

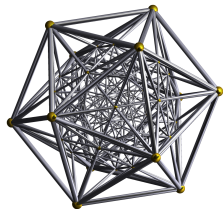
ACDL Summer Course 2023

## Lecture 4: Robust Learning of Overparameterized Networks via Low-Dimensional Models

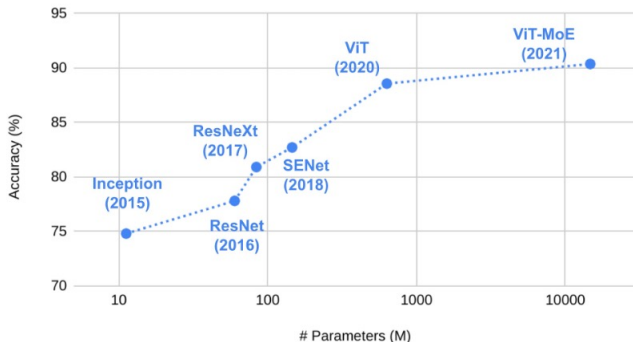
**Qing Qu**

EECS, University of Michigan

June 10th, 2023



# The Trend of Large Models...



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

~23 million    >>    ~1 million  
(# Parameters in ResNet-50)    (# Samples in ImageNet)

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

# The Curse of Overparameterization: Robustness

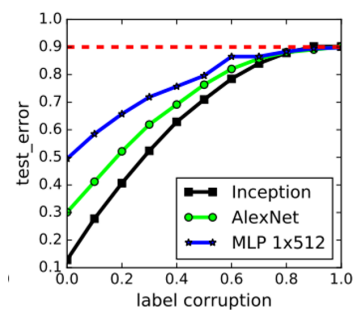


Figure: Label memorization.

# The Curse of Overparameterization: Robustness

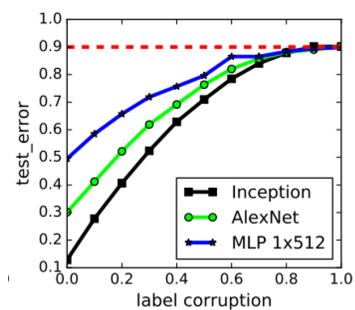


Figure: Label memorization.

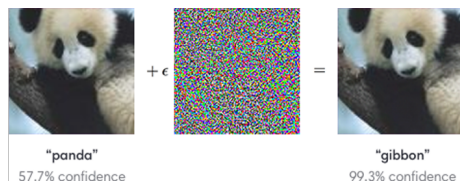


Figure: Adversarial attack.



# Outline

- ① Robust Classification under Noisy Labels
  - A Sparse Over-Parameterization Method
  - Theoretical Justification based on Simple Models
  - Experimental Results
- ② Extension to Robust Image Recovery
- ③ Conclusion

# Neural Collapse → Overfitting to Corruptions!

Label noise is common and often unavoidable

- Some proportion of the labels are incorrect (5-80%?)
- We don't know which labels are correct/incorrect

Inputs



cat

dog



ship

plane



bird

bird

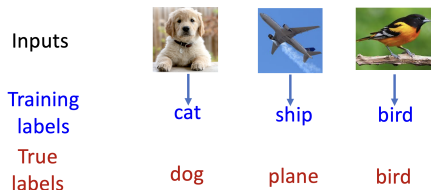
Training labels

True labels

# Neural Collapse → Overfitting to Corruptions!

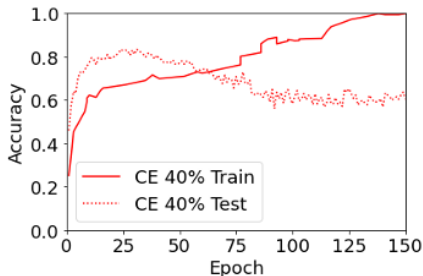
Label noise is common and often unavoidable

- Some proportion of the labels are incorrect (5-80%?)
- We don't know which labels are correct/incorrect



Neural Collapse always happens

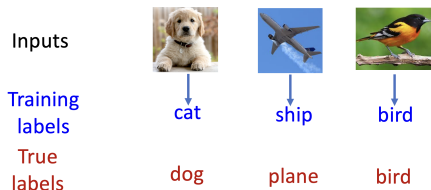
- Perfectly fits noisy labels (overfitting)
- Cannot predict well on new images



# Neural Collapse → Overfitting to Corruptions!

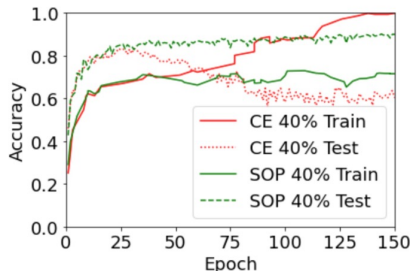
Label noise is common and often unavoidable

- Some proportion of the labels are incorrect (5-80%?)
- We don't know which labels are correct/incorrect



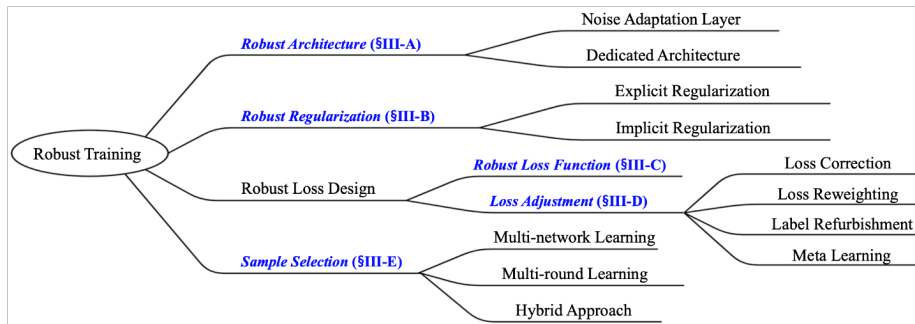
Neural Collapse always happens

- Perfectly fits noisy labels (overfitting)
- Cannot predict well on new images



# Prior Work on Robust Deep Learning for Noisy Labels


Various (heuristic or principled) methods have been proposed<sup>1</sup>



<sup>1</sup>Song et al., Learning from noisy labels with deep neural networks: A survey, IEEE TNLS, 2022.

# A Sparse Over-Parameterization (SOP) Method

**Observation:** Only a small fraction of the labels are corrupted, so that the label noise is **sparse**.



$$\begin{array}{c} \text{"cat"} \\ \text{"dog"} \\ \vdots \\ 0 \end{array} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = f(x; \theta) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$


Corrupted label      True label      Sparse noise

<sup>2</sup>Wright et al., Robust face recognition via sparse representation, TPAMI, 2008.

Candes et al., Robust principal component analysis? JACM, 2011.

# A Sparse Over-Parameterization (SOP) Method

**Observation:** Only a small fraction of the labels are corrupted, so that the label noise is **sparse**.



$$\begin{array}{c} \text{"cat"} \\ \text{"dog"} \end{array} \begin{bmatrix} y \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = f(x; \theta) \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + s \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

Corrupted label
True label
Sparse noise

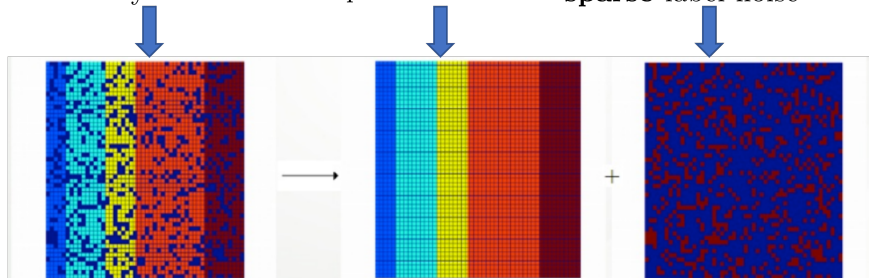
**Idea from the past:** we developed principled methods for dealing with sparse corruption in Compressive Sensing: Robust PCA<sup>2</sup>

<sup>2</sup>Wright et al., Robust face recognition via sparse representation, TPAMI, 2008.

Candes et al., Robust principal component analysis? JACM, 2011.

# A Sparse Over-Parameterization (SOP) Method

$$\underset{\text{noisy label}}{\mathbf{y}_i} = f\left(\underset{\text{input}}{\mathbf{x}_i}; \underset{\text{params.}}{\Theta^*}\right) + \underset{\text{sparse label noise}}{\mathbf{s}_i^*}$$



Exact Separation of Sparse Corruption with Incoherence between Data and Noise



# A Sparse Over-Parameterization (SOP) Method

**Our approach:**<sup>3</sup> minimize the distance between  $\mathbf{y}$  and  $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } s_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

<sup>3</sup>Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

# A Sparse Over-Parameterization (SOP) Method

**Our approach:**<sup>3</sup> minimize the distance between  $\mathbf{y}$  and  $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } \mathbf{s}_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

- Here the over-parameterization  $\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i$  introduces implicit algorithmic regularization [Vaskevicius et al.'19, Zhao et al.'19]

$$\text{variational form } \|\mathbf{s}\|_1 = \min_{\mathbf{s}=\mathbf{u}\odot\mathbf{u}-\mathbf{v}\odot\mathbf{v}} \frac{1}{2}(\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2)$$

<sup>3</sup>Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

# A Sparse Over-Parameterization (SOP) Method

**Our approach:**<sup>3</sup> minimize the distance between  $\mathbf{y}$  and  $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } \mathbf{s}_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

- Here the over-parameterization  $\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i$  introduces implicit algorithmic regularization [Vaskevicius et al.'19, Zhao et al.'19]

$$\text{variational form } \|\mathbf{s}\|_1 = \min_{\mathbf{s}=\mathbf{u}\odot\mathbf{u}-\mathbf{v}\odot\mathbf{v}} \frac{1}{2}(\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2)$$

- Why not use explicit regularization?

$$\min_{\boldsymbol{\theta}, \{\mathbf{s}_i\}} \frac{1}{N} \sum_{i=1}^N \underbrace{\mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \mathbf{s}_i, \mathbf{y}_i)}_{\rightarrow 0} + \underbrace{\lambda \|\mathbf{s}_i\|_1}_{\rightarrow 0}$$

<sup>3</sup>Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

# A Sparse Over-Parameterization (SOP) Method

**Our approach:**<sup>4</sup> minimize the distance between  $\mathbf{y}$  and  $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } s_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

**Training:** gradient descent with a discrepant learning rate:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \tau \frac{\partial}{\partial \boldsymbol{\theta}} L(\{\mathbf{u}_i, \mathbf{v}_i\}; \boldsymbol{\theta})$$

$$\mathbf{u}_i \leftarrow \mathbf{u}_i - \alpha \tau \frac{\partial}{\partial \boldsymbol{\theta}} L(\{\mathbf{u}_i, \mathbf{v}_i\}; \mathbf{u}_i)$$

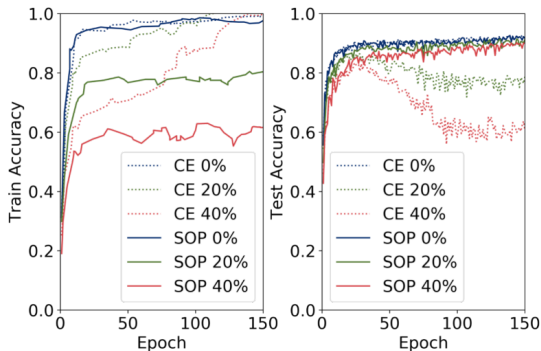
$$\mathbf{v}_i \leftarrow \mathbf{v}_i - \alpha \tau \frac{\partial}{\partial \boldsymbol{\theta}} L(\{\mathbf{u}_i, \mathbf{v}_i\}; \mathbf{v}_i)$$

Ideally, the implicit regularization drives the GD dynamics to the desired solution.

<sup>4</sup>Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

## A Sparse Over-Parameterization (SOP) Method

$\{0\%, 20\%, 40\%\}$  percent of labels for CIFAR-10 training data are randomly flipped uniformly to another class. Use ResNet34.



**Observation:** Compared to vanilla training, SOP does not overfit to wrong labels and obtain better generalization performance.

# Theoretical Justification on SOP

A simple model: assume  $f(\mathbf{x}; \boldsymbol{\theta})$  is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

# Theoretical Justification on SOP

A simple model: assume  $f(\mathbf{x}; \boldsymbol{\theta})$  is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

WLOG, assume  $f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta}_0 \rangle = 0$ . For  $N$  training samples,

$$\begin{bmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{x}_N; \boldsymbol{\theta}) \end{bmatrix} \approx \begin{bmatrix} \nabla f(\mathbf{x}_1; \boldsymbol{\theta}_0)^\top \\ \vdots \\ \nabla f(\mathbf{x}_N; \boldsymbol{\theta}_0)^\top \end{bmatrix} \boldsymbol{\theta} = \mathbf{J} \cdot \boldsymbol{\theta}$$

## Theoretical Justification on SOP

A simple model: assume  $f(\mathbf{x}; \boldsymbol{\theta})$  is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

WLOG, assume  $f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta}_0 \rangle = 0$ . For  $N$  training samples,

$$\begin{bmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{x}_N; \boldsymbol{\theta}) \end{bmatrix} \approx \begin{bmatrix} \nabla f(\mathbf{x}_1; \boldsymbol{\theta}_0)^\top \\ \vdots \\ \nabla f(\mathbf{x}_N; \boldsymbol{\theta}_0)^\top \end{bmatrix} \boldsymbol{\theta} = \mathbf{J} \cdot \boldsymbol{\theta}$$

This leads to the following corrupted observation problem

$$\mathbf{y} = \mathbf{J} \cdot \boldsymbol{\theta}_\star + \mathbf{s}_\star$$

where  $\boldsymbol{\theta}_\star$  is the underlying groundtruth parameter, and  $\mathbf{s}_\star$  is sparse.



## Theoretical Justification on SOP

We over-parameterize the sparse noise by  $\mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v}$  and solve

$$\min_{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}} g(\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}) = \frac{1}{2} \|\mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v} - \mathbf{y}\|_2^2$$

using gradient descent with *discrepant learning rates*

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mu \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t), \quad \begin{bmatrix} \mathbf{u}_{t+1} \\ \mathbf{v}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_t \\ \mathbf{v}_t \end{bmatrix} - \alpha \mu \begin{bmatrix} \nabla_{\mathbf{u}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \nabla_{\mathbf{v}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \end{bmatrix}$$

## Theoretical Justification on SOP

We over-parameterize the sparse noise by  $\mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v}$  and solve

$$\min_{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}} g(\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}) = \frac{1}{2} \|\mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v} - \mathbf{y}\|_2^2$$

using gradient descent with *discrepant learning rates*

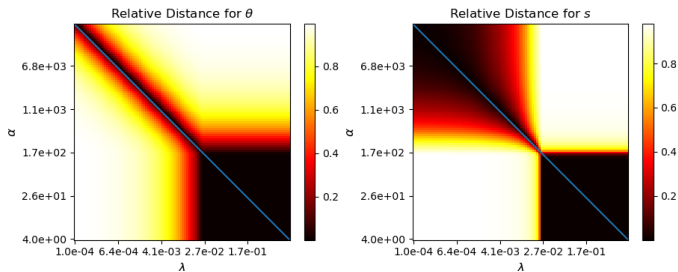
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mu \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t), \quad \begin{bmatrix} \mathbf{u}_{t+1} \\ \mathbf{v}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_t \\ \mathbf{v}_t \end{bmatrix} - \alpha \mu \begin{bmatrix} \nabla_{\mathbf{u}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \nabla_{\mathbf{v}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \end{bmatrix}$$

**Theorem (informal)** If gradient descent with infinitesimally small initialization and step size  $\mu$  converges to  $(\hat{\boldsymbol{\theta}}, \hat{\mathbf{u}}, \hat{\mathbf{v}})$ , then  $(\hat{\boldsymbol{\theta}}, \hat{\mathbf{u}} \odot \hat{\mathbf{u}} - \hat{\mathbf{v}} \odot \hat{\mathbf{v}})$  is an optimal solution to the following convex problem

$$\min_{\boldsymbol{\theta}, \mathbf{s}} \frac{1}{2} \|\boldsymbol{\theta}\|_2^2 + \lambda \|\mathbf{s}\|_1, \quad \text{s.t. } \mathbf{y} = \mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{s},$$

solving which exactly recovers  $(\boldsymbol{\theta}_*, \mathbf{s}_*)$  when  $\mathbf{J}$  is *incoherent* [Candes & Tao'05].

# Theoretical Justification on SOP



**Figure:** The SOP and the convex problem produce the same solutions with  $\alpha = -\frac{\log \gamma}{2\lambda}$ .

## Superior Performance with Training Efficiency

Methods	CIFAR-10				CIFAR-100			
	Symmetric			Asym	Symmetric			Asym
	20%	50%	80%	40%	20%	50%	80%	40%
CE	87.2	80.7	65.8	82.2	58.1	47.1	23.8	43.3
MixUp	93.5	87.9	72.3	-	69.9	57.3	33.6	-
DivideMix	96.1	94.6	93.2	93.4	77.1	74.6	60.2	72.1
ELR+	95.8	94.8	93.3	93.0	77.7	73.8	60.8	77.5
<b>SOP+</b>	<b>96.3</b>	<b>95.5</b>	<b>94.0</b>	<b>93.8</b>	<b>78.8</b>	<b>75.9</b>	<b>63.3</b>	<b>78.0</b>

**Ours**

## Superior Performance with Training Efficiency

Methods	CIFAR-10				CIFAR-100			
	Symmetric			Asym	Symmetric			Asym
	20%	50%	80%	40%	20%	50%	80%	40%
CE	87.2	80.7	65.8	82.2	58.1	47.1	23.8	43.3
MixUp	93.5	87.9	72.3	-	69.9	57.3	33.6	-
DivideMix	96.1	94.6	93.2	93.4	77.1	74.6	60.2	72.1
ELR+	95.8	94.8	93.3	93.0	77.7	73.8	60.8	77.5
<b>SOP+</b>	<b>96.3</b>	<b>95.5</b>	<b>94.0</b>	<b>93.8</b>	<b>78.8</b>	<b>75.9</b>	<b>63.3</b>	<b>78.0</b>

**Ours**






CE	Co-teaching+	DivideMix	ELR+	SOP (ours)	SOP+ (ours)
0.9h	4.4h	5.4h	2.3h	<b>1.0h</b>	<b>2.1h</b>

**Table: Comparison of total training time** in hours on CIFAR-10 with 50% symmetric label noise.

# SOP on CIFAR-10 with human annotated noisy labels

**CIFAR-10N:** provide CIFAR-10 with human annotated noisy labels<sup>5</sup>

Label Set	CIFAR-10N Aggregate	CIFAR-10N Random 1	CIFAR-10N Random 2	CIFAR-10N Random 3	CIFAR-10N Worst
<b>Noise Rate</b>	9.03%	17.23%	18.12%	17.64%	40.21%

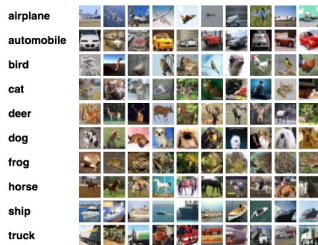
<sup>5</sup>Wei et al., Learning with noisy labels revisited: A study using real-world human annotations, ICLR 2022.     

# SOP on CIFAR-10 with human annotated noisy labels

**CIFAR-10N**: provide CIFAR-10 with human annotated noisy labels<sup>5</sup>

Label Set	CIFAR-10N Aggregate	CIFAR-10N Random 1	CIFAR-10N Random 2	CIFAR-10N Random 3	CIFAR-10N Worst
<b>Noise Rate</b>	9.03%	17.23%	18.12%	17.64%	40.21%

- Annotated by 747 independent workers
- Provide 5 noisy label sets for CIFAR-10 train images:
- **Random**  $i = 1, 2, 3$ : the  $i$ -th submitted label for each image;
- **Aggregate**: aggregation of three noisy labels by majority voting
- **Worst**: label set with the highest noise rate



<sup>5</sup>Wei et al., Learning with noisy labels revisited: A study using real-world human annotations, ICLR 2022. [↗](#) [↖](#) [↻](#) [🔍](#)

## New SOTA on CIFAR-10N

Method	CIFAR-10N					
	Clean	Aggregate	Random 1	Random 2	Random 3	Worst
CE (Standard)	92.92 ± 0.11	87.77 ± 0.38	85.02 ± 0.65	86.46 ± 1.79	85.16 ± 0.61	77.69 ± 1.55
Forward $T$ (Patrini et al., 2017)	93.02 ± 0.12	88.24 ± 0.22	86.88 ± 0.50	86.14 ± 0.24	87.04 ± 0.35	79.79 ± 0.46
Backward $T$ (Patrini et al., 2017)	93.10 ± 0.05	88.13 ± 0.29	87.14 ± 0.34	86.28 ± 0.80	86.86 ± 0.41	77.61 ± 1.05
GCE (Zhang & Sabuncu, 2018)	92.83 ± 0.16	87.85 ± 0.70	87.61 ± 0.28	87.70 ± 0.56	87.58 ± 0.29	80.66 ± 0.35
Co-teaching (Han et al., 2018)	93.35 ± 0.14	91.20 ± 0.13	90.33 ± 0.13	90.30 ± 0.17	90.15 ± 0.18	83.83 ± 0.13
To-teaching+ (Yu et al., 2019)	92.41 ± 0.20	90.61 ± 0.22	89.70 ± 0.27	89.47 ± 0.18	89.54 ± 0.22	83.26 ± 0.17
T-Revision (Xia et al., 2019)	93.35 ± 0.23	88.52 ± 0.17	88.33 ± 0.32	87.71 ± 1.02	87.79 ± 0.67	80.48 ± 1.20
Peer Loss (Liu & Guo, 2020)	93.99 ± 0.13	90.75 ± 0.25	89.06 ± 0.11	88.76 ± 0.19	88.57 ± 0.09	82.00 ± 0.60
ELR (Liu et al., 2020)	93.45 ± 0.65	92.38 ± 0.64	91.46 ± 0.38	91.61 ± 0.16	91.41 ± 0.44	83.58 ± 1.13
ELR+ (Liu et al., 2020)	<b>95.39 ± 0.05</b>	94.83 ± 0.10	94.43 ± 0.41	94.20 ± 0.24	94.34 ± 0.22	91.09 ± 1.60
Positive-LS (Lukasik et al., 2020)	94.77 ± 0.17	91.57 ± 0.07	89.80 ± 0.28	89.35 ± 0.33	89.82 ± 0.14	82.76 ± 0.53
F-Div (Wei & Liu, 2020)	94.88 ± 0.12	91.64 ± 0.34	89.70 ± 0.40	89.79 ± 0.12	89.55 ± 0.49	82.53 ± 0.52
Divide-Mix (Li et al., 2020)	<b>95.37 ± 0.14</b>	<b>95.01 ± 0.71</b>	<b>95.16 ± 0.19</b>	<b>95.23 ± 0.07</b>	<b>95.21 ± 0.14</b>	<b>92.56 ± 0.42</b>
Negative-LS (Wei et al., 2021)	<b>94.92 ± 0.25</b>	91.97 ± 0.46	90.29 ± 0.32	90.37 ± 0.12	90.13 ± 0.19	82.99 ± 0.36
JoCoR (Wei et al., 2020)	93.40 ± 0.24	91.44 ± 0.05	90.30 ± 0.20	90.21 ± 0.19	90.11 ± 0.21	83.37 ± 0.30
CORES <sup>2</sup> (Cheng et al., 2021)	93.43 ± 0.24	91.23 ± 0.11	89.66 ± 0.32	89.91 ± 0.45	89.79 ± 0.50	83.60 ± 0.53
CORES* (Cheng et al., 2021)	94.16 ± 0.11	<b>95.25 ± 0.09</b>	94.45 ± 0.14	94.88 ± 0.31	94.74 ± 0.03	91.66 ± 0.09
VolMinNet (Li et al., 2021)	92.14 ± 0.30	89.70 ± 0.21	88.30 ± 0.12	88.27 ± 0.09	88.19 ± 0.41	80.53 ± 0.20
CAL (Zhu et al., 2021a)	94.50 ± 0.31	91.97 ± 0.32	90.93 ± 0.31	90.75 ± 0.30	90.74 ± 0.24	85.36 ± 0.16
PES (Semi) (Bai et al., 2021)	94.76 ± 0.2	94.66 ± 0.18	<b>95.06 ± 0.15</b>	<b>95.19 ± 0.23</b>	<b>95.22 ± 0.13</b>	<b>92.68 ± 0.22</b>
SOP (Liu et al., 2022)	N/A	<b>95.61 ± 0.13</b>	<b>95.28 ± 0.13</b>	<b>95.31 ± 0.10</b>	<b>95.39 ± 0.11</b>	<b>93.24 ± 0.21</b>

Two-network based

Two-network based

Semi-supervised  
OursSparse modeling gives super performance again label noise<sup>6</sup><sup>6</sup>Wei et al., Learning with noisy labels revisited: A study using real-world human annotations, ICLR 2022.



# Outline

- ① Robust Classification under Noisy Labels
  - A Sparse Over-Parameterization Method
  - Theoretical Justification based on Simple Models
  - Experimental Results
- ② Extension to Robust Image Recovery
- ③ Conclusion

## Deep Image Prior<sup>7</sup>

- **Goal:** given a corrupted image  $\mathbf{y} = \mathbf{x}_* + \mathbf{s}$ , recover the clean image  $\mathbf{x}_*$  from the noisy observation

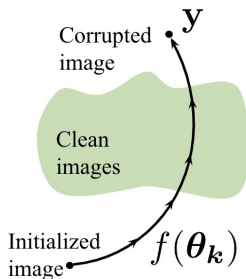
---

<sup>7</sup>Ulyanov D, Vedaldi A, Lempitsky V. Deep image prior[J]. International Journal of Computer Vision, 2020, 128(7).

## Deep Image Prior<sup>7</sup>

- **Goal:** given a corrupted image  $\mathbf{y} = \mathbf{x}_* + \mathbf{s}$ , recover the clean image  $\mathbf{x}_*$  from the noisy observation
- **Idea:** using a deep network  $f(\boldsymbol{\theta})$  to fit the observation  $\mathbf{y}$ :

$$\min_{\boldsymbol{\theta}} \ell(\underset{\text{corrupted image}}{\mathbf{y}}, \underset{\text{recovered image}}{f(\boldsymbol{\theta})})$$

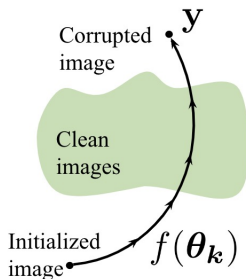


<sup>7</sup>Ulyanov D, Vedaldi A, Lempitsky V. Deep image prior[J]. International Journal of Computer Vision, 2020, 128(7).

## Deep Image Prior<sup>7</sup>

- **Goal:** given a corrupted image  $\mathbf{y} = \mathbf{x}_* + \mathbf{s}$ , recover the clean image  $\mathbf{x}_*$  from the noisy observation
- **Idea:** using a deep network  $f(\boldsymbol{\theta})$  to fit the observation  $\mathbf{y}$ :

$$\min_{\boldsymbol{\theta}} \ell(\underset{\text{corrupted image}}{\mathbf{y}}, \underset{\text{recovered image}}{f(\boldsymbol{\theta})})$$



- **Early stopping:** As the network is highly **overparameterized**, early stopping is needed.

<sup>7</sup>Ulyanov D, Vedaldi A, Lempitsky V. Deep image prior[J]. International Journal of Computer Vision, 2020, 128(7).

## A Case Study: Robust Image Recovery with Sparse Noise



$$\min_{\theta} \|$$

$$\underbrace{y}_{\text{corrupted input}}$$

$$-$$

$$\underbrace{f(\theta)}_{\text{recovered image}}$$

$$\|_1$$

sparse corruption

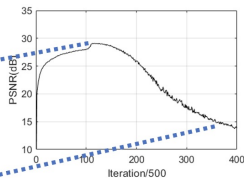
# A Case Study: Robust Image Recovery with Sparse Noise



Early termination solution  
(impractical!)



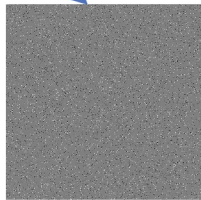
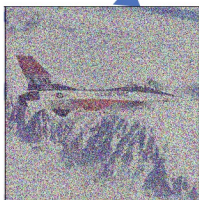
Global solution:  $f(\theta) \approx y$   
(overfitting!)



# Robust Recovery without Overfitting?

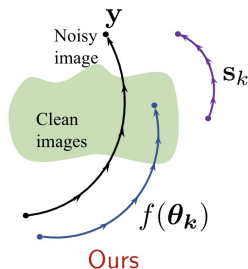
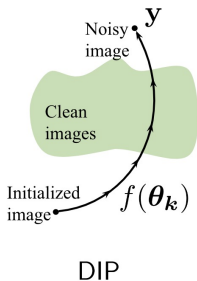
**Method:** sparse (double) overparameterization:<sup>8</sup>

$$\min_{\theta, g, h} \left\| \mathbf{y} - \left( f(\theta) + \underbrace{g \odot g - h \odot h}_{\text{noise}} \right) \right\|_2^2$$



<sup>8</sup>You C, Zhu Z, Qu Q, Ma Y. Robust recovery via implicit bias of discrepant learning rates for double over-parameterization. NeurIPS'20.

# Sparse Overparameterization Method



- **Optimization:** gradient descent with discrepant learning rate:

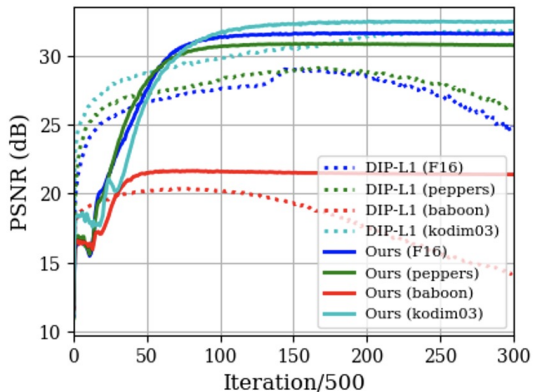
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \tau \frac{\partial}{\partial \boldsymbol{\theta}} \ell(\{\mathbf{u}, \mathbf{v}\}; \boldsymbol{\theta})$$

$$\mathbf{u} \leftarrow \mathbf{u} - \alpha \tau \frac{\partial}{\partial \boldsymbol{\theta}} \ell(\{\mathbf{u}, \mathbf{v}\}; \mathbf{u})$$

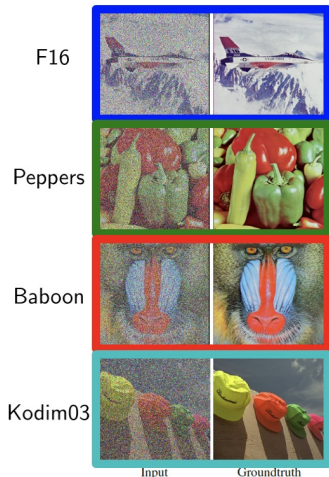
$$\mathbf{v} \leftarrow \mathbf{v} - \alpha \tau \frac{\partial}{\partial \boldsymbol{\theta}} \ell(\{\mathbf{u}, \mathbf{v}\}; \mathbf{v})$$



## Experiments on Real Images



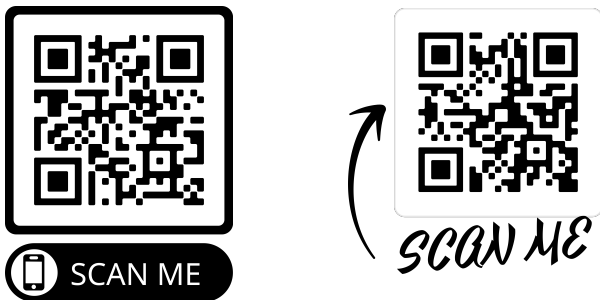
**No early stop, no parameter tuning!**



# Outline

- ① Robust Classification under Noisy Labels
  - A Sparse Over-Parameterization Method
  - Theoretical Justification based on Simple Models
  - Experimental Results
- ② Extension to Robust Image Recovery
- ③ Conclusion

## References



- 1 Liu S, Zhu Z, Qu Q, You C, Robust Training under Label Noise by Over-parameterization, ICML'22.
- 2 You C, Zhu Z, Qu Q, Ma Y. Robust recovery via implicit bias of discrepant learning rates for double over-parameterization. NeurIPS'20.
- 3 Li X, Yaras C, An Y, Ravishankar S, and Qu Q. Data-Free Deep Image Recovery from Sparsely Corrupted Incomplete Measurements, In preparation, 2023.

## Conclusion and Coming Attractions

**Take-home Message:** We can achieve better robustness in learning our overparameterized deep models by exploiting the low-dimensional structures in the data and network.

**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.**

