{\text{overcomplete } \mathbf{A}_0} \cdot \underbrace{\begin{bmatrix} \mathbf{C}_{\mathbf{x}_{1i}} \\ \vdots \\ \mathbf{C}_{\mathbf{x}_{Ki}} \end{bmatrix}}_{\text{sparse } \mathbf{X}_i},$$

# Convolutional DL vs. Overcomplete DL

Given  $Y = A_0 \cdot X_0$ , learn **overcomplete**  $A_0$  and **sparse**  $X_0$ ?

$$\left[ \begin{array}{c} \text{data } Y \\ C_{y_1} \cdots C_{y_p} \end{array} \right] = \left[ \begin{array}{c} \text{dictionary } A_0 \\ C_{a_1} \cdots C_{a_K} \end{array} \right] \left[ \begin{array}{c} \text{sparse } X_0 \\ X_1 \cdots X_p \end{array} \right]$$

**data  $Y$**                       **dictionary  $A_0$**                       **sparse  $X_0$**

# Relationship to Dictionary Learning

We can find **one column** of  $\mathbf{A}_0$  via

$$\min_{\mathbf{q}} f_{DL}(\mathbf{q}) = - \|\mathbf{Y}^\top \mathbf{q}\|_4^4, \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1.$$

The underlying reasoning is that, in expectation

$$\mathbb{E}_{\mathbf{X}} \left[ \|\mathbf{Y}^\top \mathbf{q}\|_4^4 \right] = \mathbb{E}_{\mathbf{X}} \left[ \|\mathbf{X}^\top \mathbf{A}_0^\top \mathbf{q}\|_4^4 \right] = c_1 \|\mathbf{A}_0^\top \mathbf{q}\|_4^4 + c_2$$

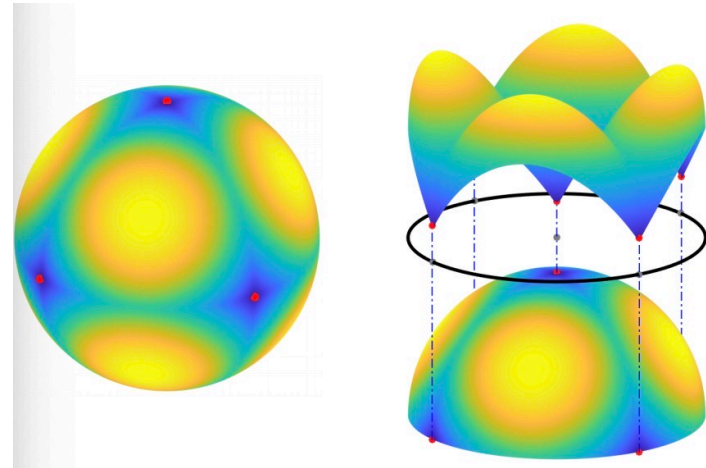
for  $\mathbf{X}$  following some sparse zero-mean distributions (e.g., Bernoulli-Gaussian)



# Relationship to Dictionary Learning

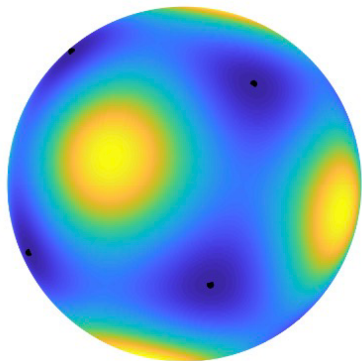
Given  $\mathbf{A}_0 = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_m] \in \mathbb{R}^{n \times m}$

$$\min_{\mathbf{q}} - \|\mathbf{A}_0^\top \mathbf{q}\|_4^4, \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1.$$



- When  $m \leq n$ , and  $\{\mathbf{a}_i\}_{i=1}^m$  are **orthogonal**, existing result [Ge'15] has shown that the function is a **strict saddle function** with benign optimization landscape, all global solutions are approximately  $\{\pm \mathbf{a}_i\}_{i=1}^m$ .
- The analysis of orthogonal case **cannot** be generalized to overcomplete settings.

# Global Landscape of Overcomplete DL



$$\min_{\mathbf{q}} f_{DL}(\mathbf{q}) = - \|\mathbf{Y}^\top \mathbf{q}\|_4^4, \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1.$$

**Theorem (Informal)** *Suppose that (i)  $K = m/n$  is a constant, (ii)  $\mathbf{A}_0$  is near orthogonal, and (iii)  $p \geq \Omega(\text{poly}(n))$ . Then with high probability every critical point of  $f(\mathbf{q})$  is either*

- *a **strict saddle point** exhibits negative curvature;*
- *or close to a **target solution**: one column  $\mathbf{a}_i$  of  $\mathbf{A}$ .*

# Assumptions on $A$ (Near Orthogonal)

- Row orthogonal: unit norm **tight frame** (UNTF)

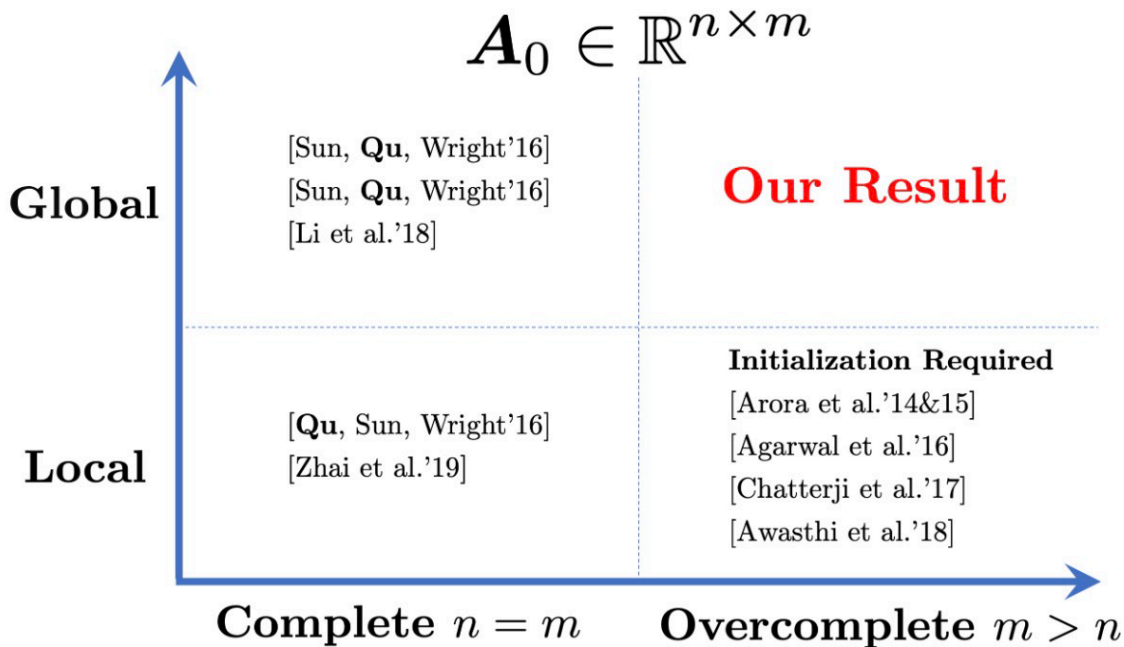
$$\sqrt{\frac{n}{m}} \mathbf{A}_0 \mathbf{A}_0^\top = \mathbf{I}, \quad \|\mathbf{a}_i\|_2 = 1.$$

- **Incoherence** of the columns (near orthogonal)

$$\max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle| \leq \mu.$$

# Overcomplete Dictionary Learning

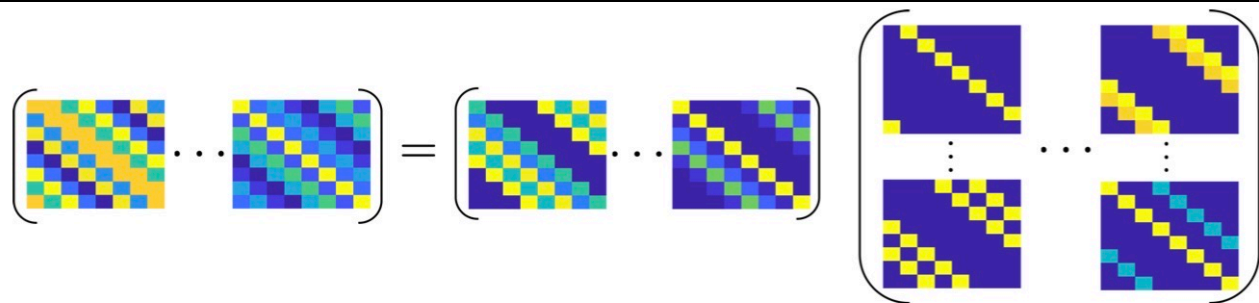
Given  $Y = A_0 \cdot X_0$ , learn **overcomplete**  $A_0$  and **sparse**  $X_0$ ?



# From Overcomplete DL to Convolutional DL

Find one shift of the kernel  $\mathbf{a}_i$  via

$$\min_{\mathbf{q}} f_{CDL}(\mathbf{q}) = - \|\mathbf{q}^\top \mathbf{Y}\|_4^4, \quad \text{s.t. } \mathbf{q} \in \mathbb{S}^{n-1}$$



$$\begin{bmatrix} \mathbf{C}_{\mathbf{y}_1} & \cdots & \mathbf{C}_{\mathbf{y}_p} \end{bmatrix} = \begin{bmatrix} \mathbf{C}_{\mathbf{a}_1} & \cdots & \mathbf{C}_{\mathbf{a}_K} \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_p \end{bmatrix}$$

data  $\mathbf{Y}$

dictionary  $\mathbf{A}_0$

sparse  $\mathbf{X}_0$

# Convolutional Dictionary Learning

Find one shift of the kernel  $\mathbf{a}_i$  via

$$\min_{\mathbf{q}} f_{CDL}(\mathbf{q}) = - \|\mathbf{q}^\top \mathbf{P}\mathbf{Y}\|_4^4, \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}$$

- Preconditioning matrix:

$$\mathbf{P} = \left( (\theta n K)^{-1} \mathbf{Y}\mathbf{Y}^\top \right)^{-1/2} \approx \left( \mathbf{A}_0 \mathbf{A}_0^\top \right)^{-1/2}$$

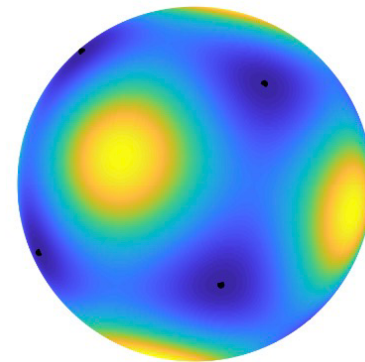
- Effective dictionary is tight frame (but not necessarily unit norm)

$$\mathbf{P}\mathbf{Y} \approx \underbrace{\left( \mathbf{A}_0 \mathbf{A}_0^\top \right)^{-1/2} \mathbf{A}_0}_{\mathbf{A}} \mathbf{X}_0 = \mathbf{A} \mathbf{X}_0$$

# Convolutional Dictionary Learning

Find one shift of the kernel  $\mathbf{a}_i$  via

$$\min_{\mathbf{q}} - \|\mathbf{q}^\top \mathbf{A} \mathbf{X}_0\|_4^4, \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}$$



- Preconditioning matrix:

$$\mathbf{P} = \left( (\theta n K)^{-1} \mathbf{Y} \mathbf{Y}^\top \right)^{-1/2} \approx \left( \mathbf{A}_0 \mathbf{A}_0^\top \right)^{-1/2}$$

- Effective dictionary is tight frame (but not necessarily unit norm)

$$\mathbf{P} \mathbf{Y} \approx \underbrace{\left( \mathbf{A}_0 \mathbf{A}_0^\top \right)^{-1/2} \mathbf{A}_0}_{\mathbf{A}} \mathbf{X}_0 = \mathbf{A} \mathbf{X}_0$$

# Local Landscape of Convolutional DL

**Theorem (Informal)** *Suppose that (i)  $K = m/n$  is a constant, (ii)  $\mathbf{A}$  is near orthogonal, and (iii)  $p \geq \Omega(\text{poly}(n))$ . **Locally**, every critical point of  $f_{CDL}(\mathbf{q})$  is either*

- *a **strict saddle point** exhibits negative curvature;*
- *or close to a **target solution**: a precond. shift of  $\mathbf{a}_i$ .*

- We show the result over a local level-set

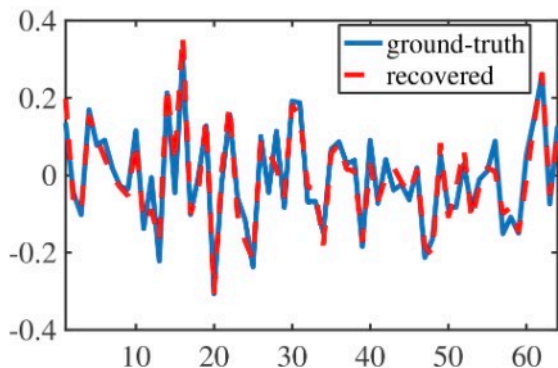
$$\mathcal{R}_{CDL} := \left\{ \mathbf{q} \in \mathbb{S}^{n-1} \mid \mathbb{E}_{\mathbf{X}} [f_{CDL}(\mathbf{q})] \leq -c \|\mathbf{A}^\top \mathbf{q}\|_3^2 \right\},$$

- We can cook up smart yet simple initializations, with all future iterations stay in the region.

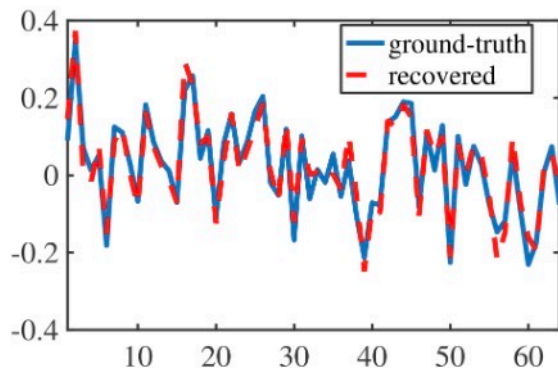


# Learning Random Filters

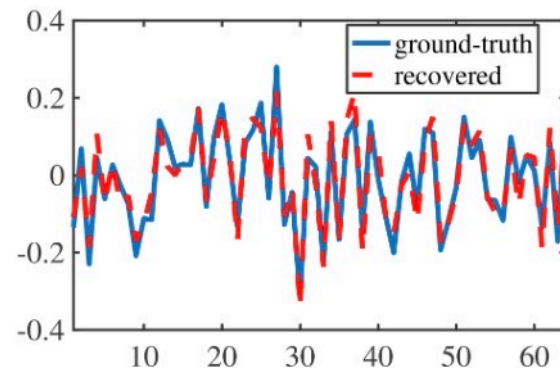
$$\min_{\mathbf{q}} f_{CDL}(\mathbf{q}) = - \|\mathbf{q}^\top \mathbf{P}\mathbf{Y}\|_4^4, \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}$$



Filter 1



Filter 2



Filter 3

**Learning 3 random filters by the proposed approach.**

# From Theory to Practical Methods

- Recovering one filter:

$$\min_{\mathbf{q}} f_{CDL}(\mathbf{q}) = - \|\mathbf{q}^\top \mathbf{P}\mathbf{Y}\|_4^4, \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}$$

- Finding all filters via Bilinear Lasso formulation:

$$\min_{\mathbf{a}_k, \mathbf{x}_k} \frac{1}{2} \left\| \mathbf{y} - \sum_{k=1}^N \mathbf{a}_k \circledast \mathbf{x}_k \right\|^2 + \lambda \sum_{k=1}^N \|\mathbf{x}_k\|_1, \quad \text{s.t.} \quad \|\mathbf{a}_k\| = 1.$$

# From Theory to Practical Methods

- Finding all filters via Bilinear Lasso formulation:

$$\min_{\mathbf{a}_k, \mathbf{x}_k} \frac{1}{2} \left\| \mathbf{y} - \sum_{k=1}^N \mathbf{a}_k \circledast \mathbf{x}_k \right\|^2 + \lambda \sum_{k=1}^N \|\mathbf{x}_k\|_1, \quad \text{s.t. } \|\mathbf{a}_k\| = 1.$$

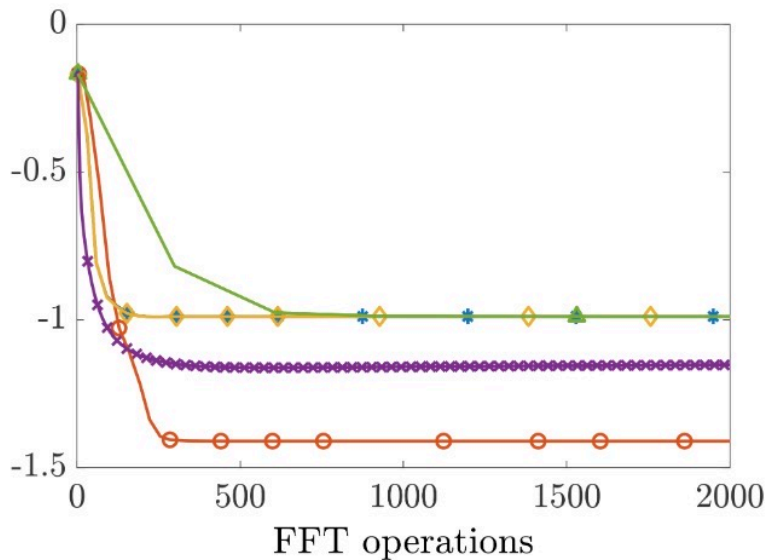
- Optimization. Alternating descent method

$$\mathbf{x} \leftarrow \text{prox}(\mathbf{x} - \tau \cdot \nabla_{\mathbf{x}} \varphi_{\text{BL}}(\mathbf{a}, \mathbf{x}))$$

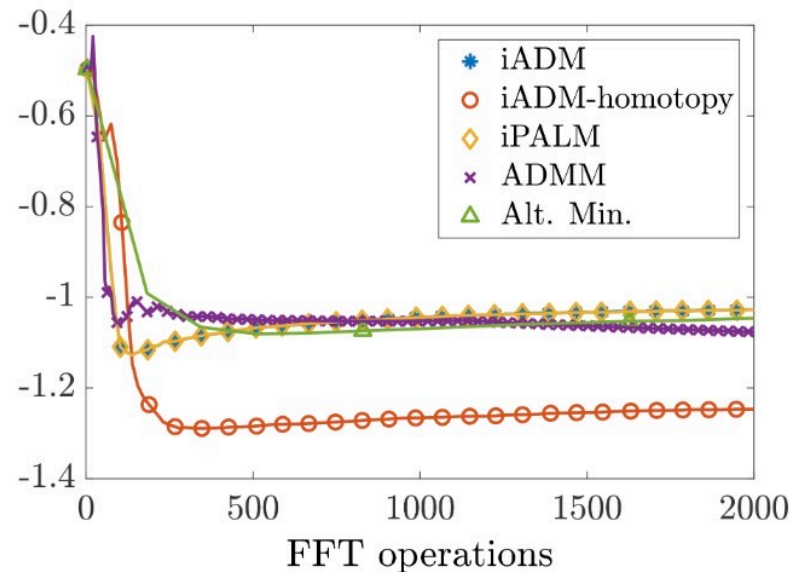
$$\mathbf{a} \leftarrow \mathcal{P}_{\mathbb{S}^{n-1}}(\mathbf{a} - t \cdot \text{grad}_{\mathbf{a}} \varphi_{\text{BL}}(\mathbf{a}, \mathbf{x})),$$

with few extra caveats.

# From Theory to Practical Methods



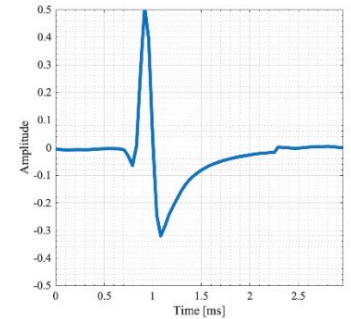
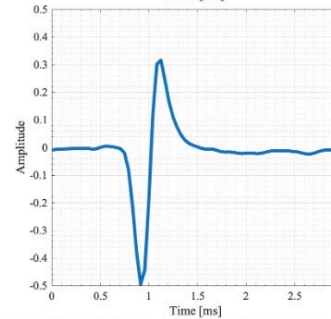
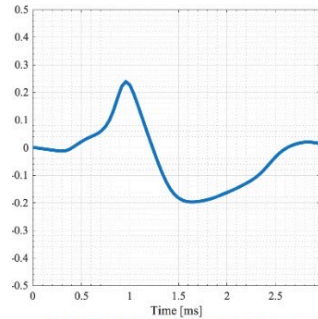
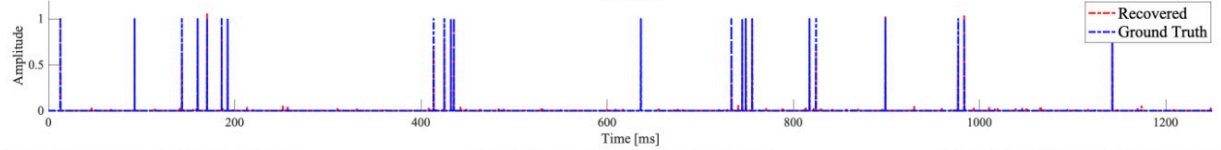
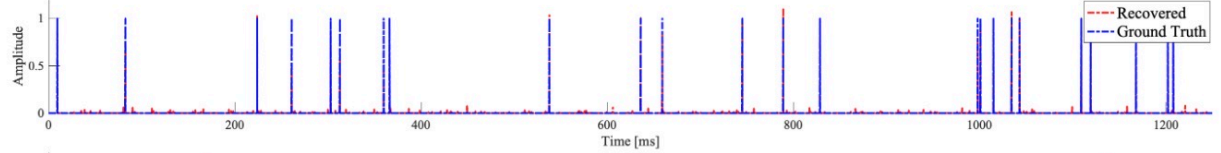
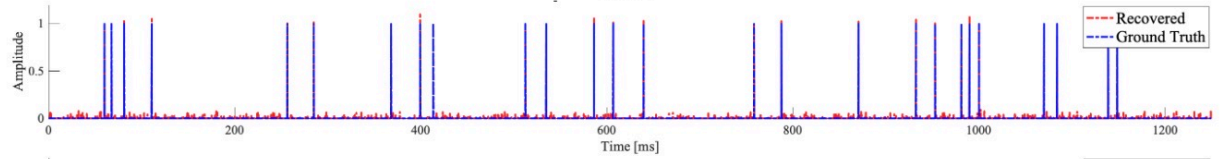
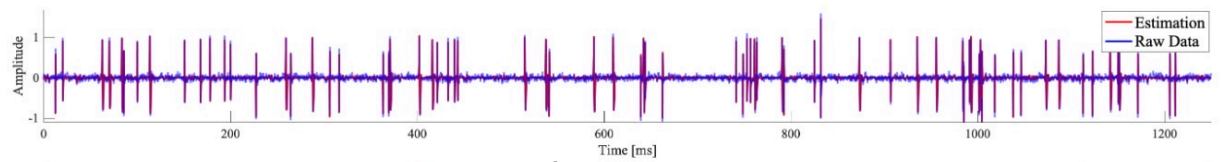
all-pass filter



low-pass filter

Comparison of convergence (time).

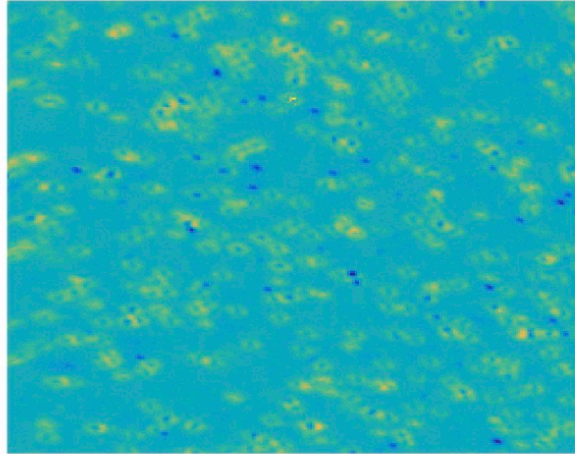
# Spike Sorting



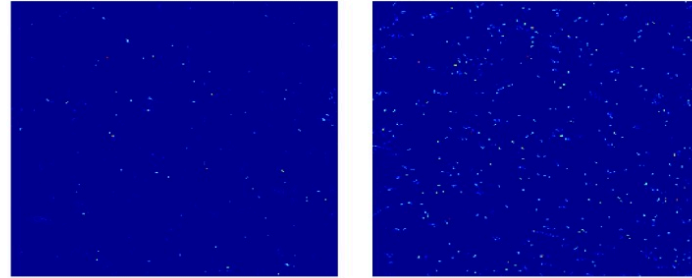
- [https://vis.caltech.edu/~rodri/Wave\\_clus/Wave\\_clus\\_home.htm](https://vis.caltech.edu/~rodri/Wave_clus/Wave_clus_home.htm)

# Defects Detection in Scan Tunneling Microscopy

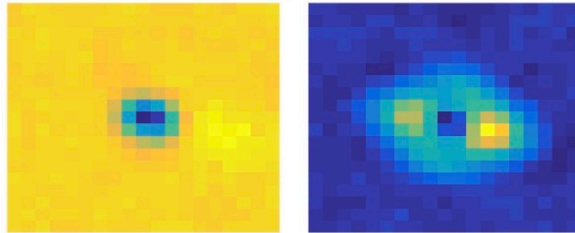
STM image  $Y$



activation map  $X_k$  ( $k = 1, 2$ )

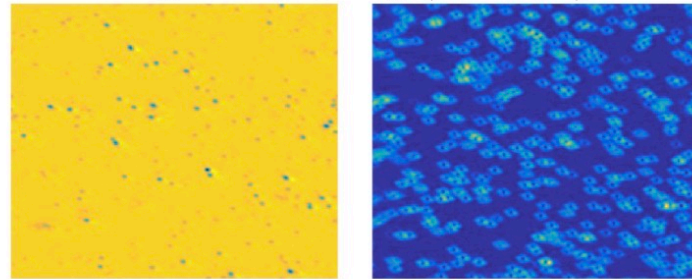


kernel  $A_k$  ( $k = 1, 2$ )



reconstruction

$$Y_k = A_k \otimes X_k \quad (k = 1, 2)$$



# Outline of this Talk

- Learning Shallow Representations:  
(Convolutional) Dictionary Learning
- **Learning Deep Representations:  
Deep Neural Collapse**
- Conclusion & Discussion

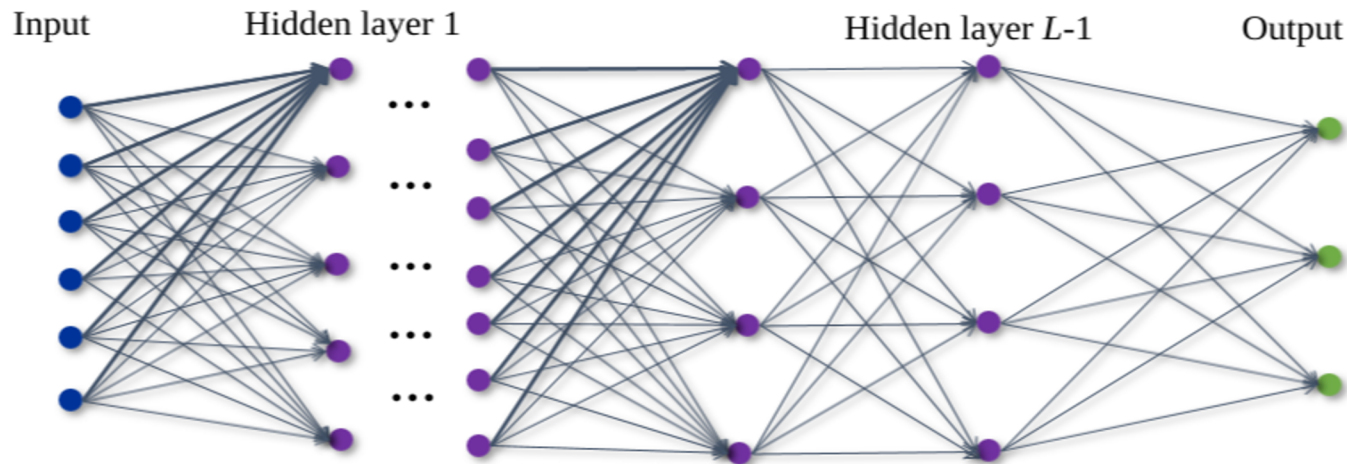
# Understanding Deep Neural Networks

Z. Zhu, T. Ding, J. Zhou, X. Li, C. You, J. Sulam, and Q. Qu, [A Geometric Analysis of Neural Collapse with Unconstrained Features](#), *arXiv Preprint arXiv:2105.02375*, May 2021.

- Analyzes the **global landscape** of the training loss based on the **unconstrained feature model**
- Explains the ubiquity of **Neural Collapse** of the learned representations of the network



# Understanding Deep Neural Networks



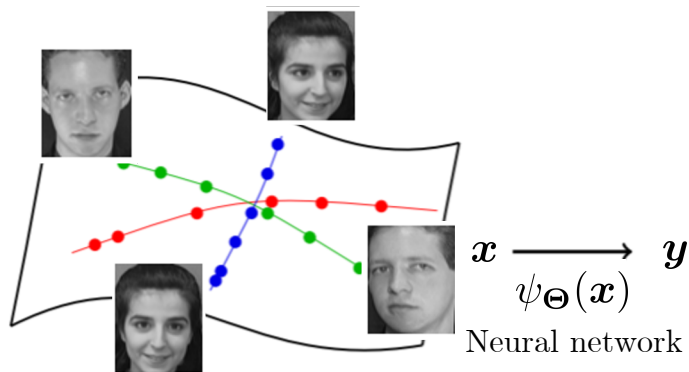
$$\psi_{\Theta}(x) = \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \cdots \sigma(\mathbf{W}_1 x + \mathbf{b}_1) + \mathbf{b}_{L-1}) + \mathbf{b}_L$$

$$\Theta := \{\mathbf{W}_\ell, \mathbf{b}_\ell\}_{\ell=1}^L \quad \sigma(\cdot): \text{nonlinear activations}$$

↑                      ↑  
weights                bias

# Terminology for Classification

- Labels:  $k = 1, \dots, K$ 
  - $K = 10$  classes (MNIST, CIFAR10, etc.)
  - $K = 1000$  classes (ImageNet)

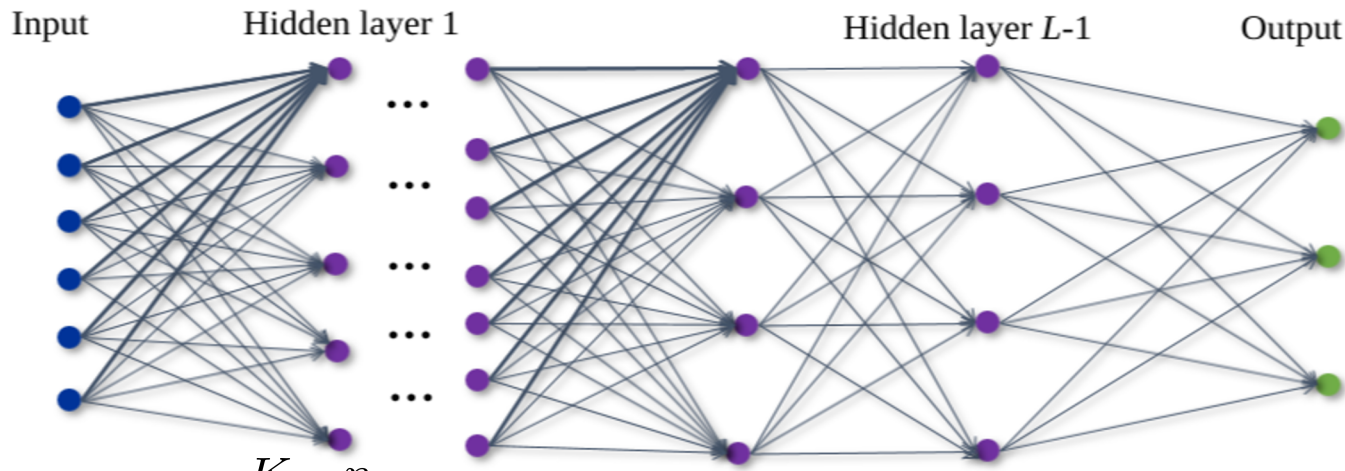


Data in the input space

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

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

# Understanding Deep Neural Networks



$$\min_{\Theta} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}_{\text{CE}}(\psi_{\Theta}(\mathbf{x}_{k,i}), \mathbf{y}_k) + \lambda \|\Theta\|_F^2$$

$\mathbf{x}_{k,i}$ :  $i$ -th input in the  $k$ -th class

$\mathbf{y}_k$ : One-hot vector for the  $k$ -th class

# Mysteries in Deep Learning

- **Architecture design** (before training):

- Feature dimensionality
- Network width and depth
- Activation functions

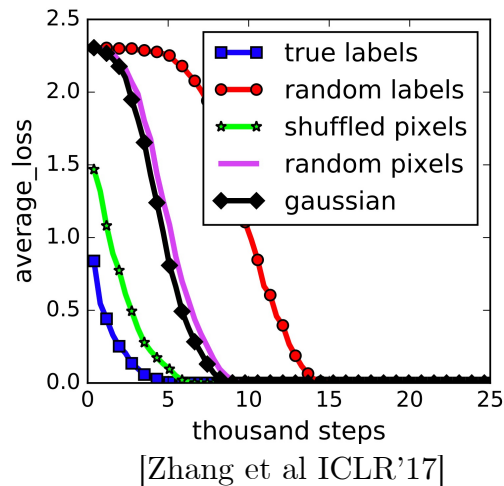
- **Optimization** (during training):

- Choices of loss functions
- Optimization algorithms, normalization

- **Properties of learned network**

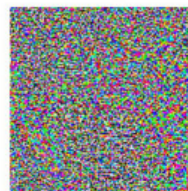
(after training):

- Generalization
- Robustness



$x$   
"panda"  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$   
"nematode"  
8.2% confidence

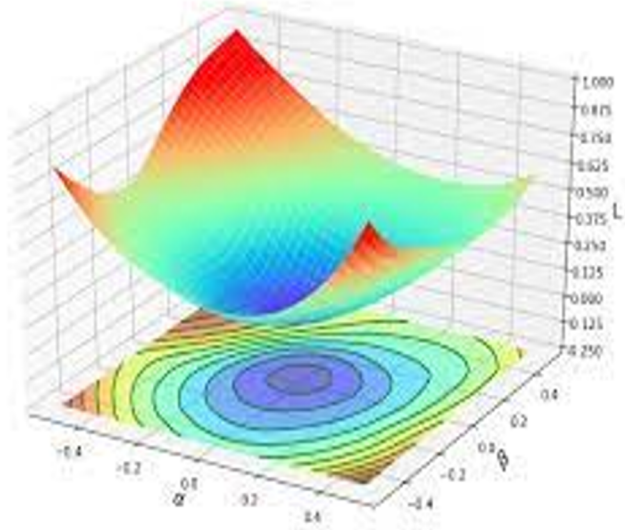
=



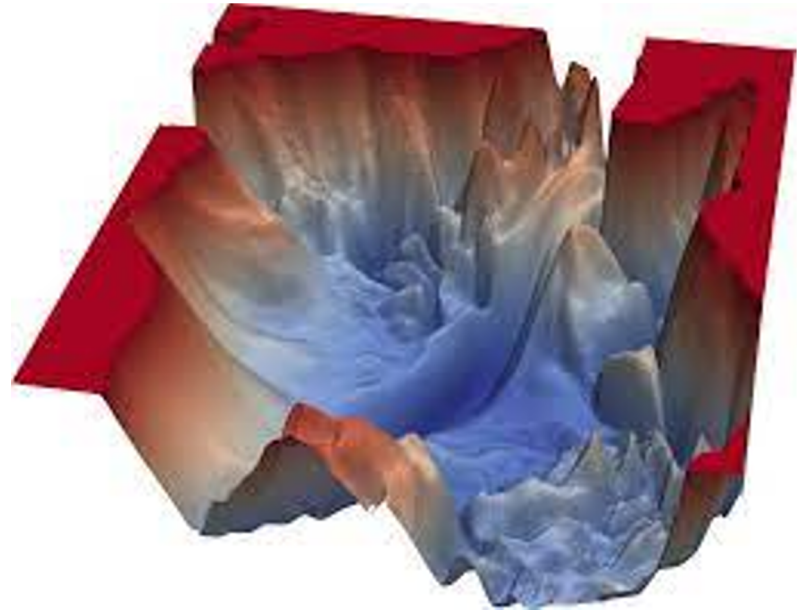
$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
"gibbon"  
99.3% confidence

Goodfellow et al ICLR'15

# Fundamental Challenges: Optimization



Landscape in **Classical** Optimization  
(abundant algorithms & theory)



Landscape of **Modern** Deep Neural Networks  
Credited to [Li'17]

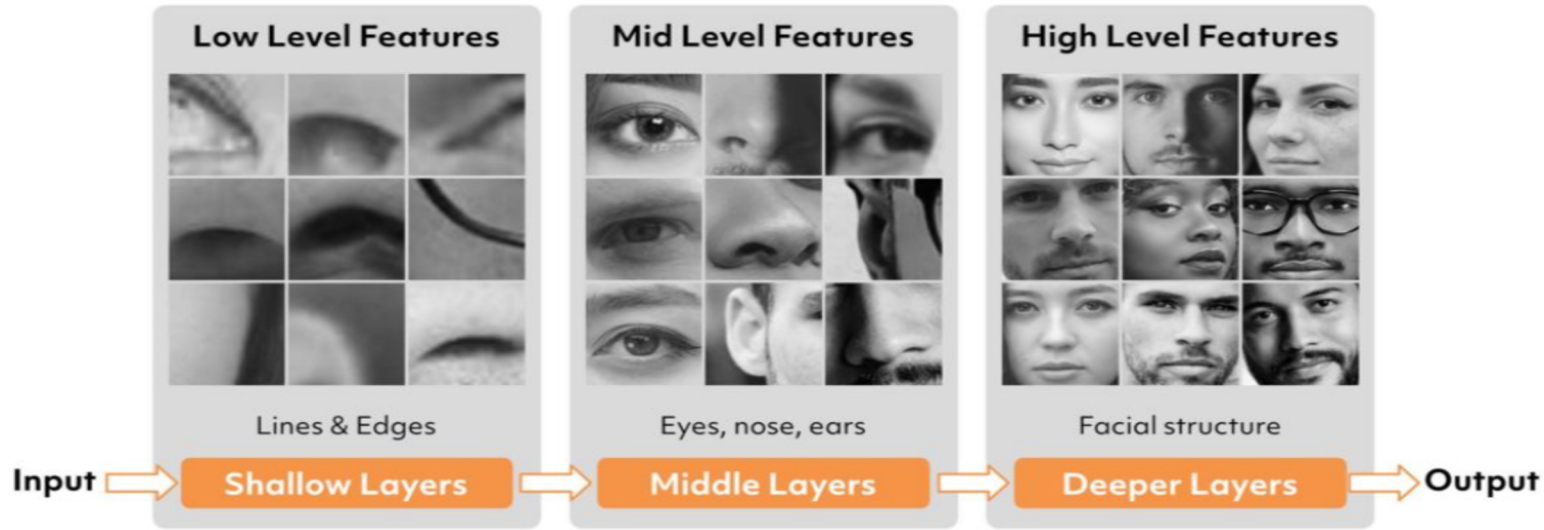
# Optimization: Existing Results

Existing analysis are based on various **simplifications**:

- **Go Linear:** deep linear networks [Kawaguchi'16], deep matrix factorizations [Arora'19], etc.
- **Go Shallow:** Two-layer neural networks [Safran'18, Liang'18], etc.
- **Go Wide:** Neural tangent kernels [Jacot'18, Allen-Zhu'18, Du'19], mean-field analysis [Mei'19, Sirignano'19], etc.

Most of results *hardly* provide much insights for **practical** neural networks.

# Features – What NNs (Conceptually) Designed to Learn



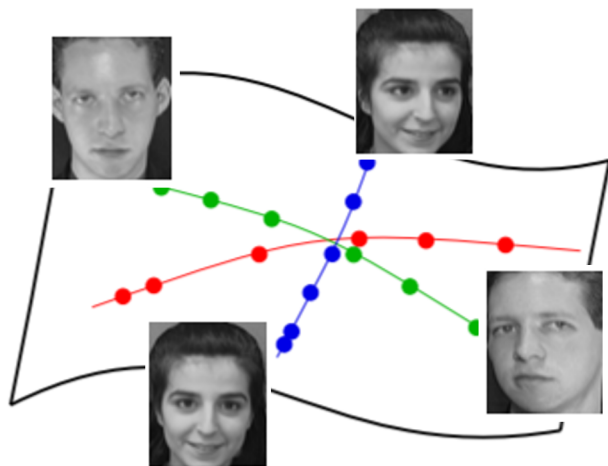
*Wishful Design: NNs learn rich feature representations across different levels?*

# Neural Collapse in Classification

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

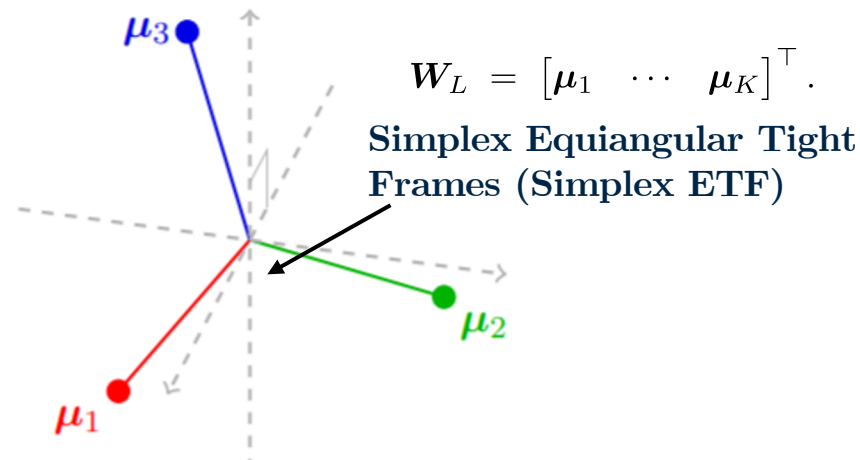
Last-layer classifier

Last-layer feature



Data in the Input Space

$\phi_{\theta}(\cdot)$



Neural Collapse  
in the Feature Space



# Neural Collapse in Classification

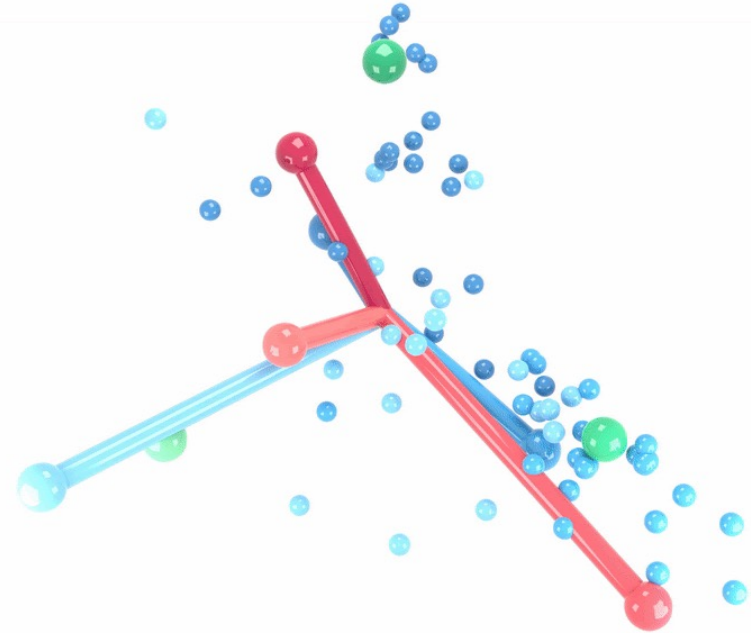
## Prevalence of neural collapse during the terminal phase of deep learning training

 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)



# Neural Collapse in Classification

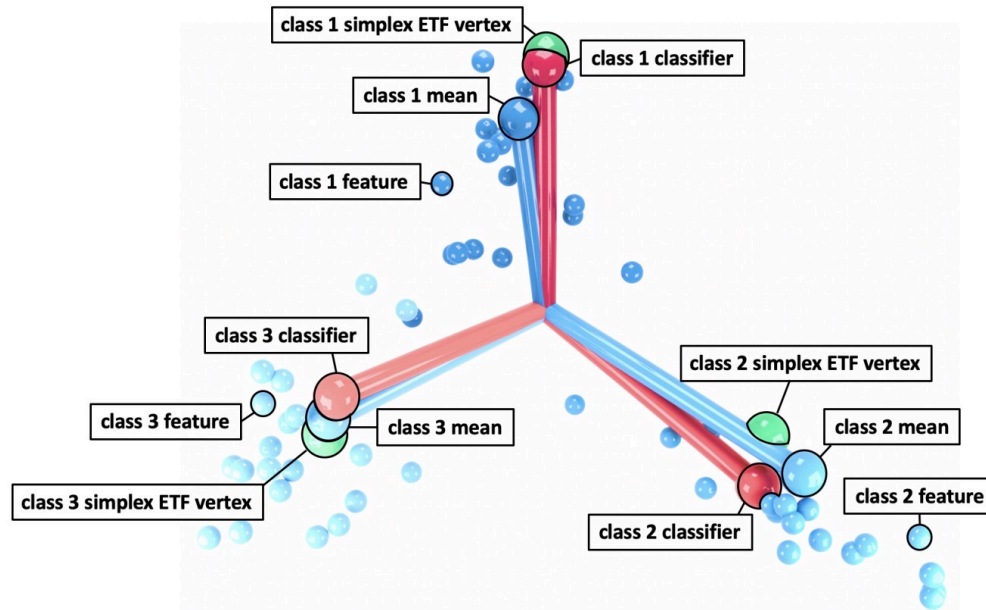


Image credited to Han et al. "Neural Collapse Under MSE Loss: Proximity to and Dynamics on the Central Path"

- Reveals common outcome of network training **across a variety of architectures** (ResNet, VGG) and **dataset** (CIFAR, ImageNet)
- **Precise mathematical structures** within the features and classifier

# Neural Collapse: Symmetry and Structures

Balanced training dataset with  $n = n_1 = n_2 = \dots = n_K$ , and

$$W := W_L, \quad H := [h_{1,1} \quad \dots \quad h_{K,n}].$$

Neural Collapse (NC) means that

- 1) *Within-Class Variability Collapse on H*: features of each class collapse to class-mean with **zero** variability;

$$h_{k,i} \rightarrow \bar{h}_k, \quad \forall k \in [K], i \in [n].$$

- 2) *Convergence to Simplex ETF on H*: the class means are **linearly separable**, and **maximally distant**;

$$M^\top M = \frac{K}{K-1} \left( \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right), \quad M = \alpha U \bar{H}$$

# Neural Collapse: Symmetry and Structures

Balanced training dataset with  $n = n_1 = n_2 = \dots = n_K$ , and

$$\mathbf{W} := \mathbf{W}_L, \quad \mathbf{H} := [\mathbf{h}_{1,1} \quad \dots \quad \mathbf{h}_{K,n}].$$

Neural Collapse (NC) means that

3) *Convergence to Self-Duality* ( $\mathbf{W}, \mathbf{H}$ ): the last-layer classifiers are **perfected matched** with the class-means of features.

$$\mathbf{w}^k = \beta \bar{\mathbf{h}}_k, \quad \forall k \in [K].$$

4) *Simple Decision Rule* via Nearest Class-Center decision.

# Simplification: Unconstrained Features

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

↖ Last-layer classifier

← Last-layer feature

Treat  $\mathbf{H} = [\mathbf{h}_{1,1} \quad \cdots \quad \mathbf{h}_{K,n}]$  as a **free** optimization variable

# Simplification: Unconstrained Features

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

↖ Last-layer classifier

← Last-layer feature

Treat  $\mathbf{H} = [\mathbf{h}_{1,1} \quad \cdots \quad \mathbf{h}_{K,n}]$  as a **free** optimization variable

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

# Simplification: Unconstrained Features

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

Last-layer classifier

$\phi_{\theta}(\mathbf{x}) =: \mathbf{h}$  ← Last-layer feature

Treat  $\mathbf{H} = [\mathbf{h}_{1,1} \quad \cdots \quad \mathbf{h}_{K,n}]$  as a **free** optimization variable

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

- **Validity:** Modern networks are highly **overparameterized**, that can **approximate any point** in the feature space [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.

# Prior Work on Unconstrained Features

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

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

- [E et al'20] considers  $\min_{\mathbf{W}, \mathbf{H}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k), \text{ s. t. } \|\mathbf{W}\|_2 \leq 1, \|\mathbf{h}_{k,i}\|_2 \leq 1$
- [Fang et al'21] studies  $\min_{\mathbf{W}, \mathbf{H}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k), \text{ s. t. } \|\mathbf{W}\|_F^2 \leq C_W, \|\mathbf{H}\|_F^2 \leq C_H$
- These work show that any **global** solution has NC, but
  - What about local minima?
  - The constrain formulation are **not aligned with practice**
- [Mixon et al'21, Han et al'21] studies NC under the MSE loss

J. Lu and S. Steinerberger, Neural collapse with cross-entropy loss, 2020

W. E and S. Wojtowytsch, On the emergence of tetrahedral symmetry in the final and penultimate layers of neural network classifiers, 2020

D. Mixon, H. Parshall, J. Pi. Neural collapse with unconstrained features, 2020

C. Fang, H. He, Q. Long, W. Su, Layer-peeled model: Toward understanding well-trained deep neural networks, 2021

X. Han, V. Papayan, D. Donoho, Neural Collapse Under MSE Loss: Proximity to and Dynamics on the Central Path, 2021



# Our Main Theoretical Results

**Theorem (Informal)** Consider the nonconvex loss with unconstrained feature model with  $K < d$  and balanced data

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

- **(Global Optimality)** Any global solution  $(\mathbf{W}_*, \mathbf{H}_*)$  satisfies the NC properties (1-4).
- **(Benign Global Landscape)** The function has **no spurious** local minimizer and is a **strict saddle function**, with negative curvature for non-global critical point.

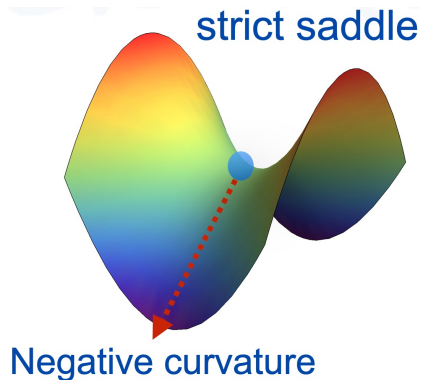
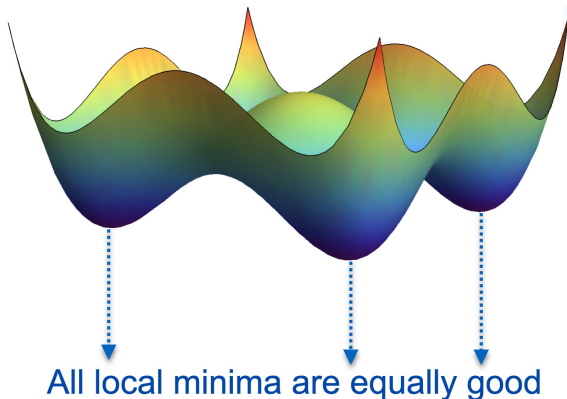
# Our Main Theoretical Results

**Theorem (Informal)** *Consider the nonconvex loss with unconstrained feature model with  $K < d$  and balanced data*

- *(Global Optimality) Any global solution  $(\mathbf{W}_*, \mathbf{H}_*)$  satisfies the NC properties (1-4).*
- *(Benign Global Landscape) The function has **no spurious** local minimizer and is a **strict saddle function**, with negative curvature for nonglobal critical point.*

*Message: deep networks always learn Neural Collapse features and classifiers, provably*

# Interpretations of Our Results



- **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.**

# Interpretations of Our Results

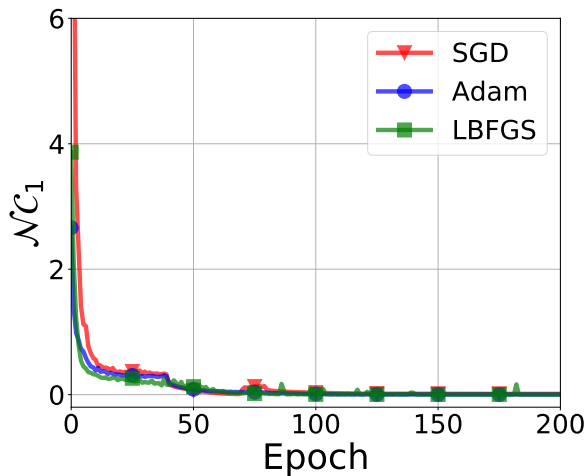
$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W} \mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

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]

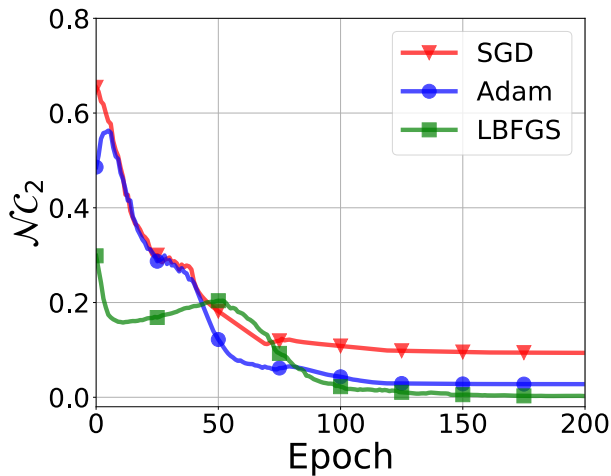
- **Difference in tasks:** classification training vs recovery
- **Difference in global solutions.**
- **Difference in loss functions, statistical properties:** cross-entropy vs least-squares; randomness or statistical properties of the sensing matrices

# Experiment: NC is Algorithm Independent

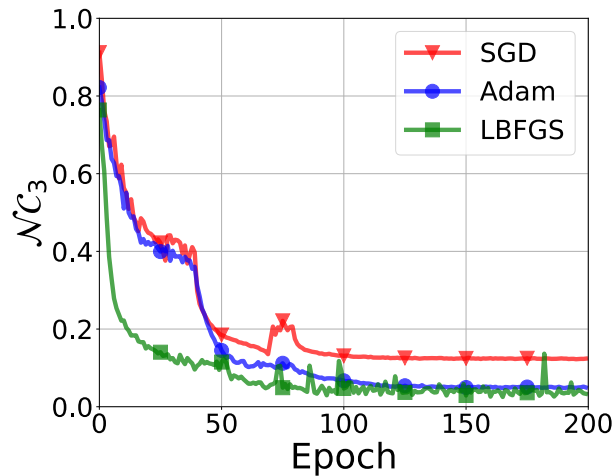
CIFAR-10 Dataset, ResNet18, with **different training algorithms**



Measure of Within-Class Variability



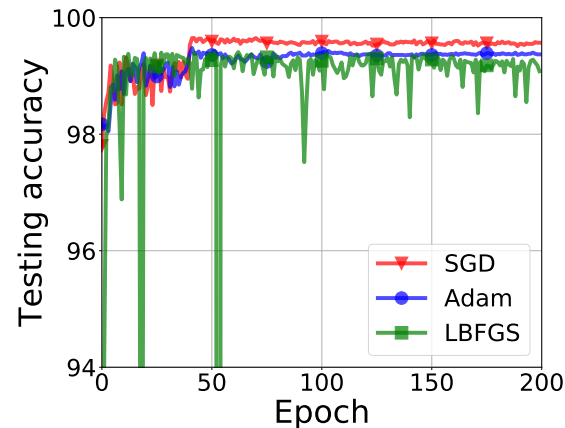
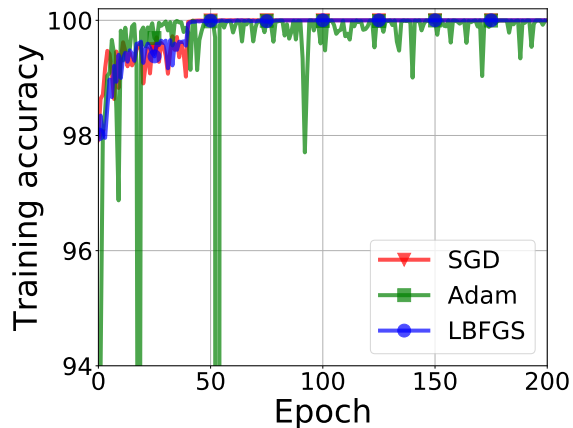
Measure of Between-Class Separation



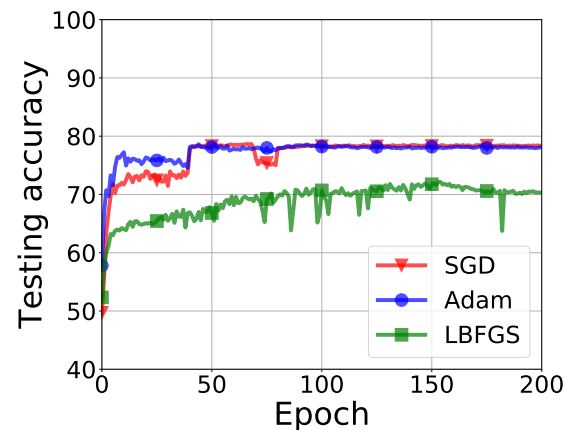
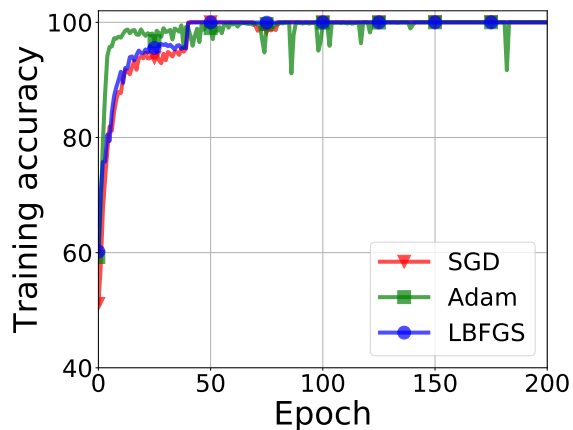
Measure of Self-Duality Collapse

# Generalization is Algorithm dependent

MINST

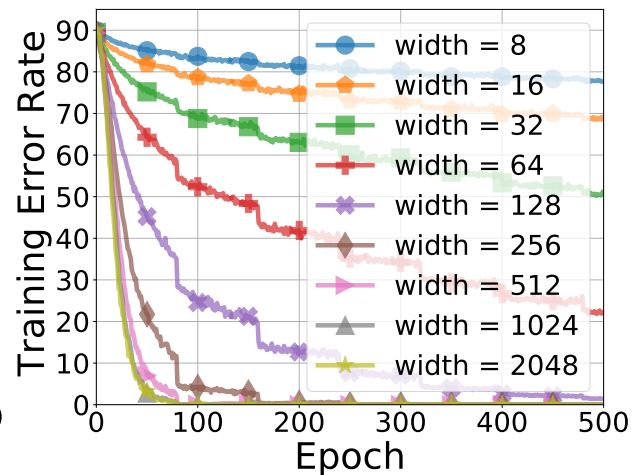
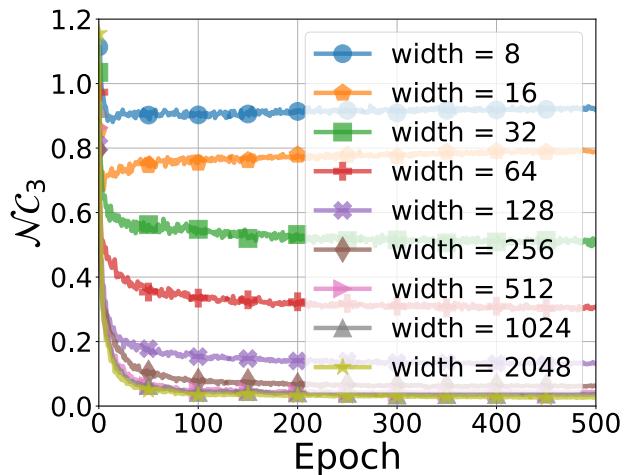
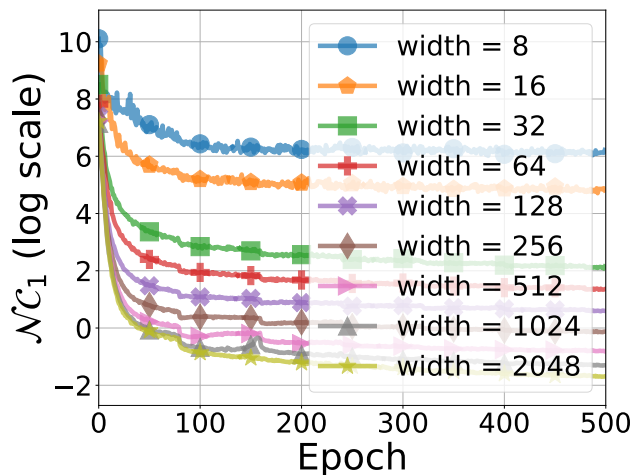


CIFAR-10



# Experiment: NC Occurs for Random Labels

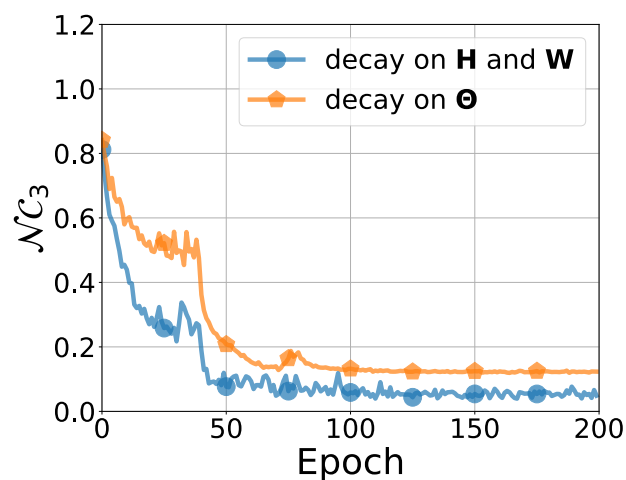
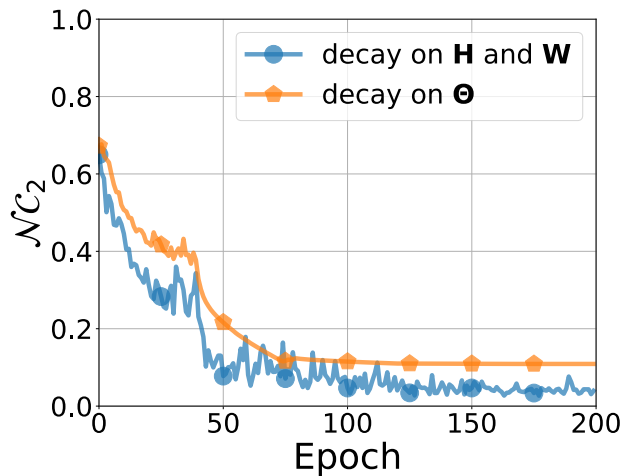
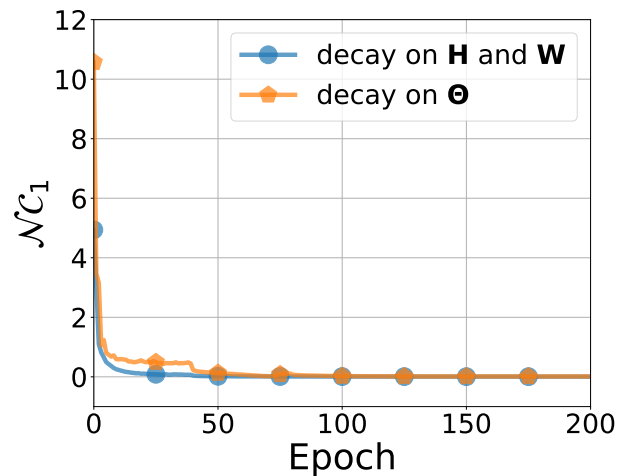
CIFAR-10 Dataset, MLP, random labels with varying network width



**Validity of Unconstrained Feature Model:** Learned last-layer features and classifiers seems to be **independent of input!**

# Experiment: NC with Different Weight Decays

CIFAR-10 Dataset, ResNet18, **different weight decay**



Measure of Within-Class Variability

Measure of Between-Class Separation

Measure of Self-Duality Collapse

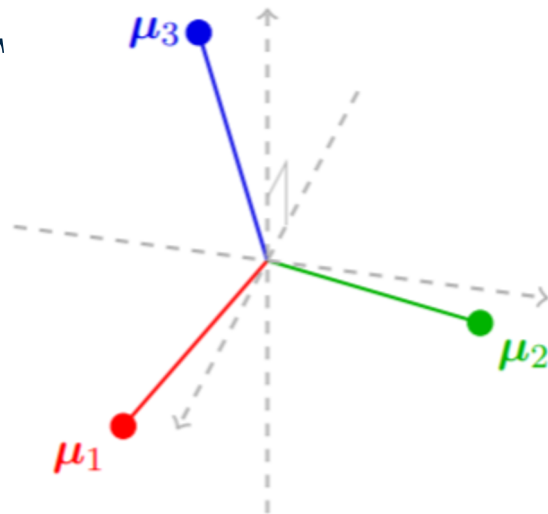
- Test Accuracy: 99.57% vs. 99.60% (MINST); 77.92% vs. 78.42% (CIFAR-10)
- Weight decay on the parameters (implicitly) regularizes the features



# Implications for Practical Network Training

**Observation:** For NC features, when  $K \leq d$  the best classifier is given by the Simplex ETF

$$\mathbf{W}_\star = [\boldsymbol{\mu}_1 \quad \cdots \quad \boldsymbol{\mu}_K]^\top.$$



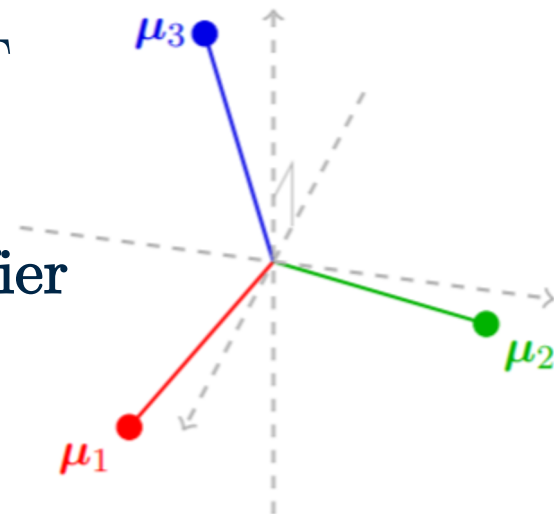
# Implications for Practical Network Training

**Observation:** For NC features, when  $K \leq d$  the best classifier is given by the Simplex ETF

$$W_{\star} = [\mu_1 \quad \cdots \quad \mu_K]^{\top}.$$

- **Implication 1: No need to learn the classifier**

- Just fix them as a Simplex ETF
- Save **8%**, **12%**, and **53%** parameters for ResNet50, DenseNet169, and ShuffleNet!



# Implications for Practical Network Training

**Observation:** For NC features, when  $K \leq d$  the best classifier is given by the Simplex ETF

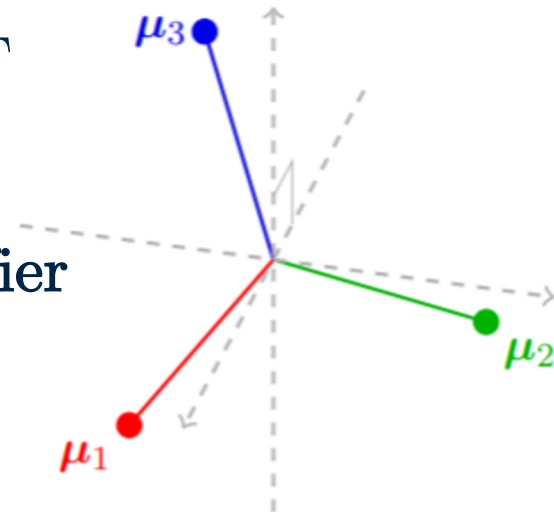
$$W_{\star} = [\mu_1 \quad \cdots \quad \mu_K]^{\top}.$$

- **Implication 1: No need to learn the classifier**

- ❑ Just fix them as a Simplex ETF
- ❑ Save **8%**, **12%**, and **53%** parameters for ResNet50, DenseNet169, and ShuffleNet!

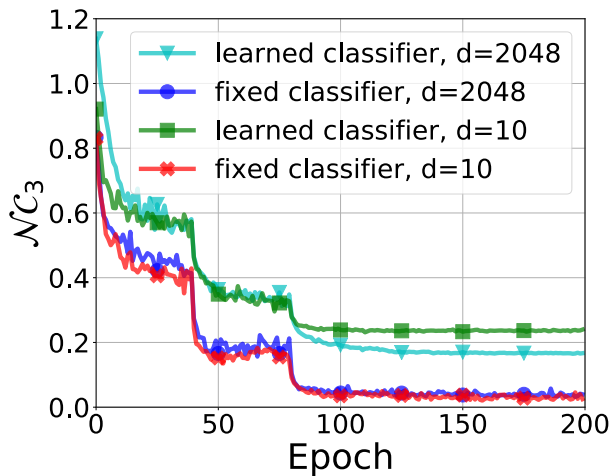
- **Implication 2: No need of large feature dimension  $d$**

- ❑ Just use feature dim  $d = \# \text{class } K$  (e.g.,  $d=10$  for CIFAR10)
- ❑ Further saves **21%** and **4.5%** parameters for ResNet18 and ResNet50!

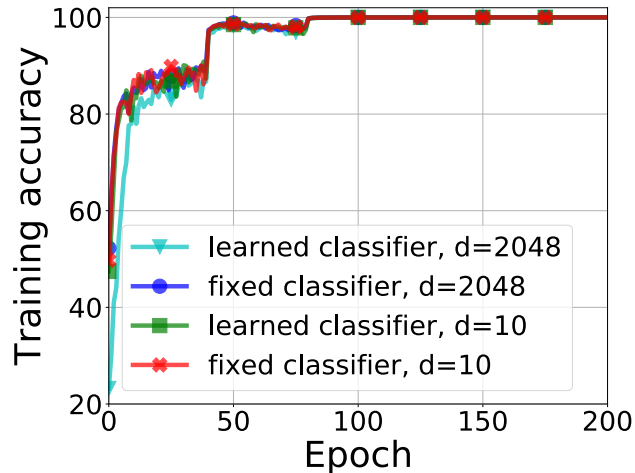


# Experiment: Fixed Classifier with $d = K$

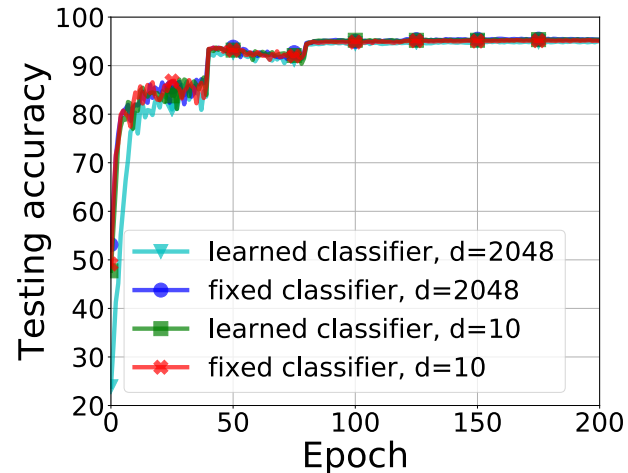
ResNet50, CIFAR10, Comparison of **Learned vs. Fixed Classifiers of  $W$**



Measure of Between-Class Separation



Training Accuracy



Testing Accuracy

Training with fixed last-layer classifiers achieves **on-par performance** with learned classifiers.

# Summary and Discussion

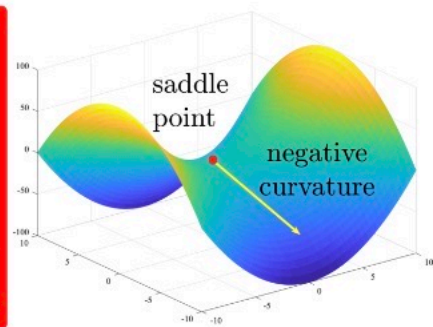
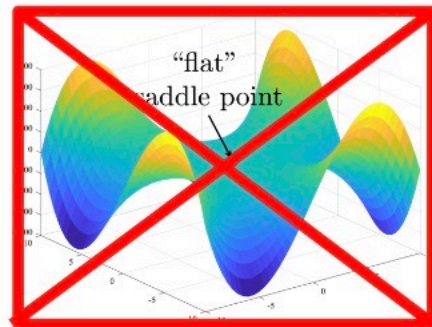
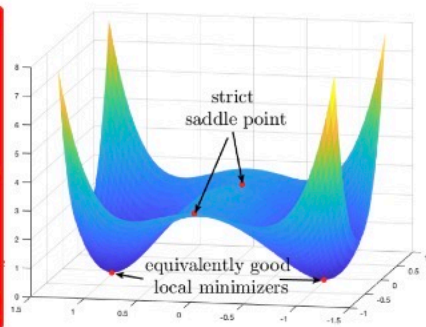
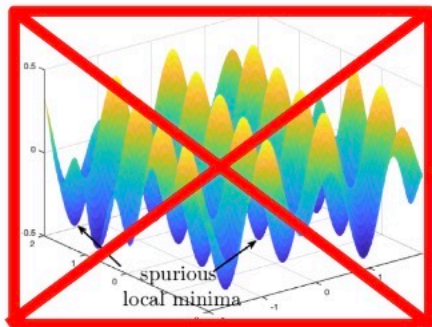
Z. Zhu, T. Ding, J. Zhou, X. Li, C. You, J. Sulam, and Q. Qu, [A Geometric Analysis of Neural Collapse with Unconstrained Features](#), *arXiv Preprint arXiv:2105.02375*, May 2021.

- Through landscape analysis under unconstrained feature model, we provide a **complete characterization of learned representation** of deep networks.
- The understandings of learned representations could shed lights on **generalization, robustness, and transferability**.

# Outline of this Talk

- Learning Shallow Representations:  
(Convolutional) Dictionary Learning
- Learning Deep Representations:  
Deep Neural Collapse
- **Conclusion & Discussion**

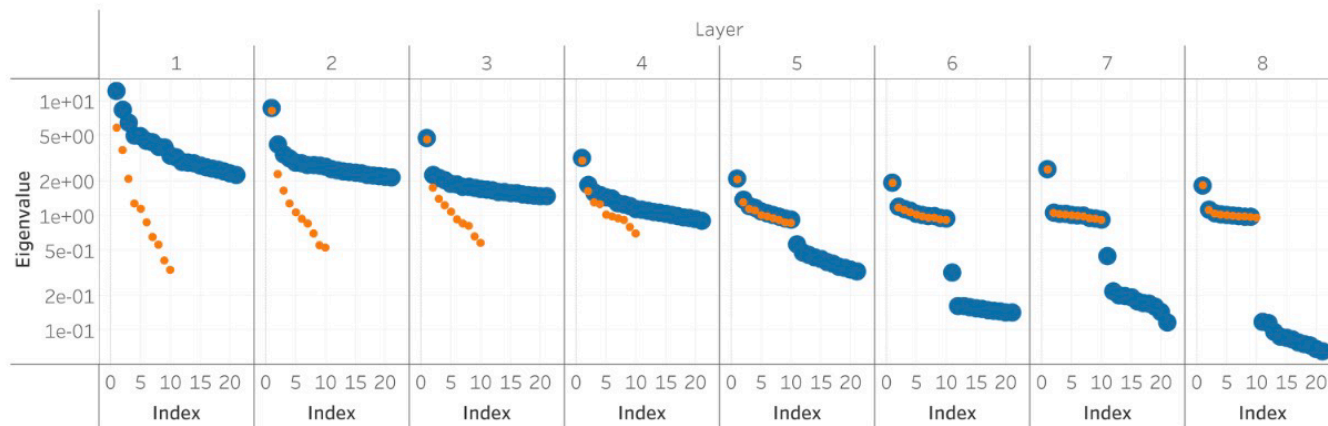
# Summary and Discussion



1. **Q. Qu**, Y. Zhai, X. Li, Y. Zhang, Z. Zhu, Analysis of optimization landscapes for overcomplete learning, *ICLR'20*, (oral, top 1.9%)
2. Y. Lau (\*), **Q. Qu**(\*), H. Kuo, P. Zhou, Y. Zhang, J. Wright, Short-and-sparse Deconvolution – A Geometric Approach, *ICLR'20*
3. Z. Zhu, T. Ding, J. Zhou, X. Li, C. You, J. Sulam, and **Q. Qu**, [Geometric Analysis of Neural Collapse with Unconstrained Features](#), *arXiv Preprint arXiv:2105.02375*, May 2021.

# Future Directions: Beyond Last-layer Features

- Study Deeper Networks
  - Fix the last layer classifier  $W$  as the Simplex ETF, and conduct NTK analysis for the learning dynamics of features  $H$ ?
  - Recursively study the features of each layer from output?



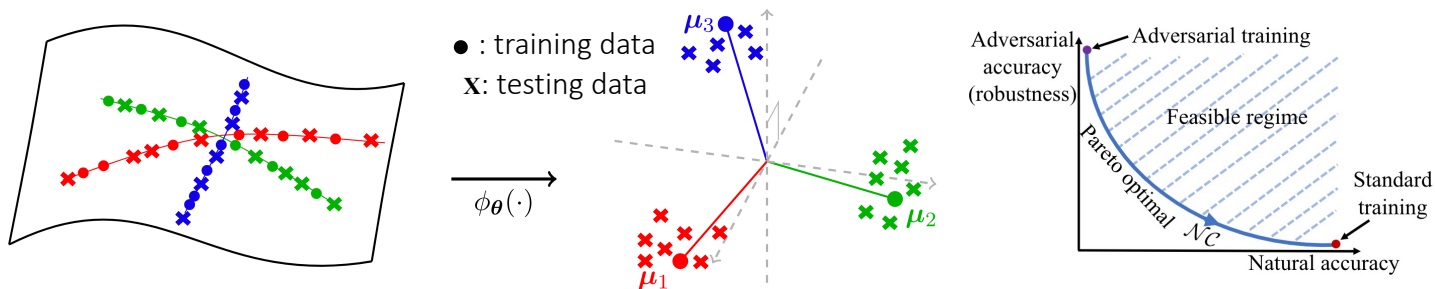
[Paypan'19]

Vardan Papyan. Traces of class/cross-class structure pervade deep learning spectra, JMLR'19.

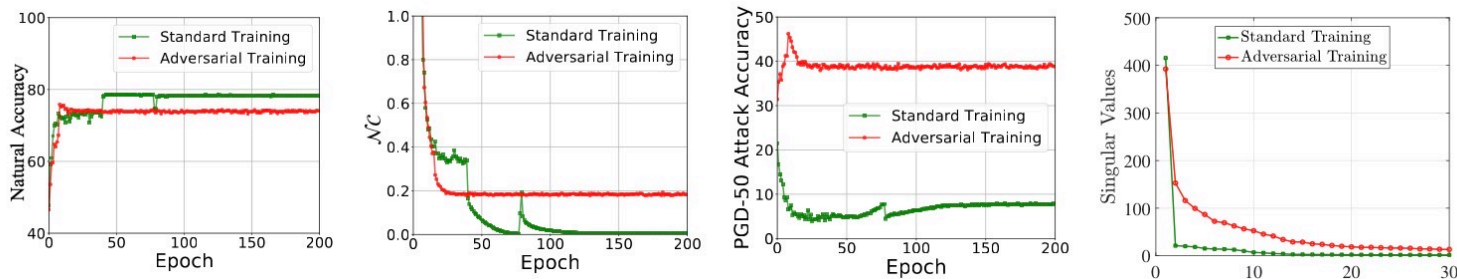


# Future Directions: Is NC a Blessing or Curse?

- Study generalization through the representation?

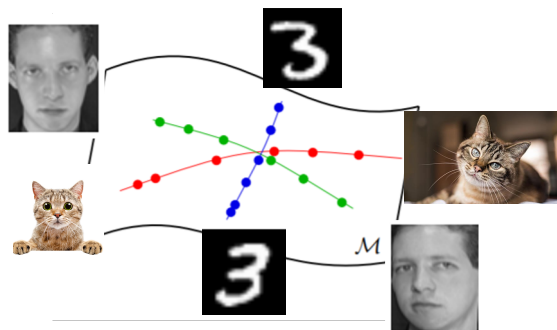


- Study tradeoff between accuracy and robustness via NC?

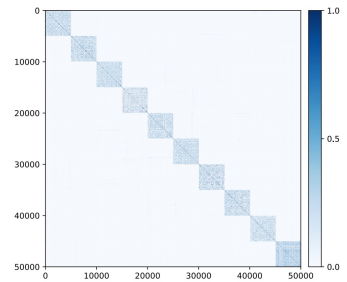
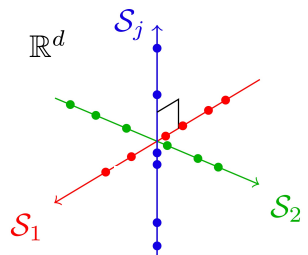


H. Zhang, Y. Yu, J. Jiao, E. Xing, L. Ghaoui, M. Jordan, Theoretically principled trade-off between robustness and accuracy, ICML2019.

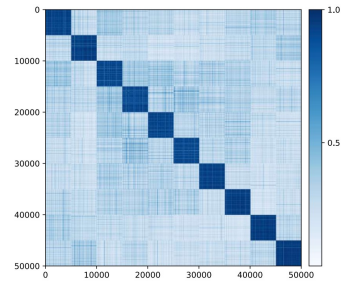
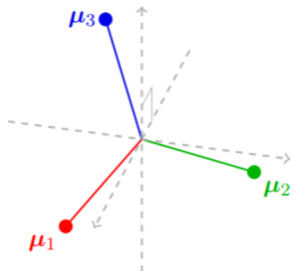
# Adaptive to the Intrinsic Data Structures



MCR<sup>2</sup>



Cross-entropy



- Can we learn diverse features that are adaptive to the intrinsic data structures?

# Acknowledgement



Tianyu Ding  
(Johns Hopkins)



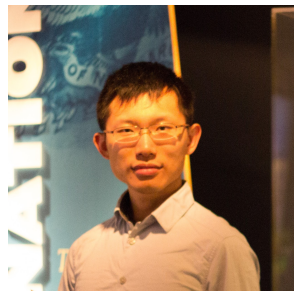
Xiao Li  
(U. Michigan)



Xiao Li  
(CUHK-Shen Zhen)



Jeremias Sulam  
(Johns Hopkins)



Chong You  
(Google Research)



Yuexiang Zhai  
(UC Berkeley)

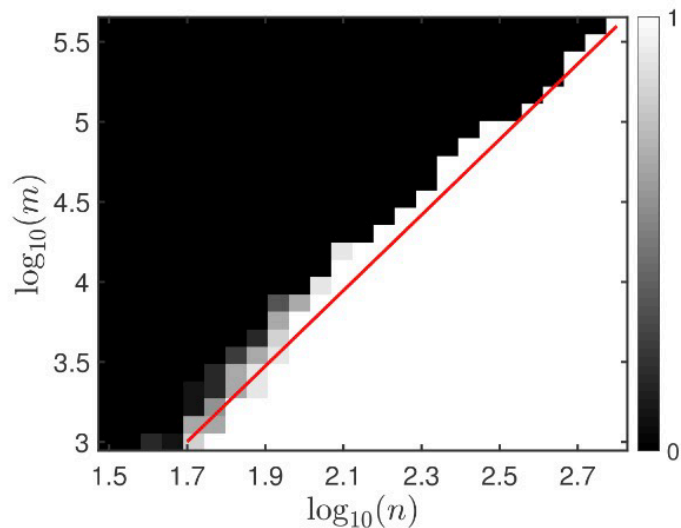


Zihui Zhu  
(University of Denver)

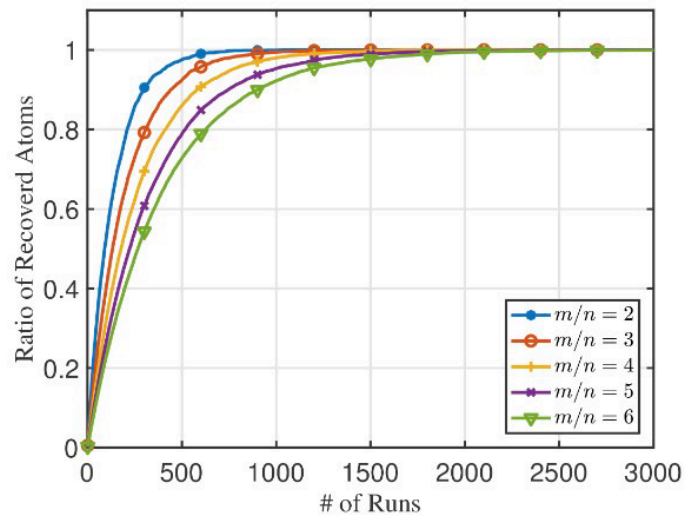


**Thank You!**

# Relationship to Dictionary Learning



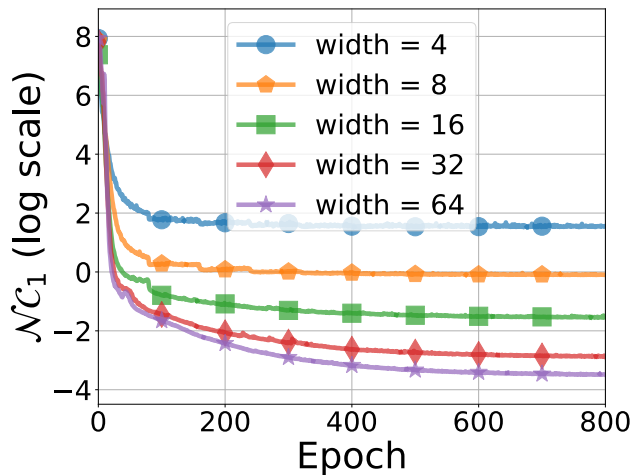
practice  $m < n^2$   
vs. theory  $m < Cn$



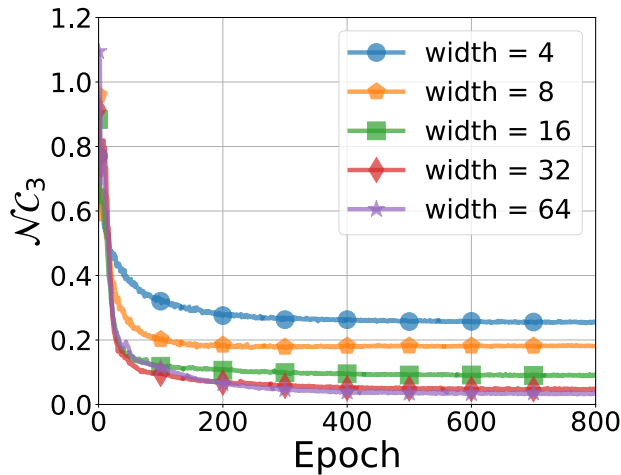
recover full  $\mathbf{A}_0$  via repeated  
independent trials

# Experiment: NC Occurs for Random Labels

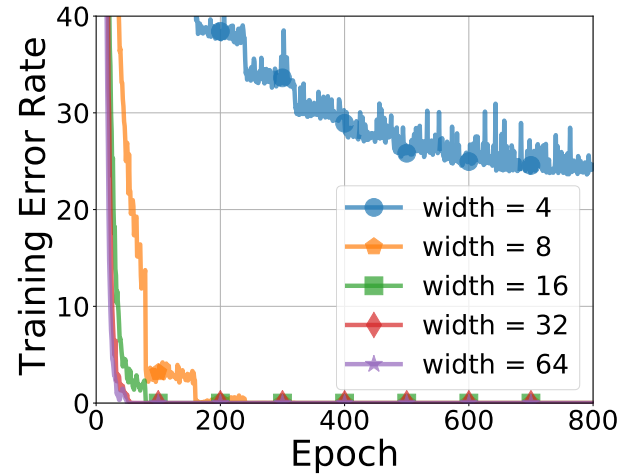
CIFAR-10 Dataset, ResNet18, random labels with varying network width



Measure of Within-Class Variability



Measure of Self-Duality Collapse



Training Error

**Validity of Unconstrained Feature Model:** Learned last-layer features and classifiers seems to be **independent of input!**

# Comparisons to MCR<sup>2</sup>

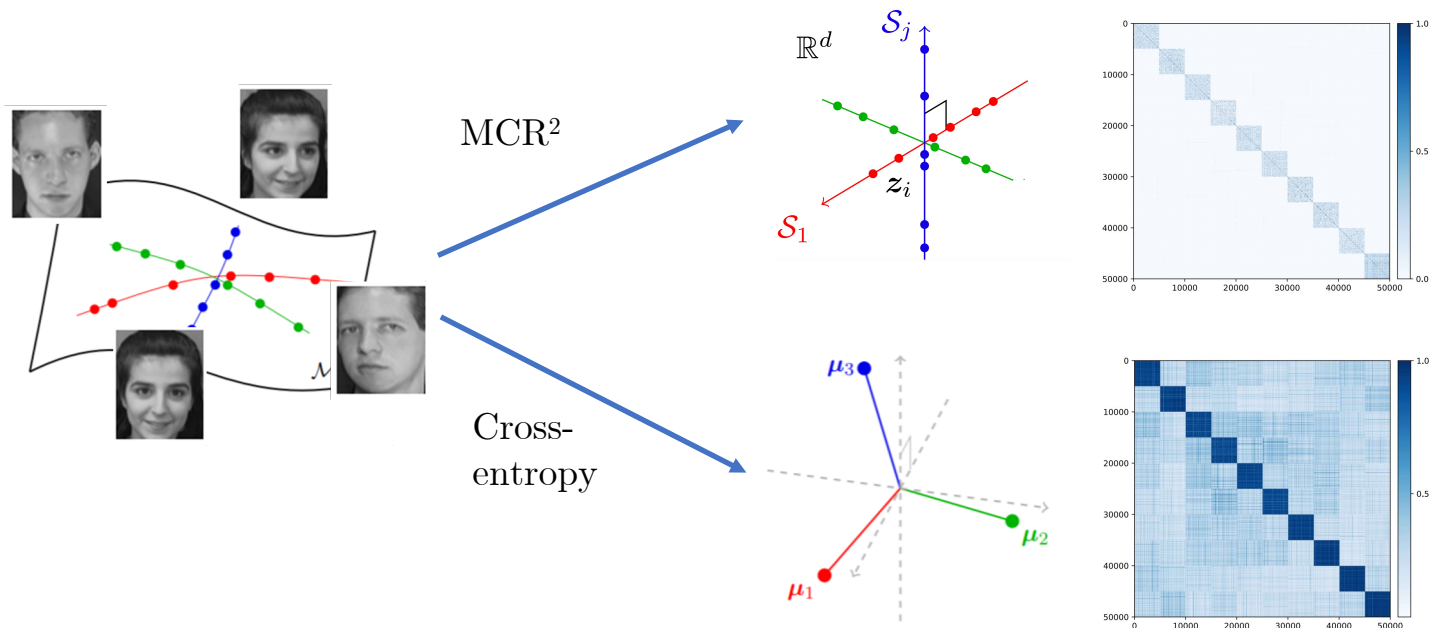
- [Yu et al, NeurIPS'20] learns not only **discriminative** but also **diverse** representations via maximizing the difference between the coding rate of all features and the average rate of features in the classes:

$$\Delta R(\mathbf{H}, \epsilon) = \underbrace{\frac{1}{2} \log \det(\mathbf{I} + \frac{d}{n\epsilon^2})}_R - \underbrace{\sum_{k=1}^K \frac{n_k}{2n} \log \det(\mathbf{I} + \frac{d}{n_k\epsilon^2} \mathbf{H}_k \mathbf{H}_k^\top)}_{R^c}$$

- $R$ : **expand** all features  $\mathbf{H}$  as **large** as possible.
- $R^c$ : **compress** all each class  $\mathbf{H}_k$  as **small** as possible.
- For balanced data, learned features  $\mathbf{H}_k$  span an **entire  $d/K$  subspace**, and the subspaces are orthogonal to each other.

# Comparisons to MCR<sup>2</sup>

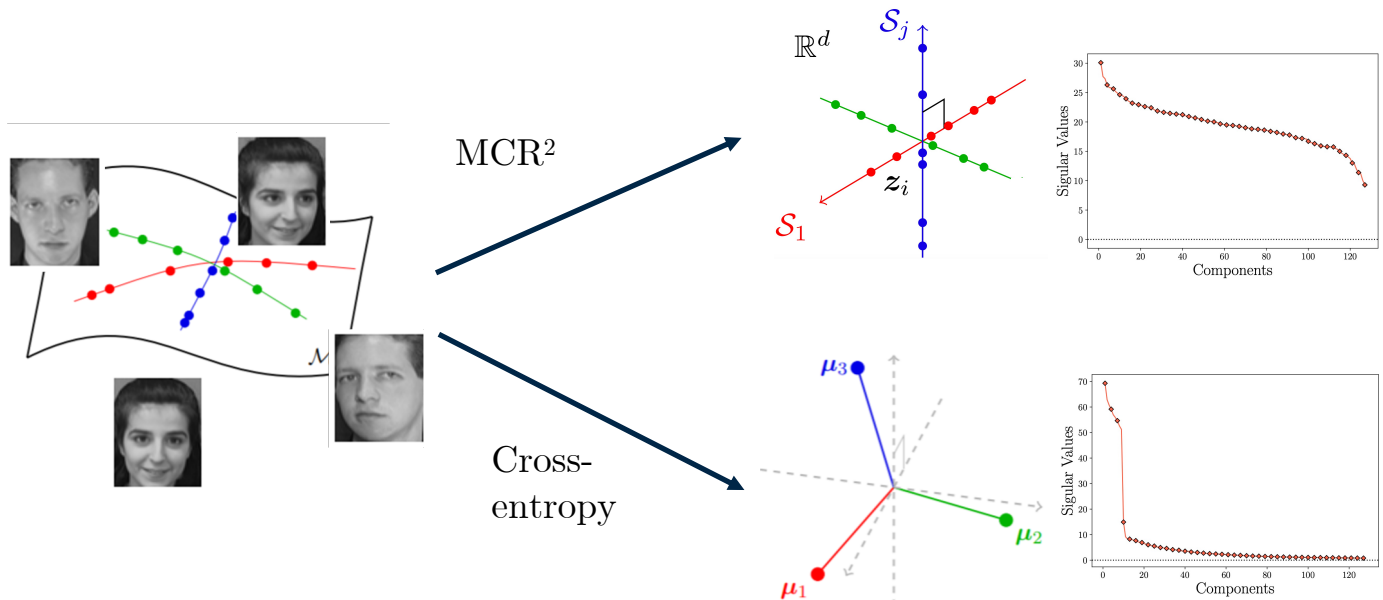
- [Yu et al, NeurIPS'20] learns not only **discriminative** but also **diverse** representations via maximizing the difference between the coding rate of all features and the average rate of features in the classes:





# Comparisons to MCR<sup>2</sup>

- [Yu et al, NeurIPS'20] learns not only **discriminative** but also **diverse** representations via maximizing the difference between the coding rate of all features and the average rate of features in the classes:



Y. Yu, K. Chan, C. You, C. Song, Y. Ma, Learning diverse and discriminative representations via the principle of maximal coding rate reduction, NeurIPS 20.

K. Chan, Y. Yu, C. You, H. Qi, J. Wright, and Y. Ma, ReduNet: A White-box Deep Network from the Principle of Maximizing Rate Reduction, 2021.