

# Deep Learning with Noisy Supervision

Dr. Bo Han

Trustworthy Machine Learning Group

Computer Science Department, Hong Kong Baptist University

<https://bhanml.github.io/>



DEPARTMENT OF  
COMPUTER SCIENCE  
HONG KONG BAPTIST UNIVERSITY  
香港浸會大學計算機科學系



# Overview of This Tutorial

- Part I: Why and What Noisy Labels
- Part II: Current Progress and Tutorial Perspectives
- Part III: Training Perspective
- Part IV: Data Perspective
- Part V: Regularization Perspective
- Part VI: Future Directions

# Part I: Why Noisy Labels

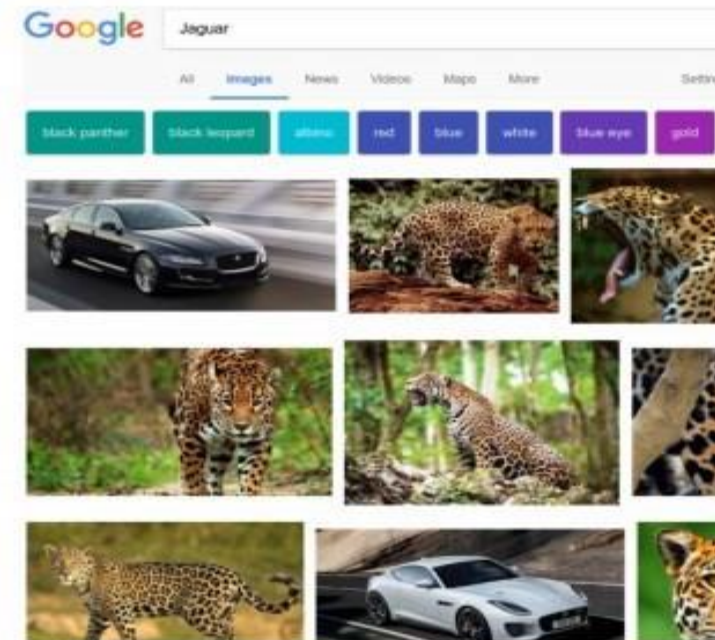
## Active label collection



In crowdsourcing,  
labels are from **non-experts**

(Credit to Amazon)

## Passive label collection



In web search,  
labels are from **users' clicks**

(Credit to Google)

# Why Noisy Labels

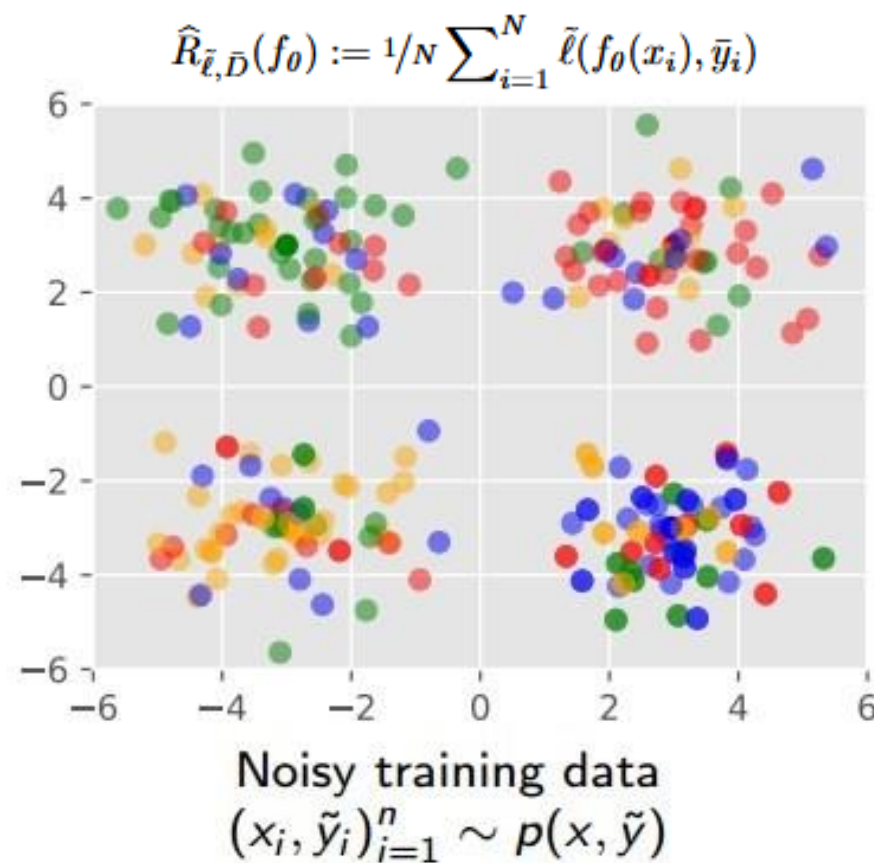
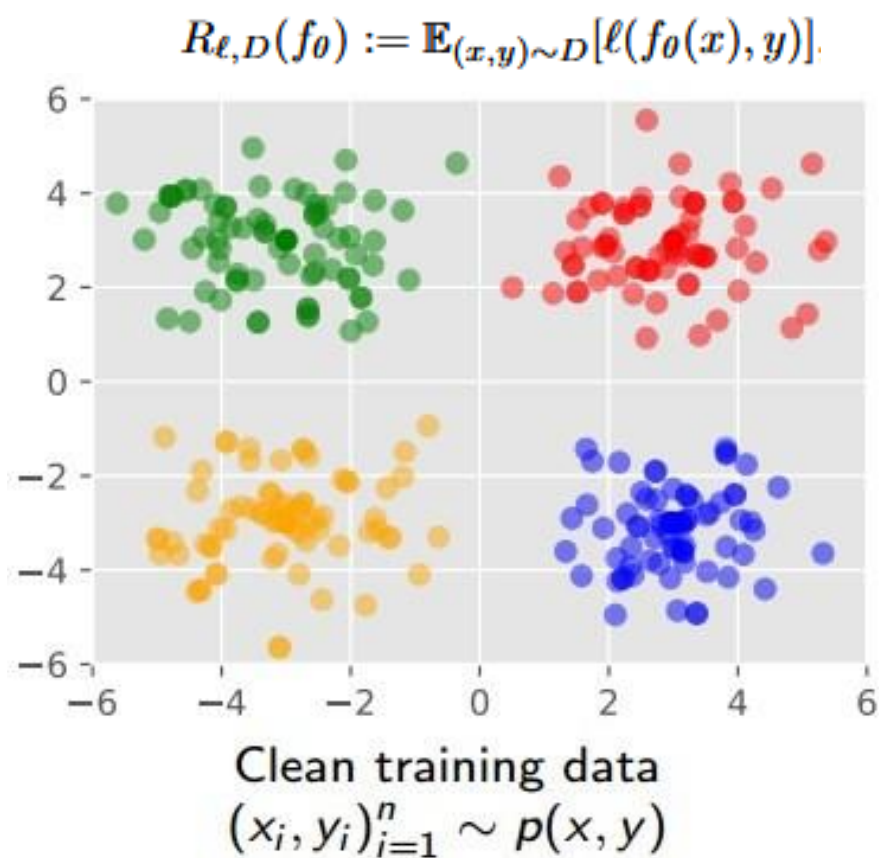


(Credit to Clothing1M)



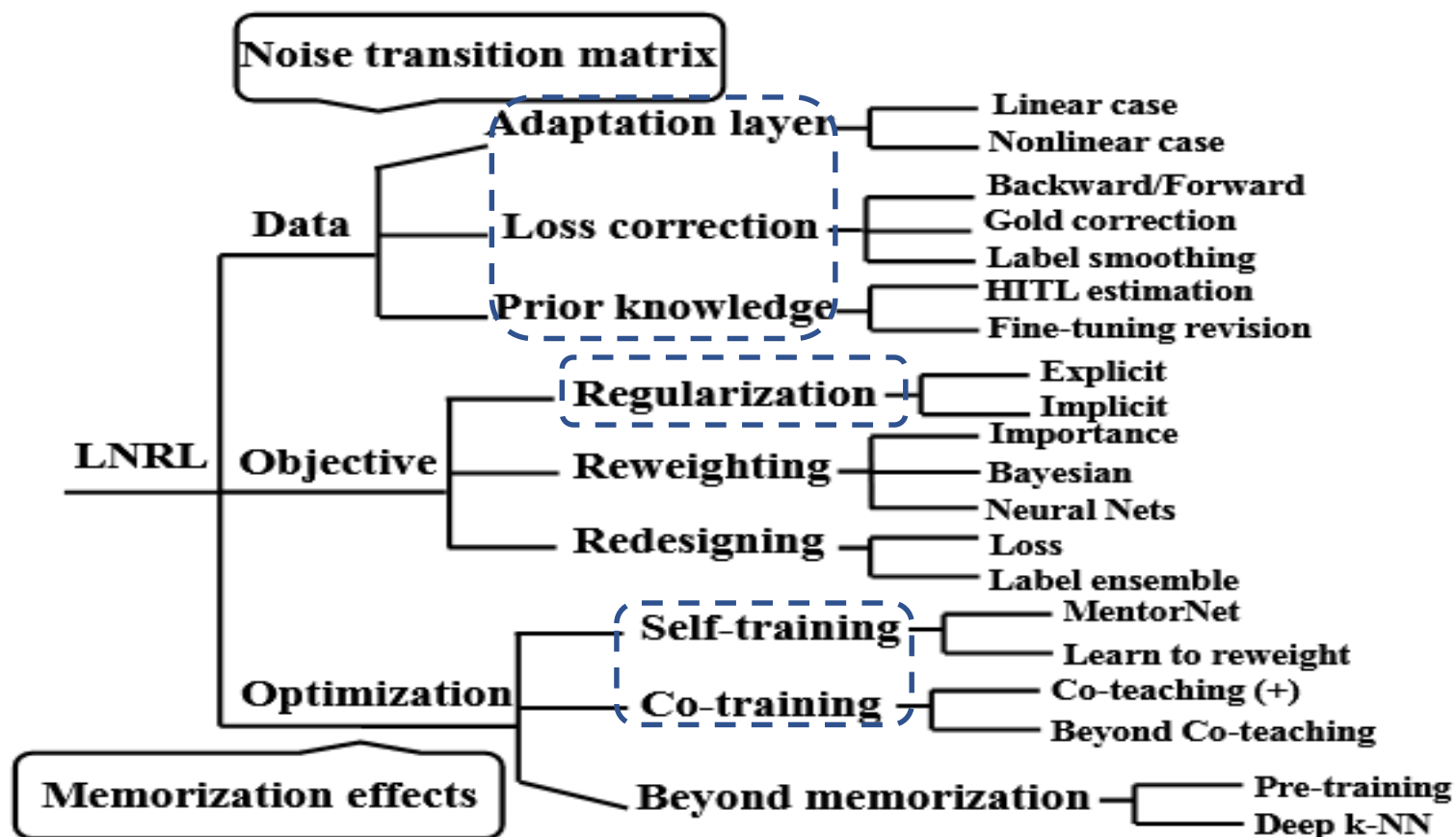
(Credit to Outlook)

# What are Noisy Labels

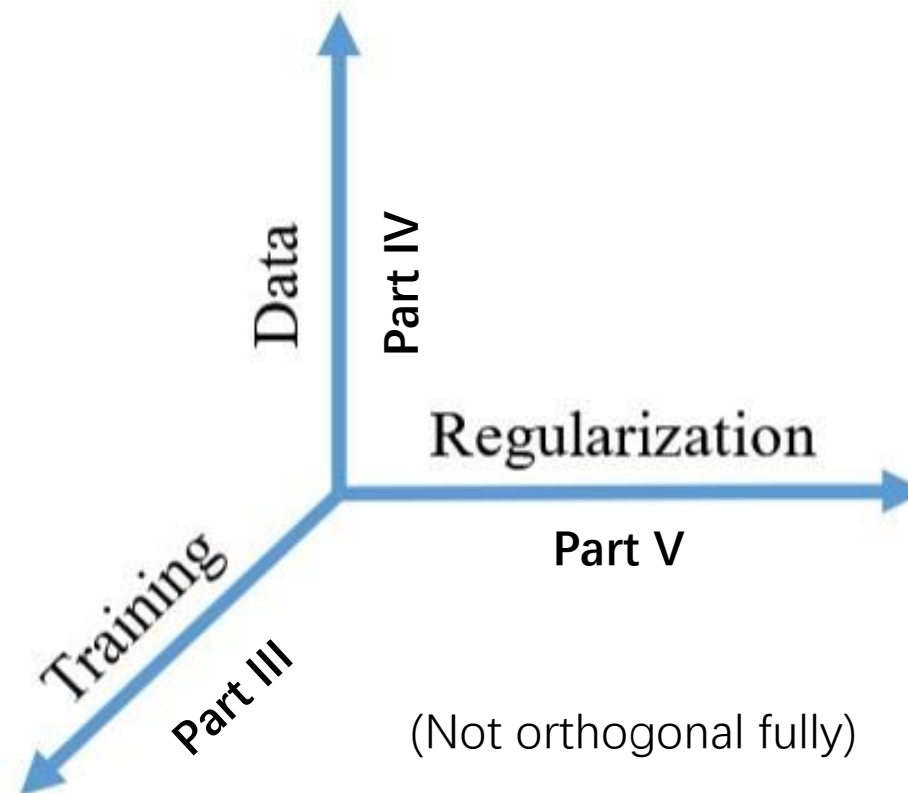


(Credit to Dr. Gang Niu)

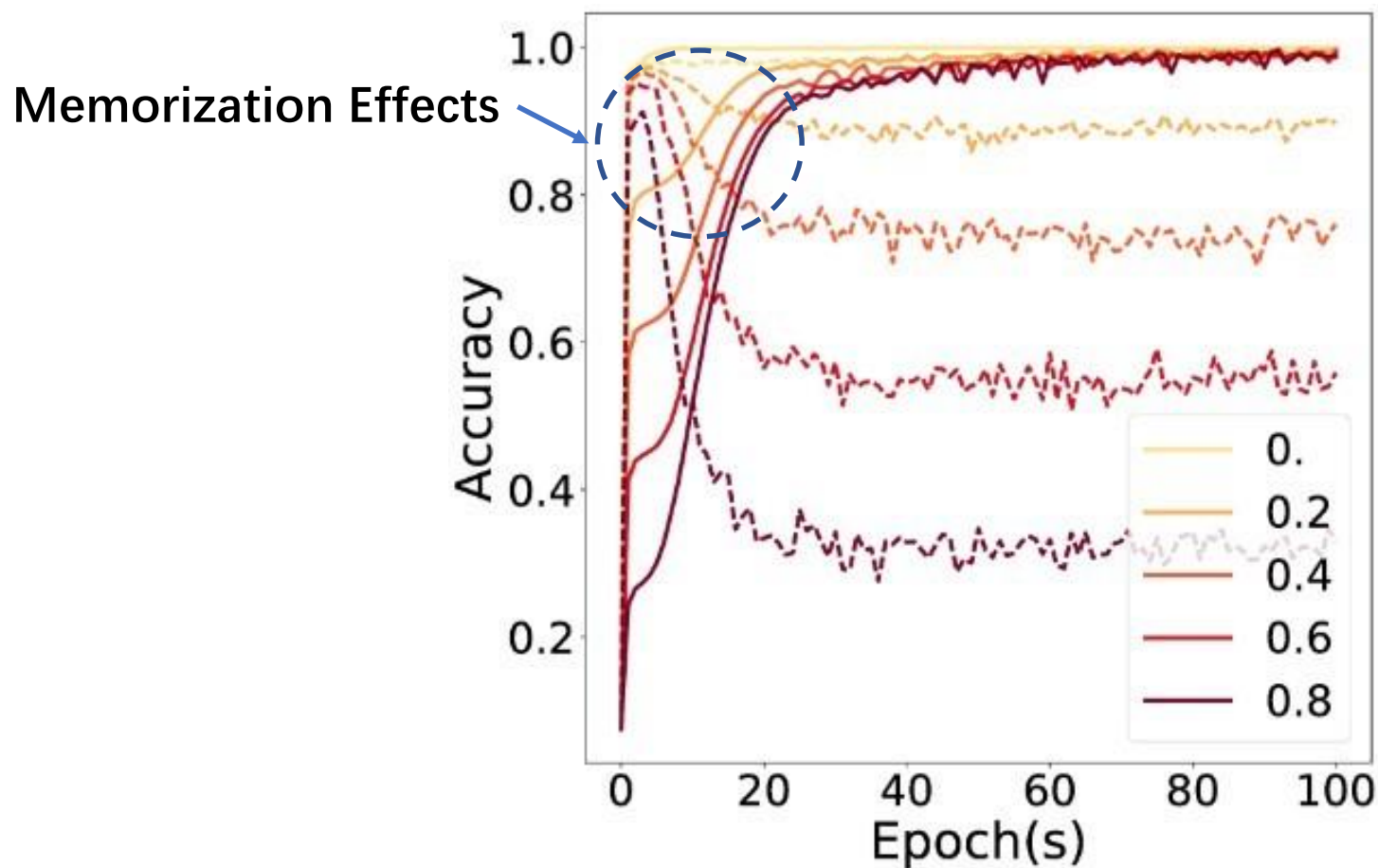
# Part II: Current Progress



# Tutorial Perspectives



# Part III: Training Perspective



# Training on Selected Samples

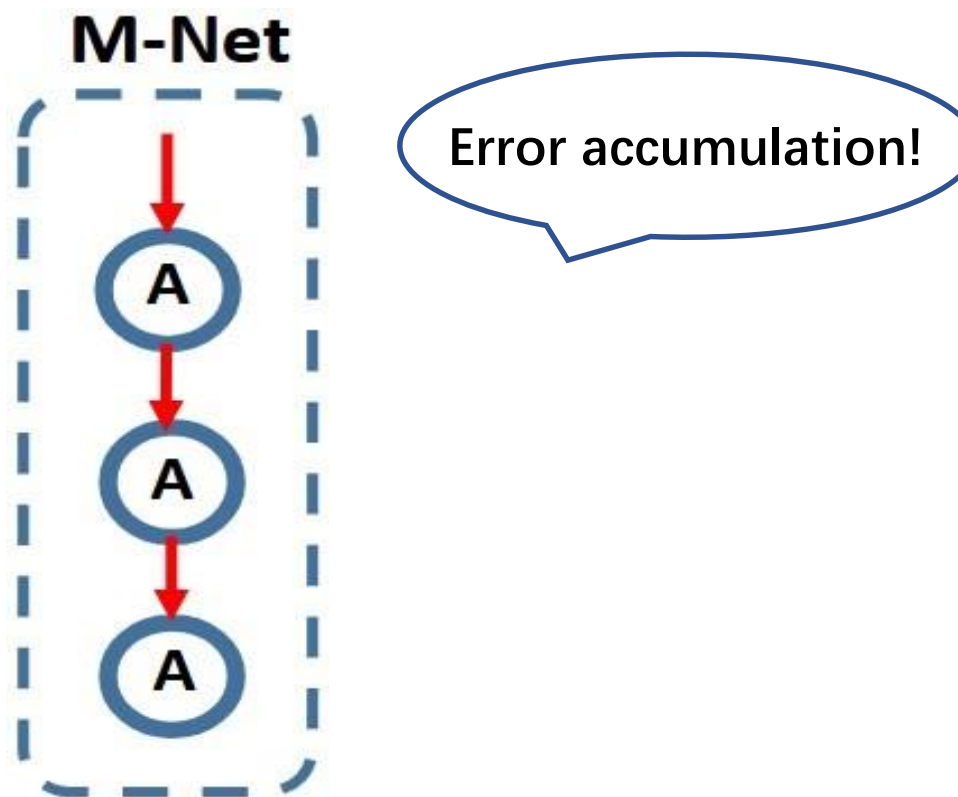
---

**Algorithm 1** General procedure on using sample selection to combat noisy labels.

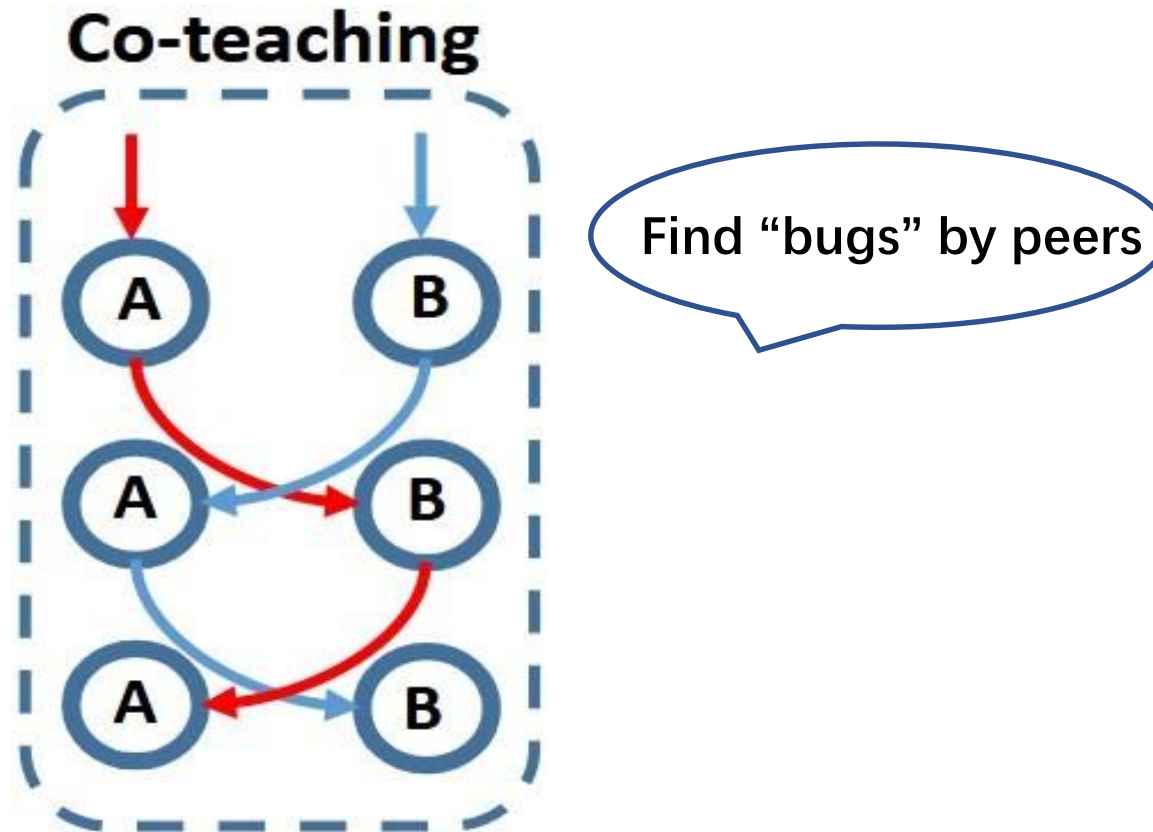
---

- 1: **for**  $t = 0, \dots, T - 1$  **do**
  - 2:   draw a mini-batch  $\bar{\mathcal{D}}$  from  $\mathcal{D}$ ;
  - 3:   select  $R(t)$  small-loss samples  $\bar{\mathcal{D}}_f$  from  $\bar{\mathcal{D}}$  based on network's predictions;
  - 4:   update network parameter using  $\bar{\mathcal{D}}_f$ ;
  - 5: **end for**
-

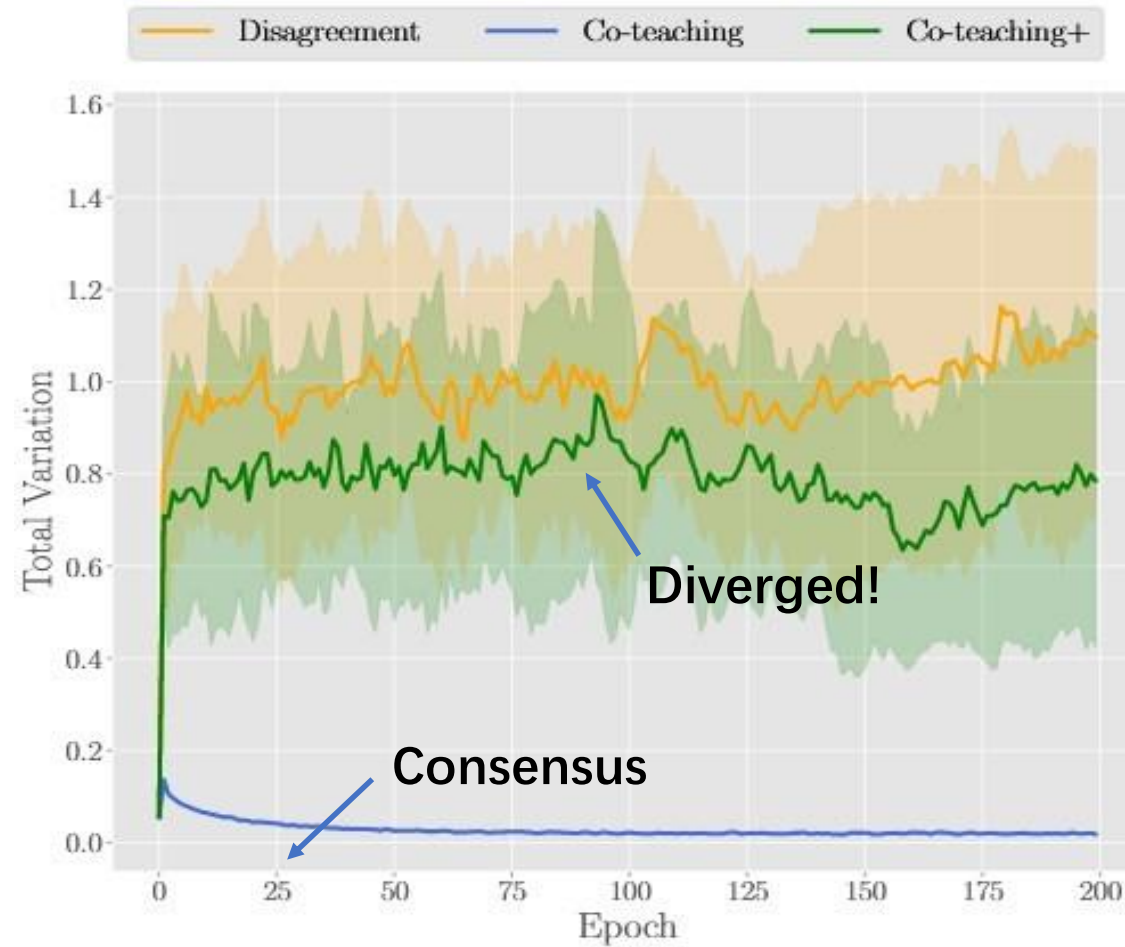
# Self-teaching (MentorNet)



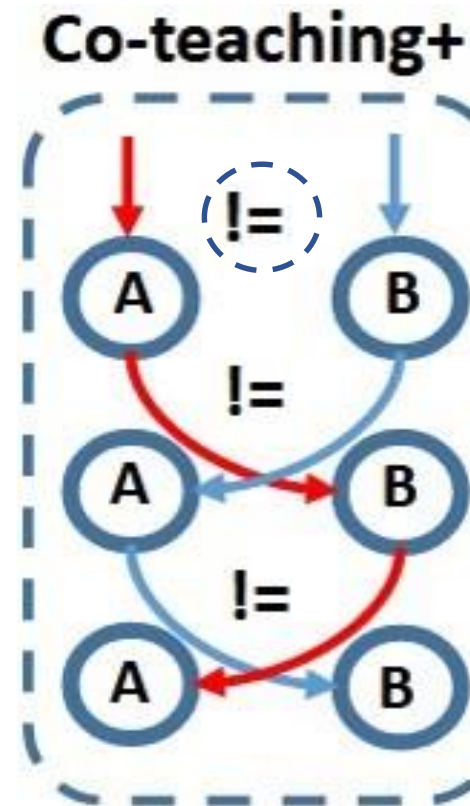
# Co-teaching



# Divergence Matters



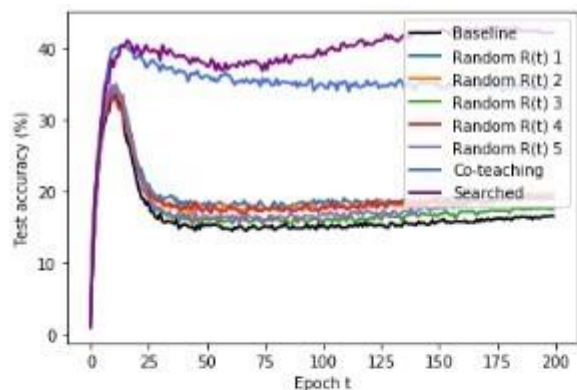
# Co-teaching+



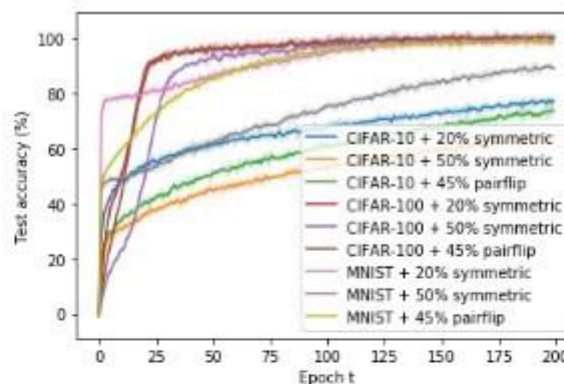
Divergence meeting  
Co-teaching

# Rethinking $R(t)$

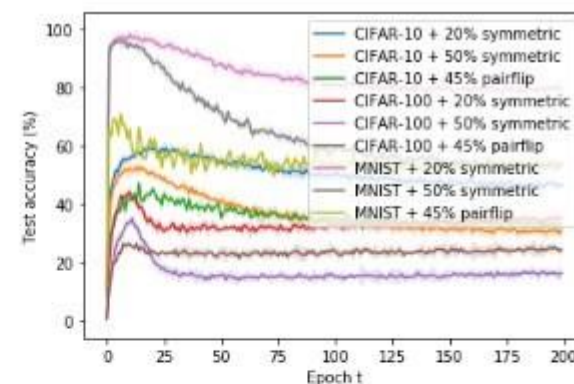
$$R(t) = 1 - \tau \cdot \min((t/t_k)^c, 1)$$



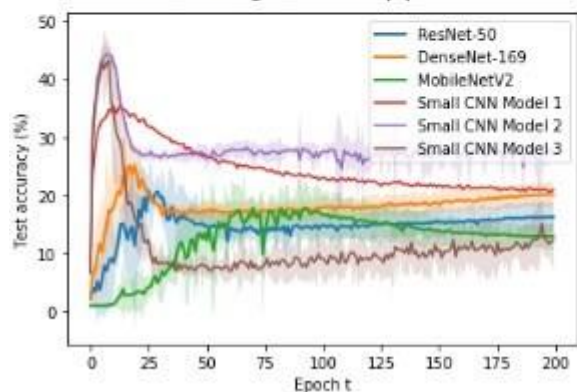
(a) Impact of  $R(t)$ .



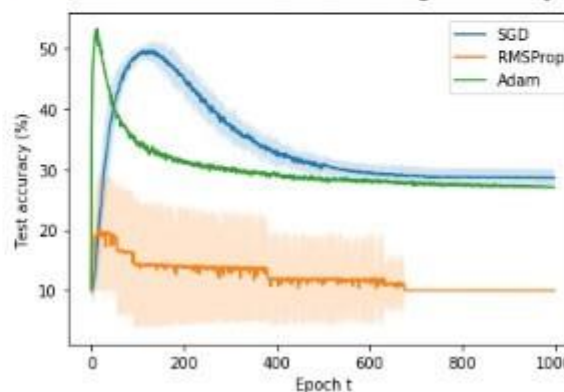
(b) Different data sets (training accuracy).



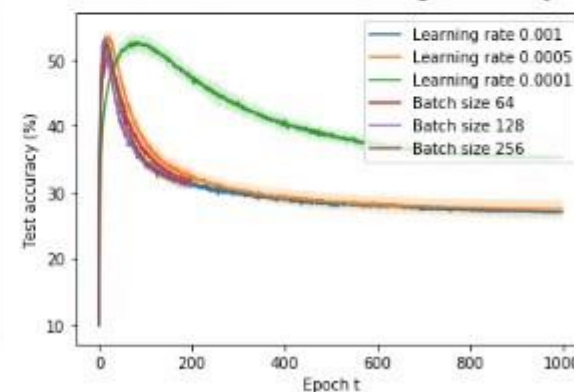
(c) Different data sets (testing accuracy).



(d) Different architectures.



(e) Different optimizers.

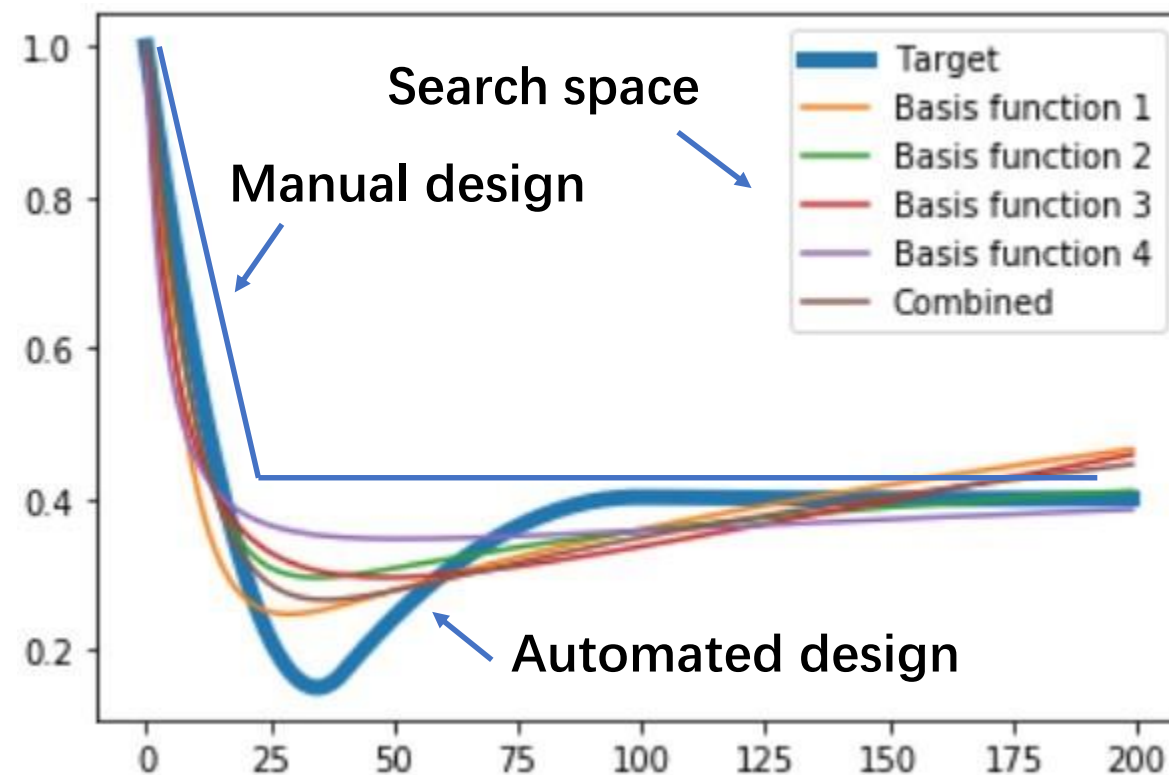


(f) Different optimizer settings.

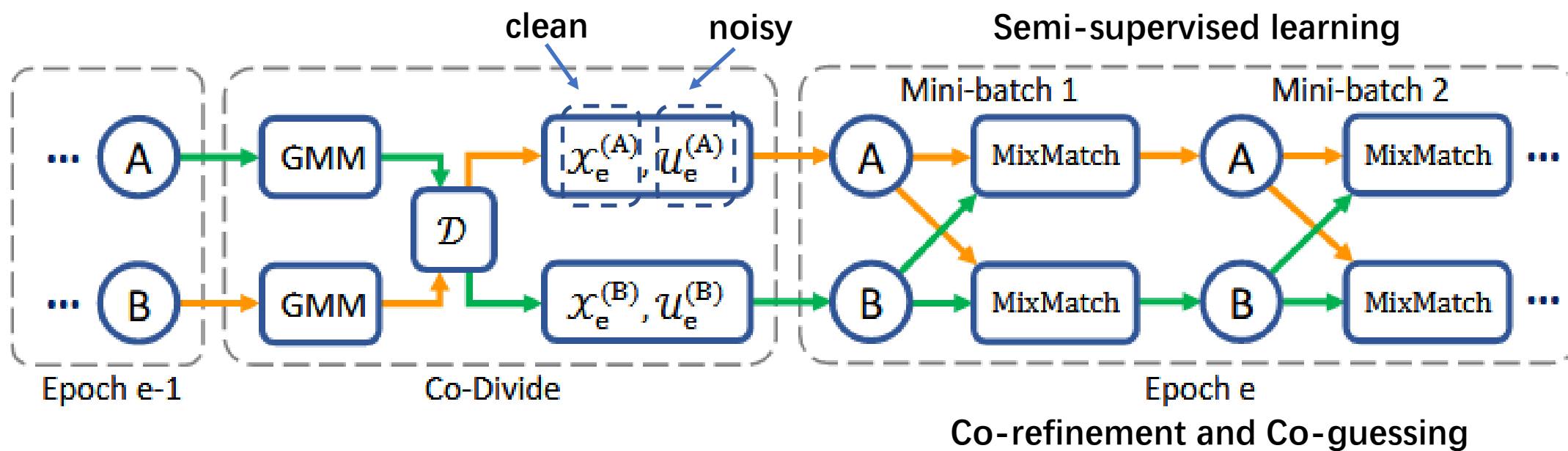
# S2E: Searching to Exploit

$$\begin{aligned}
 R^* &= \arg \min_{R(\cdot) \in \mathcal{F}} \mathcal{L}_{\text{val}}(f(\mathbf{w}^*; R), \mathcal{D}_{\text{val}}), \\
 \text{s.t. } \mathbf{w}^* &= \arg \min_{\mathbf{w}} \mathcal{L}_{\text{tr}}(f(\mathbf{w}; R), \mathcal{D}_{\text{tr}}).
 \end{aligned}$$

Bi-level Optimization

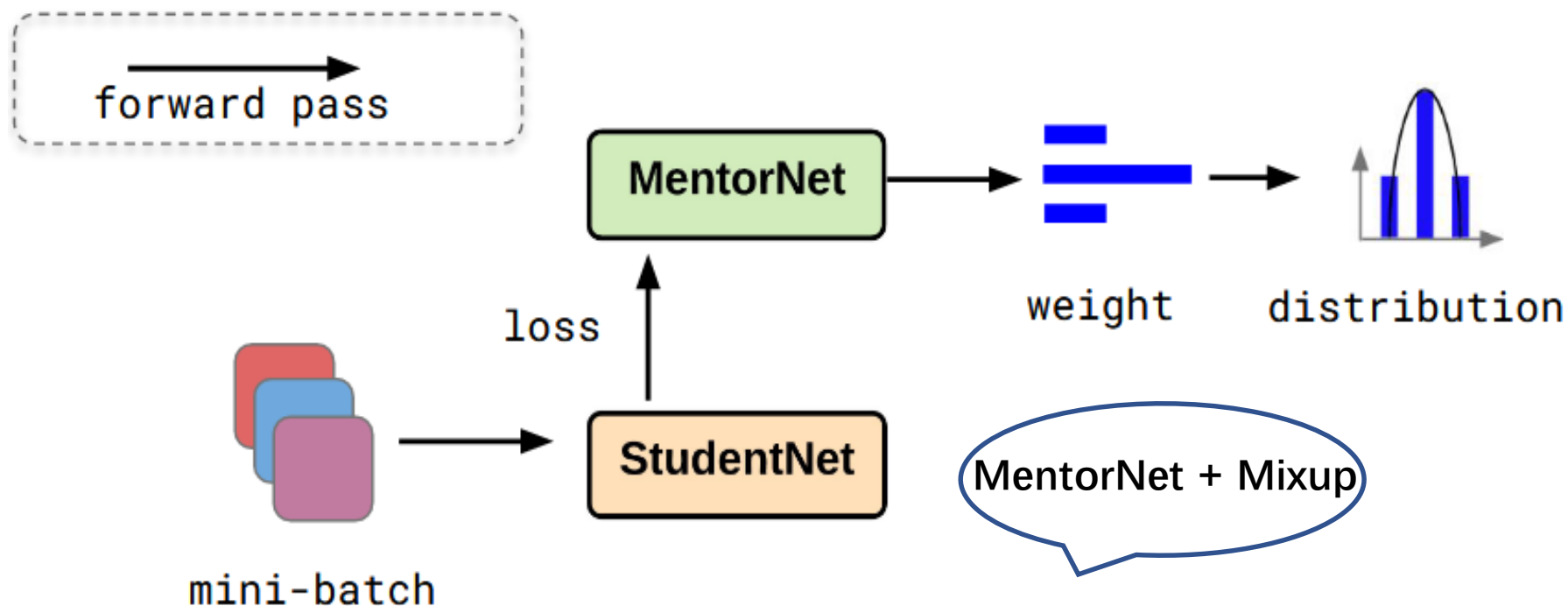


# DivideMix



# MentorMix

Weight  $\rightarrow$  Sample  $\rightarrow$  Mixup  $\rightarrow$  Weight



# Summary

- **Memorization effect** in deep learning is new and important.
- MentorNet and Co-teaching series are developed.
- Many **applications** have leveraged Co-teaching series.

# Part IV: Data Perspective

Diagram illustrating the Noise Transition Matrix for two types of noise: Sym-flipping and Pair-flipping.

**(a) Sym-flipping.**

The matrix  $y$  (labeled Dog) is shown with columns labeled Cat and Wolf. The matrix  $\tilde{y}$  is shown with columns labeled Cat and Wolf. The matrix  $y$  is a square matrix with elements:

$$y = \begin{bmatrix} 1-\tau & \frac{\tau}{n-1} & \dots & \frac{\tau}{n-1} \\ \frac{\tau}{n-1} & 1-\tau & \dots & \frac{\tau}{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\tau}{n-1} & \frac{\tau}{n-1} & \dots & 1-\tau \end{bmatrix}$$

The matrix  $\tilde{y}$  is a square matrix with elements:

$$\tilde{y} = \begin{bmatrix} 1-\tau & \tau \\ 0 & 1-\tau \\ \vdots & \vdots \\ 0 & \tau \\ \tau & 0 & \dots & 1-\tau \end{bmatrix}$$

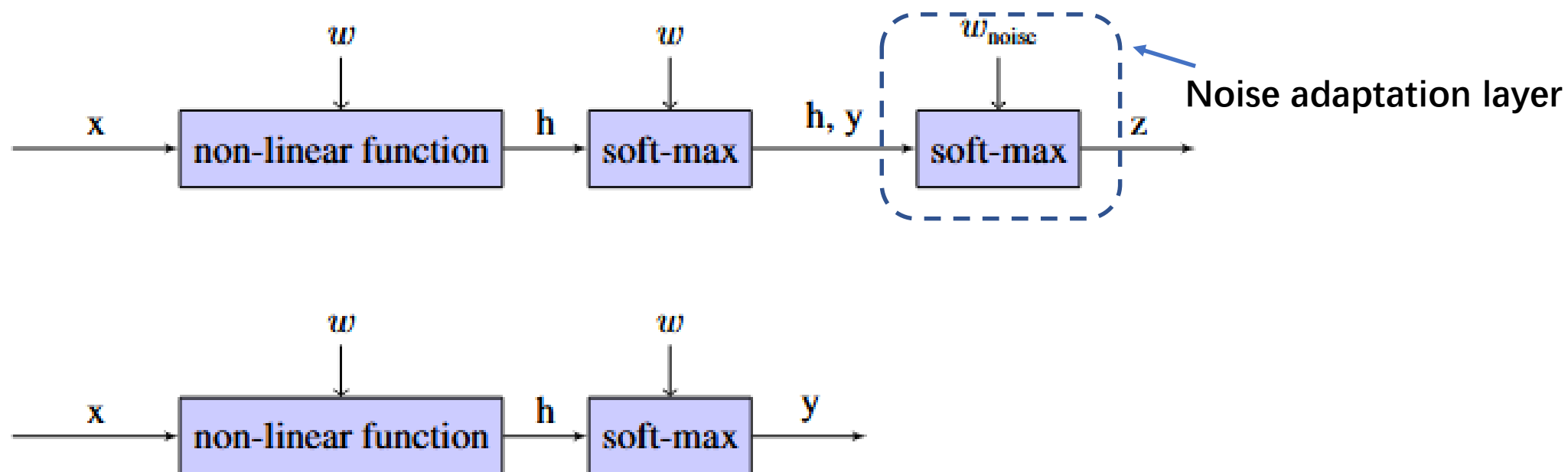
**(b) Pair-flipping.**

The matrix  $\tilde{y}$  is shown with columns labeled Cat and Wolf. The matrix  $\tilde{y}$  is a square matrix with elements:

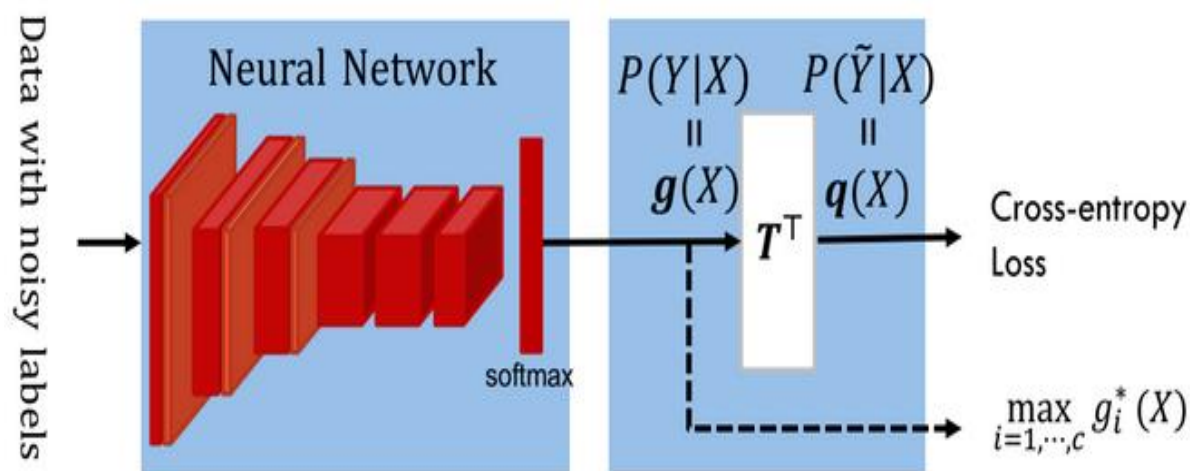
$$\tilde{y} = \begin{bmatrix} 1-\tau & \tau \\ 0 & 1-\tau \\ \vdots & \vdots \\ 0 & \tau \\ \tau & 0 & \dots & 1-\tau \end{bmatrix}$$

Noise Transition Matrix

# Adaptation layer



# Forward Correction



(Credit to Dr. Tongliang Liu)

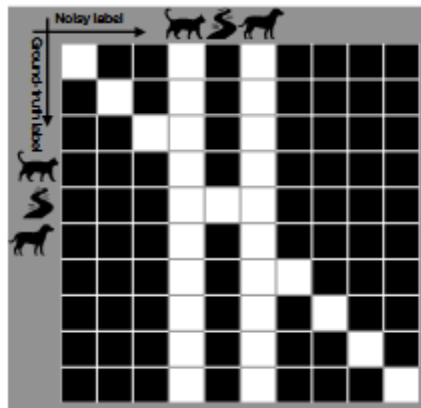
**Theorem 2.** (Forward Correction, Theorem 1 in [22]) Suppose that the label transition matrix  $T$  is non-singular, where  $T_{ij} = p(\bar{y} = j | y = i)$  given that corrupted label  $\bar{y} = j$  is flipped from clean label  $y = i$ . Given loss  $\ell$  and network function  $f$ , Forward Correction is defined as

$$\ell^{\rightarrow}(f(x), \bar{y}) = [\ell_{y|T^T f(x)}]_{\bar{y}}, \quad (6)$$

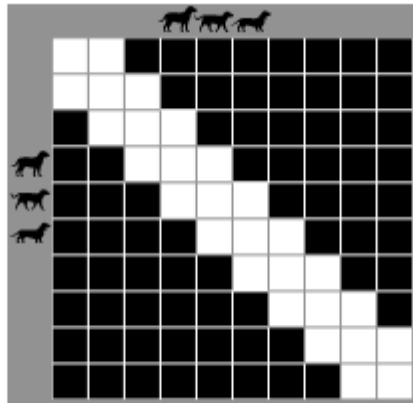
where  $\ell_{y|T^T f(x)} = (\ell(T^T f(x), 1), \dots, \ell(T^T f(x), k))$ . Then, the minimizer of the corrected loss under the noisy distribution is the same as the minimizer of the original loss under the clean distribution, namely,

$$\arg \min_f \mathbb{E}_{x, \bar{y}} \ell^{\rightarrow}(f(x), \bar{y}) = \arg \min_f \mathbb{E}_{x, y} \ell(f(x), y). \quad (7)$$

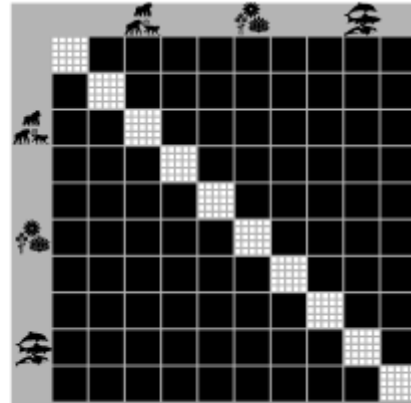
# Masking



(a) Column-diagonal



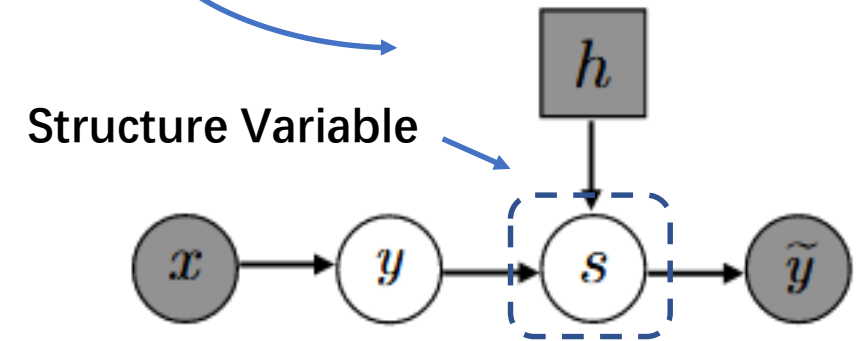
(b) Tri-diagonal



(c) Block-diagonal

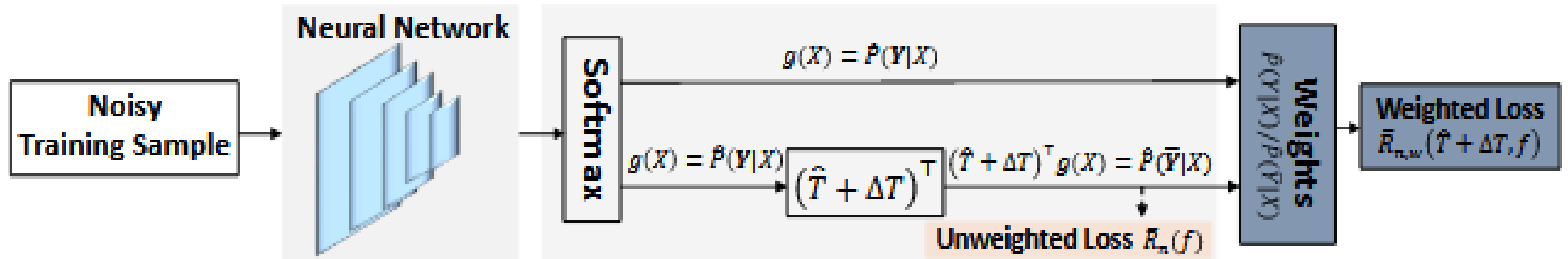


(a) Benchmark model.



(b) MASKING model.

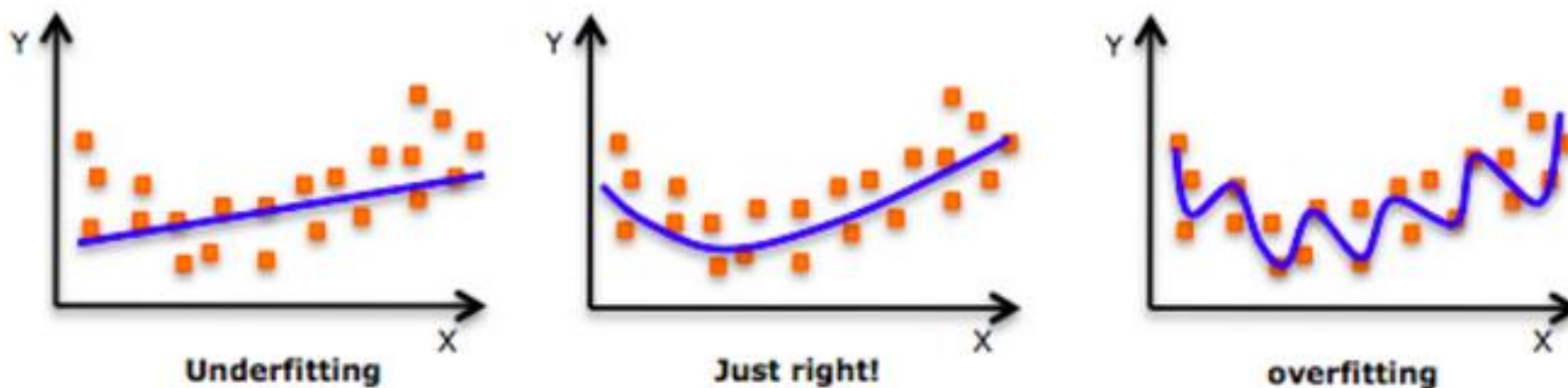
# Fine-tuning



# Summary

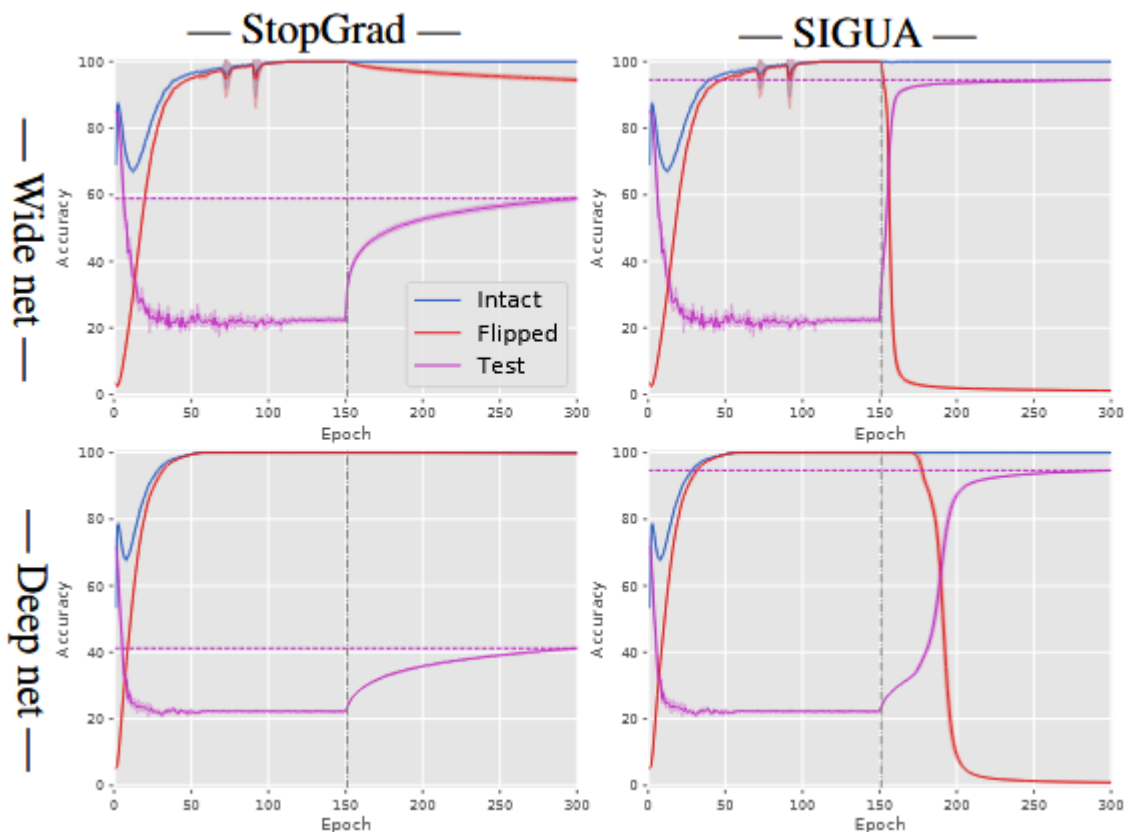
- **Noise transition matrix** is the key in data perspective.
- A potential direction is how to estimate this matrix **easily**.
- Another potential direction is how to leverage this matrix **effectively**.

# Part V: Regularization Perspective



(Credit to Analytics Vidhya)

# SIGUA



## Algorithm 1 SIGUA-prototype (in a mini-batch).

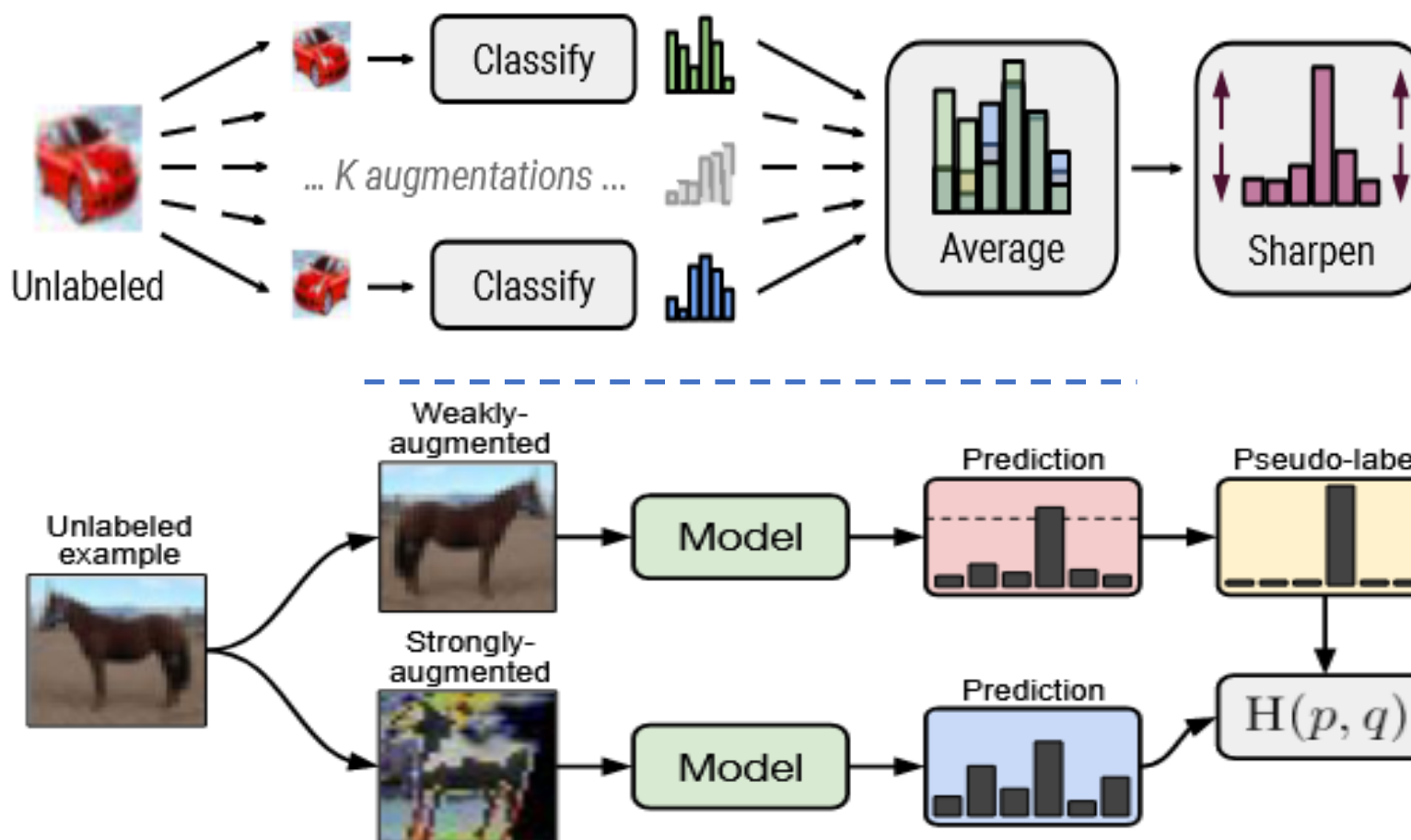
**Require:** base learning algorithm  $\mathcal{B}$ , optimizer  $\mathcal{O}$ ,  
mini-batch  $\mathcal{S}_b = \{(x_i, \tilde{y}_i)\}_{i=1}^{n_b}$  of batch size  $n_b$ ,  
current model  $f_\theta$  where  $\theta$  holds the parameters of  $f$ ,  
good- and bad-data conditions  $\mathcal{C}_{\text{good}}$  and  $\mathcal{C}_{\text{bad}}$  for  $\mathcal{B}$ ,  
underweight parameter  $\gamma$  such that  $0 \leq \gamma \leq 1$

```

1:  $\{\ell_i\}_{i=1}^{n_b} \leftarrow \mathcal{B}.\text{forward}(f_\theta, \mathcal{S}_b)$  # forward pass
2:  $\ell_b \leftarrow 0$  # initialize loss accumulator
3: for  $i = 1, \dots, n_b$  do
4:   if  $\mathcal{C}_{\text{good}}(x_i, \tilde{y}_i)$  then
5:      $\ell_b \leftarrow \ell_b + \ell_i$  # accumulate loss positively
6:   else if  $\mathcal{C}_{\text{bad}}(x_i, \tilde{y}_i)$  then  $\leftarrow$  Gradient Ascent
7:      $\ell_b \leftarrow \ell_b - \gamma \ell_i$  # accumulate loss negatively
8:   end if # ignore any uncertain data
9: end for
10:  $\ell_b \leftarrow \ell_b / n_b$  # average accumulated loss
11:  $\nabla_\theta \leftarrow \mathcal{B}.\text{backward}(f_\theta, \ell_b)$  # backward pass
12:  $\mathcal{O}.\text{step}(\nabla_\theta)$  # update model

```

# MixMatch & FixMatch



D. Berthelot et al. MixMatch: A Holistic Approach to Semi-supervised Learning. In *NeurIPS*, 2019.

K. Sohn et al. FixMatch: Simplifying Semi-supervised Learning with Consistency and Confidence. In *NeurIPS*, 2020.

# Bootstrapping

$$\ell_{\text{soft}}(q, t) = \sum_{k=1}^L [\beta t_k + (1 - \beta) q_k] \log(q_k)$$

target prediction

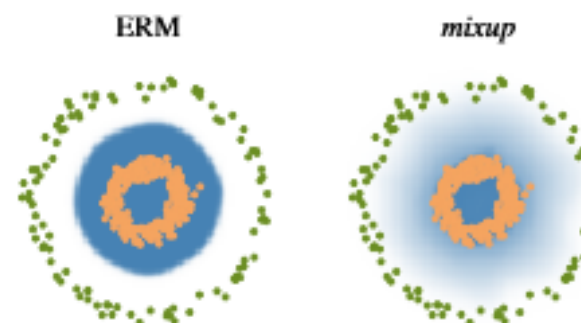
$$\ell_{\text{hard}}(q, t) = \sum_{k=1}^L [\beta t_k + (1 - \beta) z_k] \log(q_k)$$

# Mixup

```
# y1, y2 should be one-hot vectors
for (x1, y1), (x2, y2) in zip(loader1, loader2):
    lam = numpy.random.beta(alpha, alpha)
    x = Variable(lam * x1 + (1. - lam) * x2)
    y = Variable(lam * y1 + (1. - lam) * y2)
    optimizer.zero_grad()
    loss(net(x), y).backward()
    optimizer.step()
```

interpolation

(a) One epoch of *mixup* training in PyTorch.



(b) Effect of *mixup* ( $\alpha = 1$ ) on a toy problem. Green: Class 0. Orange: Class 1. Blue shading indicates  $p(y = 1|x)$ .

# Summary

- Regularization is very popular for **semi-supervised learning**.
- Explicit regularization is in the level of **objective function**.
- Implicit regularization is in the level of **algorithm** and **data**.

# Part VI: Future Directions

## A Survey of Label-noise Representation Learning: Past, Present and Future

Bo Han, Quanming Yao, Tongliang Liu, Gang Niu,  
Ivor W. Tsang, James T. Kwok, *Fellow, IEEE* and Masashi Sugiyama

**Abstract**—Classical machine learning implicitly assumes that labels of the training data are sampled from a clean distribution, which can be too restrictive for real-world scenarios. However, statistical-learning-based methods may not train deep learning models robustly with these noisy labels. Therefore, it is urgent to design Label-Noise Representation Learning (LNRL) methods for robustly training deep models with noisy labels. To fully understand LNRL, we conduct a survey study. We first clarify a formal definition for LNRL from the perspective of machine learning. Then, via the lens of learning theory and empirical study, we figure out why noisy labels affect deep models' performance. Based on the theoretical guidance, we categorize different LNRL methods into three directions. Under this unified taxonomy, we provide a thorough discussion of the pros and cons of different categories. More importantly, we summarize the essential components of robust LNRL, which can spark new directions. Lastly, we propose possible research directions within LNRL, such as new datasets, instance-dependent LNRL, and adversarial LNRL. We also envision potential directions beyond LNRL, such as learning with feature-noise, preference-noise, domain-noise, similarity-noise, graph-noise and demonstration-noise.

**Index Terms**—Machine Learning, Representation Learning, Weakly Supervised Learning, Label-noise Learning, Noisy Labels.

20 Feb 2021

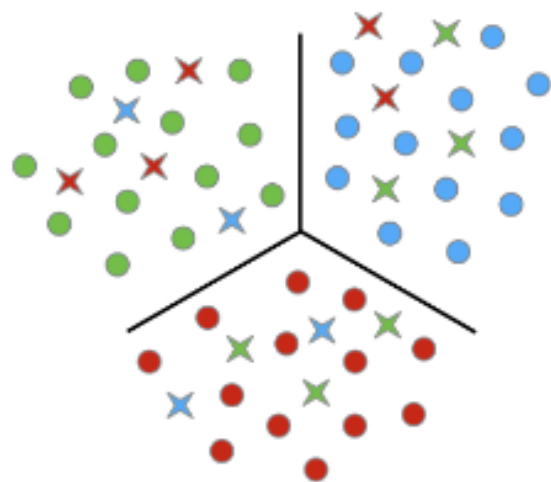
B. Han, Q. Yao, T. Liu, G. Niu, I. W. Tsang, J. T. Kwok, and M. Sugiyama.

A Survey of Label-noise Representation Learning: Past, Present and Future. *arXiv preprint: 2011.04406*, 2020.

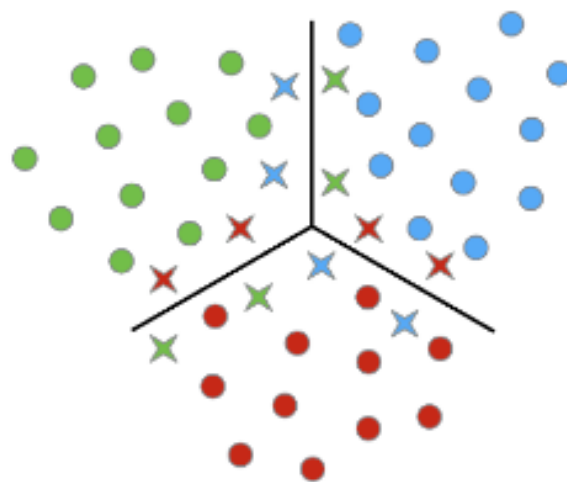
# New Datasets



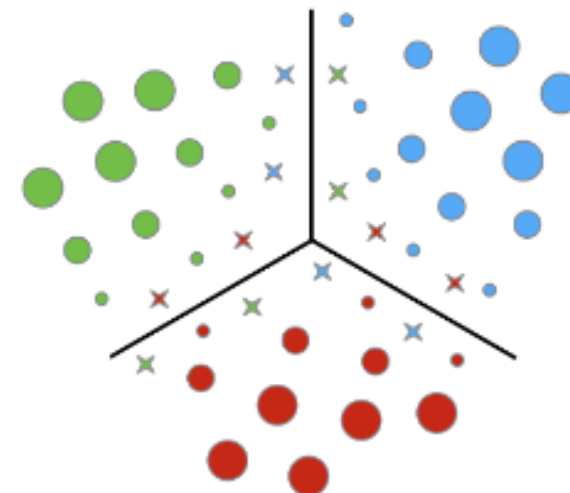
# Instance-dependent LNRL



(a) Class-conditional noise.

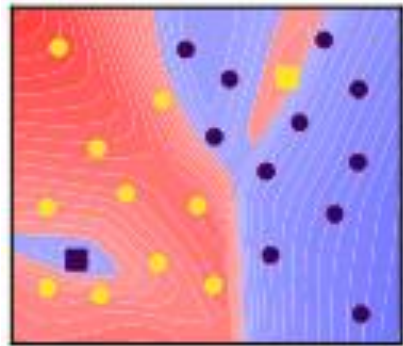


(b) Instance-dependent noise  
(boundary-consistent noise).

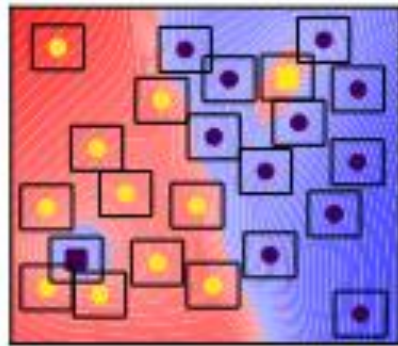


(c) Confidence-scored instance-dependent  
noise.

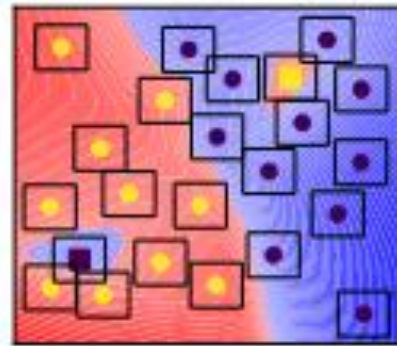
# Adversarial LNRL



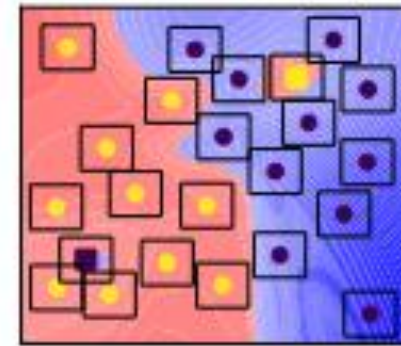
ST



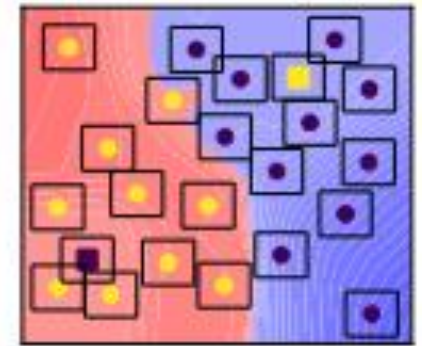
AT (PGD-1)



AT (PGD-2)



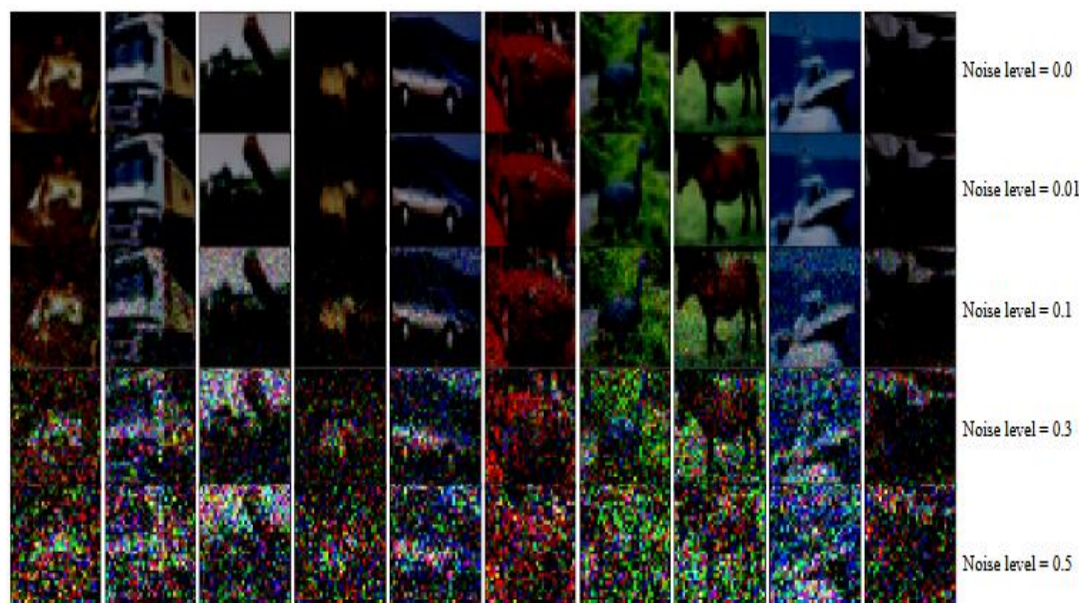
AT (PGD-3)



AT (PGD-4)

weak  $\longrightarrow$  strong

# Noisy Feature



Image

video games good for children computer games can promote problem-solving and team-building in children, say games industry experts. (Noise level = 0.0)

vedeo games good for dhildlenzcospxter games can iromote problem-sorvtng and teai-building in children, sby games industry experts. (Noise level = 0.1)

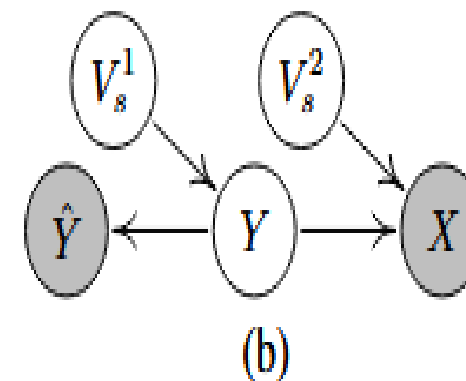
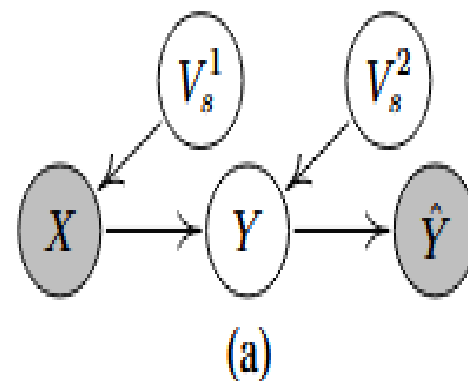
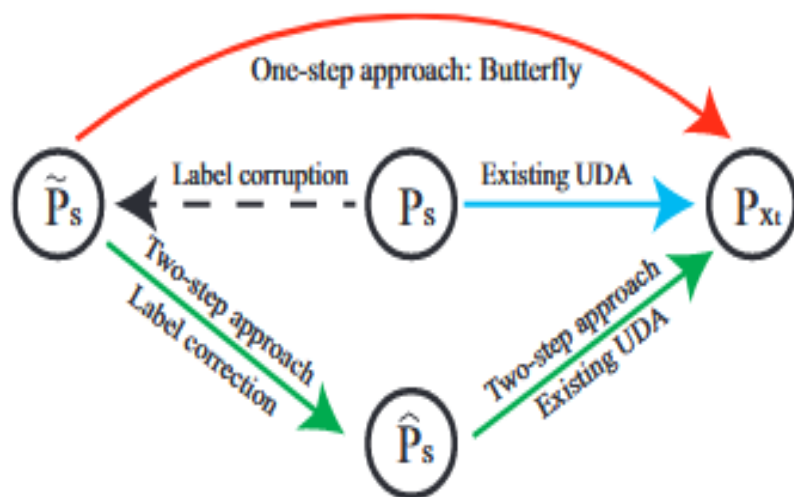
video nawvs zggood foryxhiltrengomvumer games cahcprocotubpnoblex-szbvina and tqlmmbuaddiagjin whipdren, saywgsmes ildustry exmrtrts. (Noise level = 0.3)

tmdeo gakec jgopd brr cgildrenjcoogwdeh bxdeu vanspromote xrobkeh-svlkieo and termwwwuojvinguinfcjdbdses, sacosamlt cndgstoyaagpbrus. (Noise level = 0.5)

vizwszgbrwtguihcxfatbhivrrvwq cxmpgugflziwls clfnzrommtohprtblef-solvynx rnjnyiaf-gjwlcergwklskqibdtjn,aoty gameshinzustrm expertsdm (Noise level = 0.8)

Text

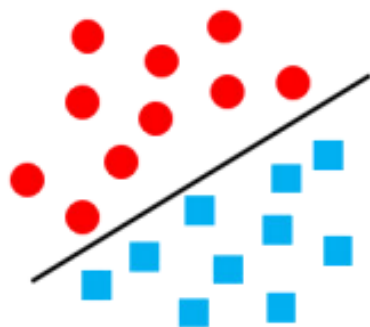
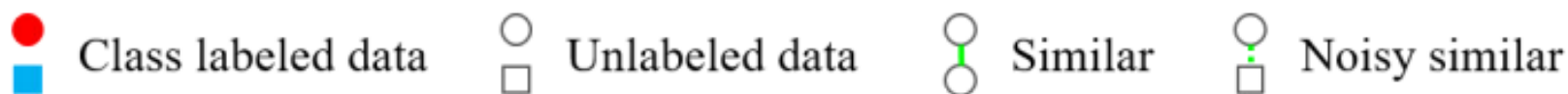
# Noisy Domain



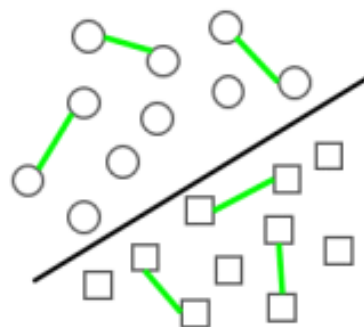
F. Liu et al. Butterfly: One-step Approach towards Wildly Unsupervised Domain Adaptation. *arXiv preprint*, 2019.

X. Yu et al. Label-noise Robust Domain Adaptation. In *ICML*, 2020.

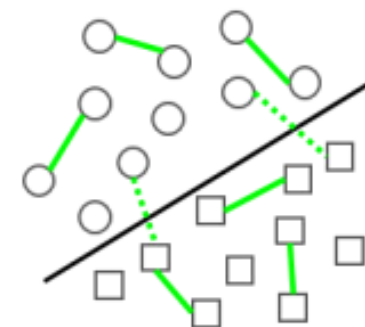
# Noisy Similarity



(a) Supervised Classification

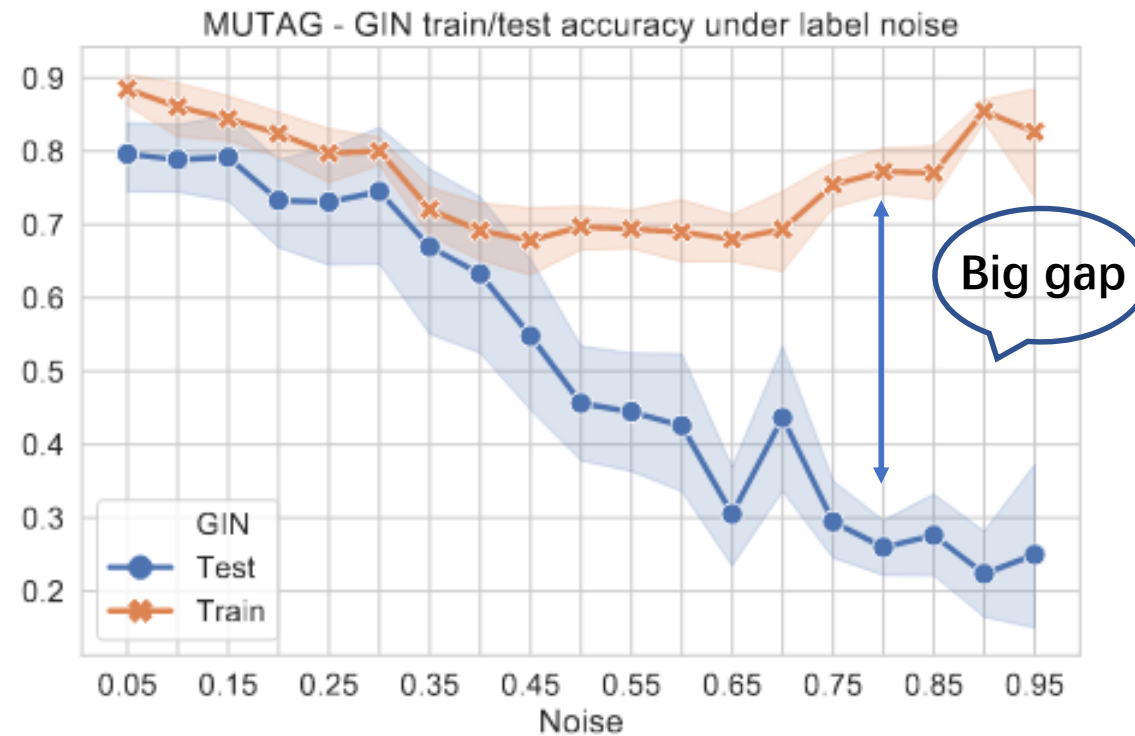


(b) SU Classification

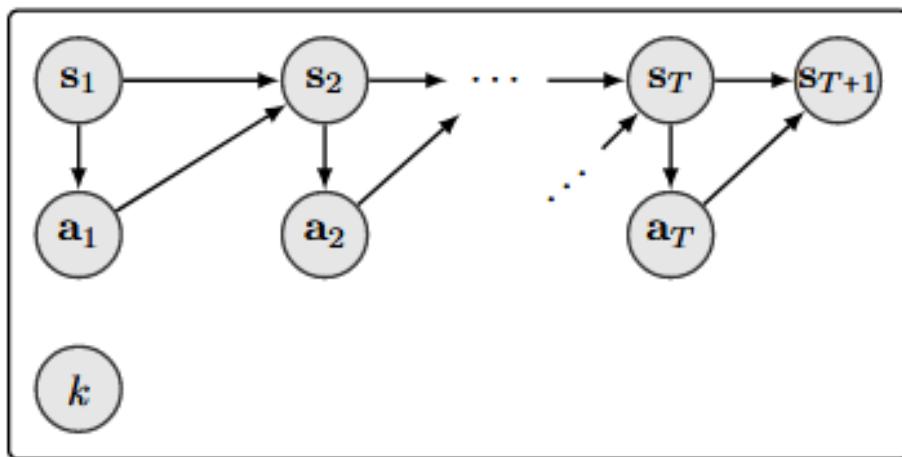


(c) NSU Classification

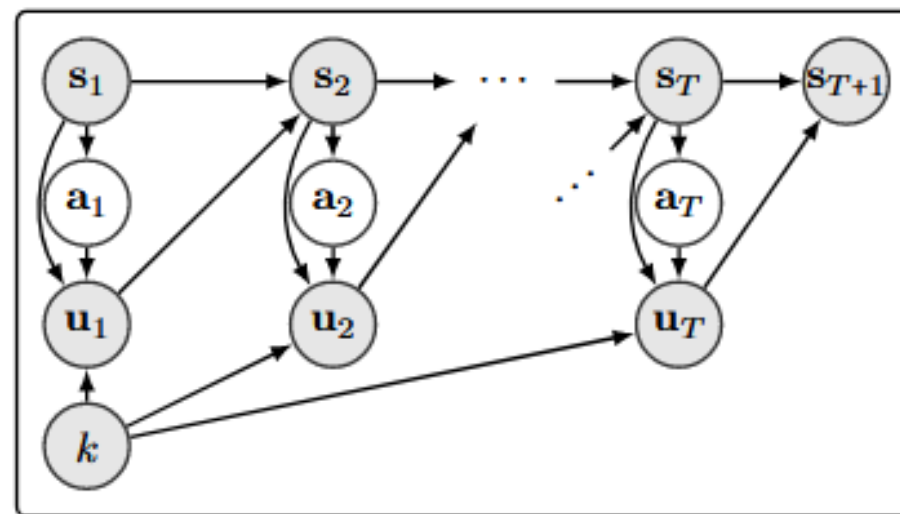
# Noisy Graph



# Noisy Demonstration



(a) Expert demonstrations



(b) Diverse-quality demonstrations

# Noisy Machine Translation

Chinese-English (ISI bitext)	
<b>Src:</b>	美国提出的报复清单是中国政府绝对不能接受的。
<b>Trg:</b>	And the Chinese side would certainly not accept the unreasonable demands put forward by the Americans concerning the protection of intellectual property rights.
<b>Human:</b>	The revenge list proposed by America will definitely not be accepted by Chinese government.

# Conclusions

- Current progress mainly focuses on **class-conditional noise**.
- The new trend focuses on **instance-dependent noise**.
- Besides noisy labels, we should pay more efforts on **noisy data**.