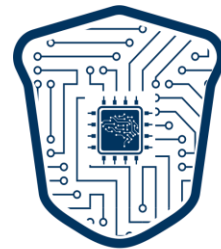# Trustworthy Machine Learning under Noisy Data

Dr. Bo Han

HKBU TMLR Group / RIKEN AIP Team

Assistant Professor / BAIHO Visiting Scientist
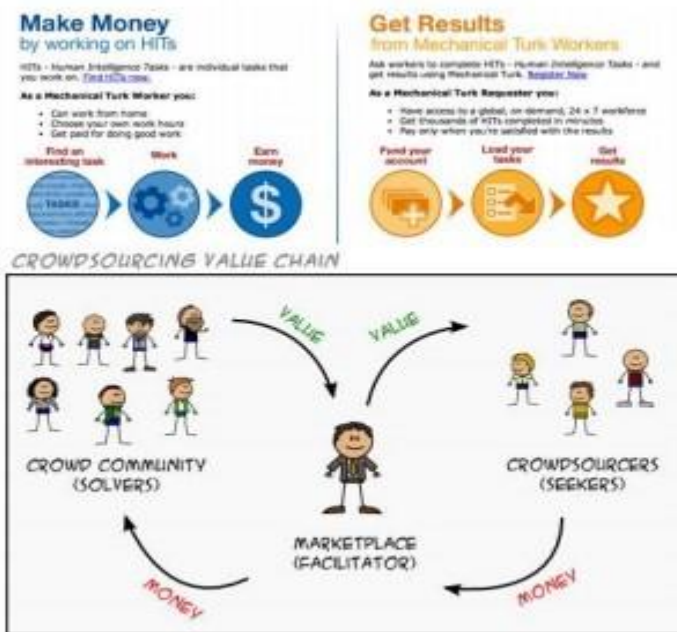
https://bhanml.github.io/

# 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

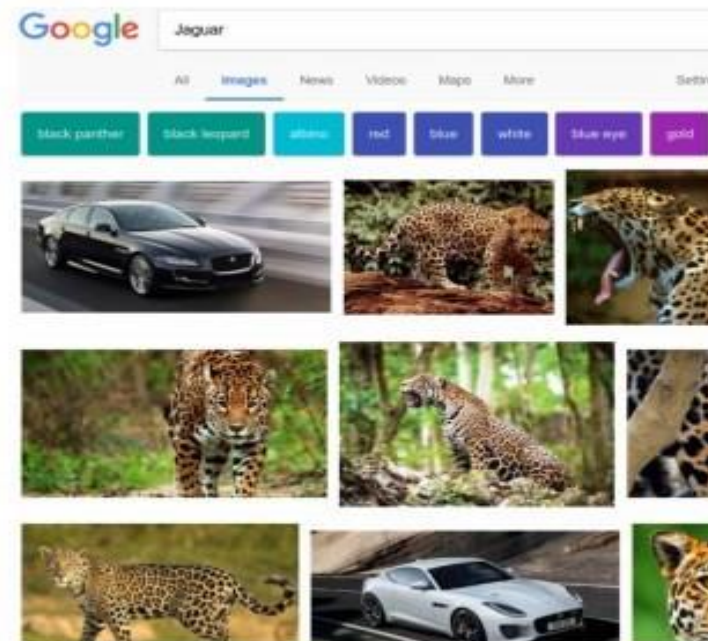- Part VI: Future Directions

# Part I: Why Noisy Labels



(Credit to Amazon)

(Credit to Google)

# Why Noisy Labels
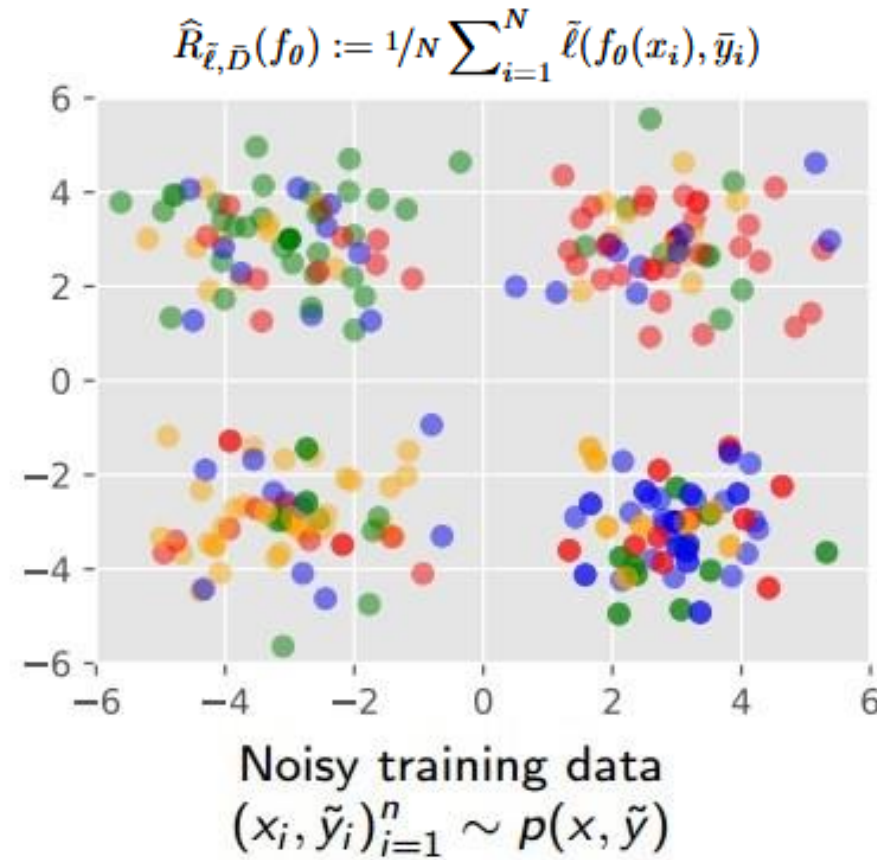


(Credit to Clothing1M)



(Credit to Outlook)

# What are Noisy Labels



$$R_{\ell,D}(f_0) := \mathbb{E}_{(x,y)\sim D}[\ell(f_0(x), y)]$$

$$\hat{R}_{\tilde{\ell},\bar{D}}(f_0) := 1/N \sum_{i=1}^{N} \tilde{\ell}(f_0(x_i), \bar{y}_i)$$

Clean training data
$(x_i, y_i)_{i=1}^{n} \sim p(x, y)$

Noisy training data
$(x_i, \tilde{y}_i)_{i=1}^{n} \sim p(x, \tilde{y})$

(Credit to Dr. Gang Niu)

# Part II: Current Progress

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.

# Tutorial Perspectives



Data

Part IV

Regularization

Part V

Training

Part III

(Not orthogonal fully)

# Part III: Training Perspective



**Memorization Effects**

D. Arpit et al. A Closer Look at Memorization in Deep Networks. In *ICML*, 2017.

# Training on Selected Samples

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

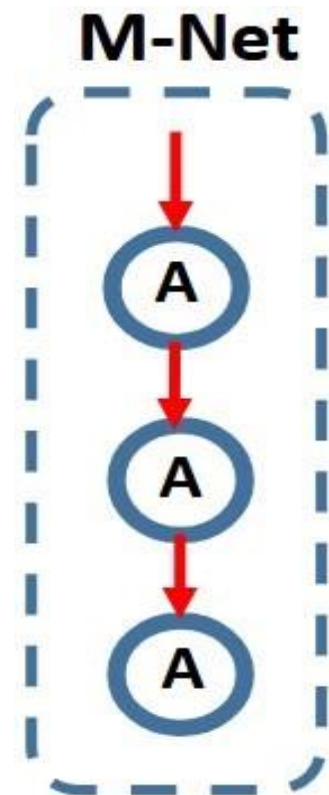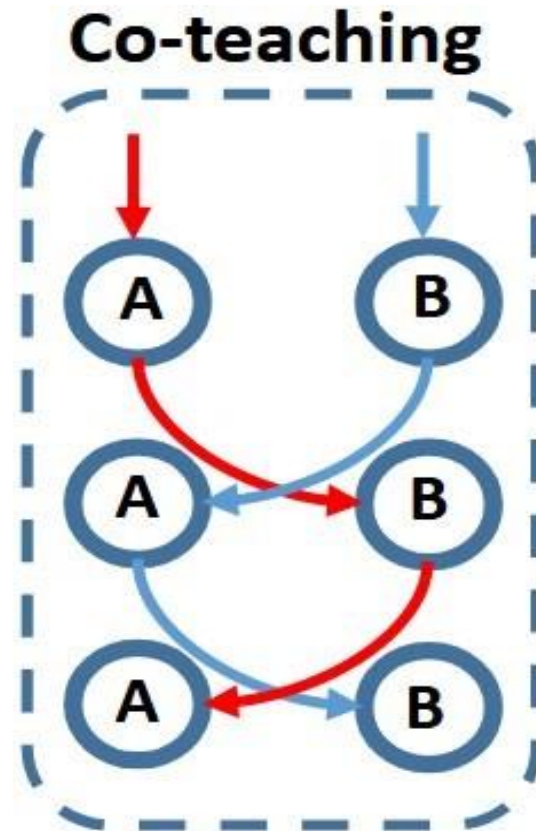1: **for** $t = 0, \ldots, T - 1$ **do**
2:      draw a mini-batch $\bar{\mathcal{D}}$ from $\mathcal{D}$;
3:      select $R(t)$ small-loss samples $\bar{\mathcal{D}}_f$ from $\bar{\mathcal{D}}$ based on network's predictions;
4:      update network parameter using $\bar{\mathcal{D}}_f$;
5: **end for**

# Self-teaching (MentorNet, 2018)

L. Jiang et al. MentorNet: Learning Data-Driven Curriculum for Very Deep Neural Networks on Corrupted Data. In *ICML*, 2018.

https://bhanml.github.io/ & https://github.com/tmlr-group

# Co-teaching (2018)



Find "bugs" by peers

B. Han et al. Co-teaching: Robust Training of Deep Neural Networks with Extremely Noisy Labels. In *NeurIPS*, 2018.

# Divergence Matters

# Co-teaching+ (2019)



Co-teaching+

!=

Divergence meeting Co-teaching

X. Yu et al. How does Disagreement Help Generalization against Label Corruption? In *ICML*, 2019.

# Rethinking R(t)

**Test accuracy depends on selecting rules**

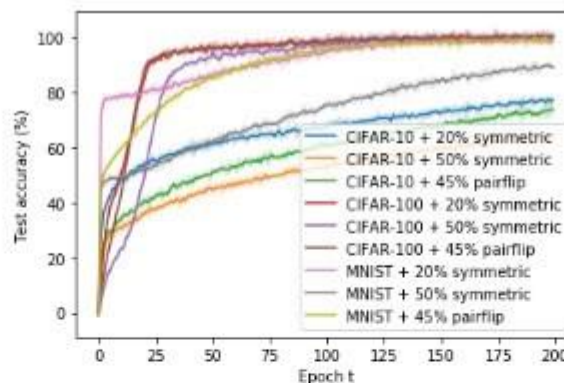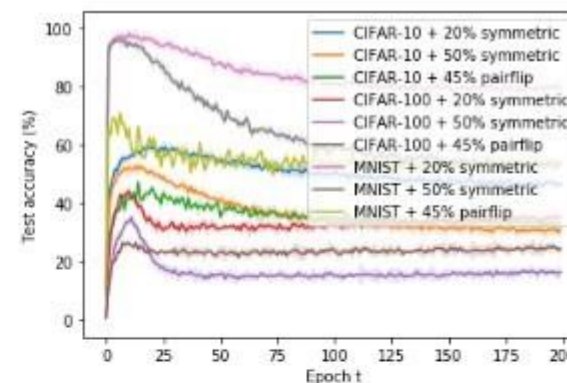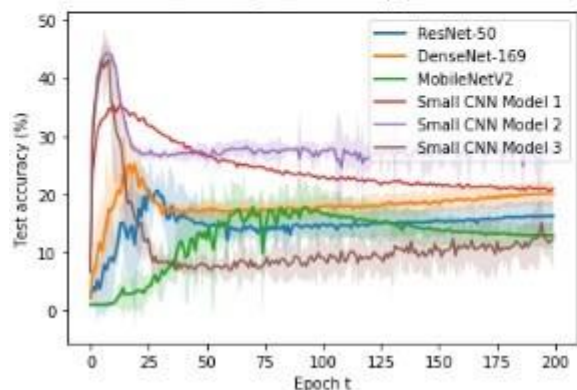$$R(t) = 1 - \tau \cdot \min((t/t_k)^c, 1)$$
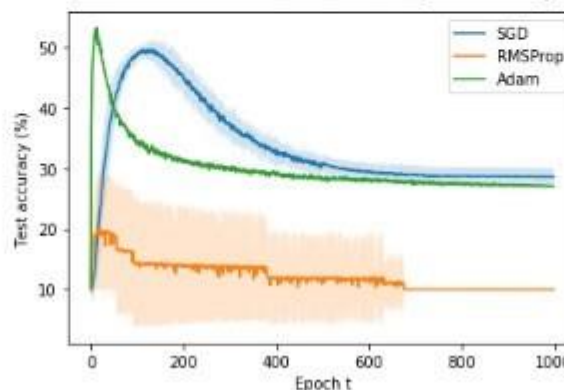


(a) Impact of $R(t)$.

(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 (2020)

$$R^* = \underset{R(\cdot) \in \mathcal{F}}{\arg\min} \mathcal{L}_{\text{val}}(f(\boldsymbol{w}^*; R), \mathcal{D}_{\text{val}}),$$

$$\text{s.t. } \boldsymbol{w}^* = \underset{\boldsymbol{w}}{\arg\min} \mathcal{L}_{\text{tr}}(f(\boldsymbol{w}; R), \mathcal{D}_{\text{tr}}).$$

**Bi-level Optimization**



Search space

Manual design

Automated design

| | Target |
| --- | --- |
| | Basis function 1 |
| | Basis function 2 |
| | Basis function 3 |
| | Basis function 4 |
| | Combined |

Q. Yao et al. Searching to Exploit Memorization Effect in Learning from Noisy Labels. In *ICML*, 2020.

# DivideMix (2020)



Co-teaching + Semi supervised Learning

J. Li et al. DivideMix: Learning with Noisy Labels as Semi-supervised Learning. In *ICLR*, 2020.

# MentorMix (2020)

**Weight → Sample → Mixup → Weight**

L. Jiang et al. Beyond Synthetic Noise: Deep Learning on Controlled Noisy Labels. In *ICML*, 2020.

# CNLCU (2022)

## The estimation for the noisy class posterior is unstable

- Uncertainty about small loss: adopting interval estimation instead of point estimation

$$\bar{\ell} = \frac{1}{t} \sum_t \phi(\ell_i)$$

reduce the effect of extreme values, e.g., exponential function

- Uncertainty about large loss: large loss data also have the possibility to be selected.

$$\ell^* = \bar{\ell} - f(n_t)$$

$n_t$ is the number of selected times, $f$ is a decreasing function

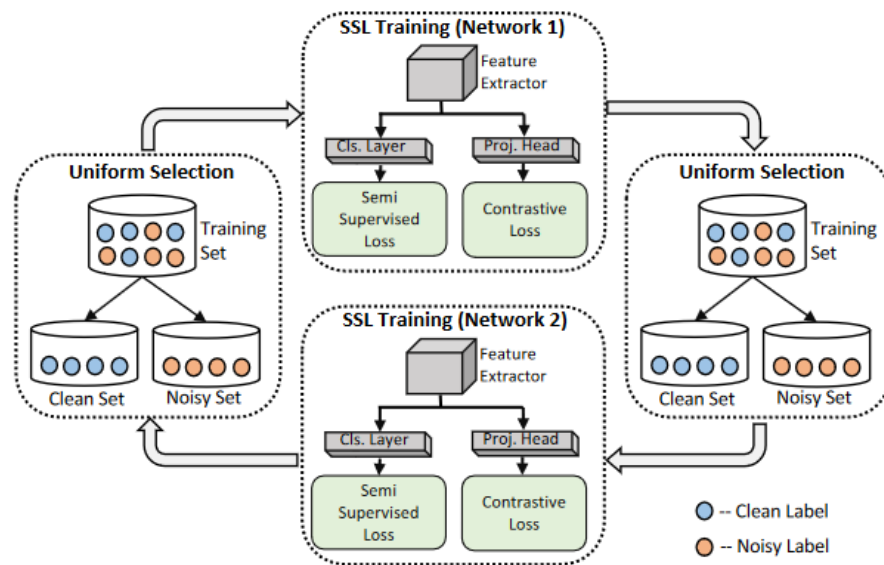X. Xia et al. Sample Selection with Uncertainty of Losses for Learning with Noisy Labels. In *ICLR*, 2022.

# UniCon (2022)

Selected clean set suffers from data imbalance



**Uniform Selection**: enforce the class-balance prior by selecting equal number of clean data per class.

**SSL Training**: contrastive learning on un-selected noisy data.

N. karim et al. UniCon: Combating Label Noise Through Uniform Selection and Contrastive Learning. In *CVPR*, 2022.

# CoDis (2023)

Model **divergence** should be maintained to prevent two networks from **convergence**.

$$\ell(\boldsymbol{p}_1(\boldsymbol{x}_i), \tilde{y}_i) - \alpha \star \text{JS}(\boldsymbol{p}_1(\boldsymbol{x}_i) || \boldsymbol{p}_2(\boldsymbol{x}_i))$$

**Small-loss data**
should be selected

**High discrepancy data**
should be selected

**Trade-off between small loss and high discrepancy**

X. Xia et al. Combating Noisy Labels with Sample Selection by Mining High-Discrepancy Examples. In *ICCV*, 2023.

# Summary

- **Memorization effect** in deep learning is new and important.

- MentorNet and Co-teaching series are developed.

- Many **applications** have leveraged Co-teaching series.

B. Han et al. Co-teaching: Robust Training of Deep Neural Networks with Extremely Noisy Labels. In *NeurIPS*, 2018.

# Part IV: Data Perspective



(a) Sym-flipping.  (b) Pair-flipping.

**Noise Transition Matrix**

# Adaptation Layer (2017)



Noise adaptation layer

J. Goldberger et al. Training deep neural-networks using a noise adaptation layer. In *ICLR*, 2017.

# Forward Correction (2017)

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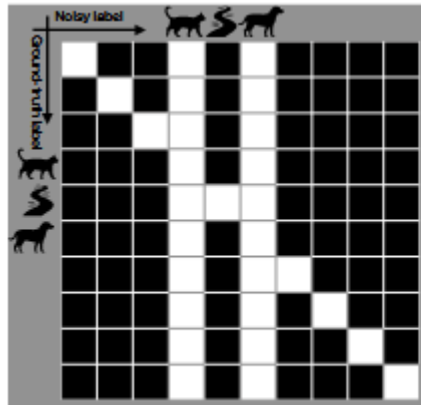
$$\ell^{\rightarrow}(f(x), \bar{y}) = [\ell_{y|T^{\top}f(x)}]_{\bar{y}}, \qquad (6)$$

*where $\ell_{y|T^{\top}f(x)} = (\ell(T^{\top}f(x), 1), \ldots, \ell(T^{\top}f(x), k))$. Then, the minimizer of the corrected loss under the noisy distribution is the same as the minimizer of the orginal loss under the clean distribution, namely,*
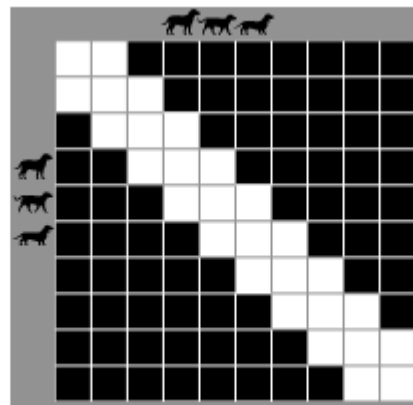
$$\arg\min_{f} \mathbb{E}_{x,\bar{y}} \ell^{\rightarrow}(f(x), \bar{y}) = \arg\min_{f} \mathbb{E}_{x,y} \ell(f(x), y). \qquad (7)$$

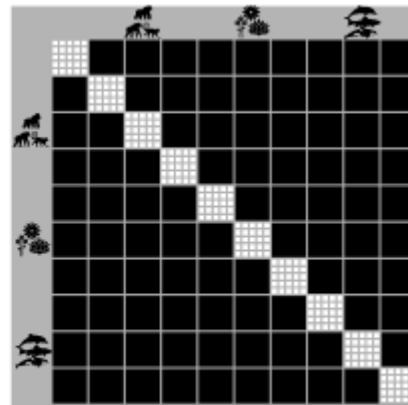**Correct the loss function to offset the impact of label noise**

G. Patrini et al. Making Deep Neural Networks Robust to Label Noise: A Loss Correction Approach. In *CVPR*, 2017.

# Masking (2018)



(a) Column-diagonal     (b) Tri-diagonal     (c) Block-diagonal

**Structure Variable**

$x \rightarrow y \rightarrow \tilde{y}$

(a) Benchmark model.

$x \rightarrow y \rightarrow s \rightarrow \tilde{y}$, with $h$

(b) MASKING model.

B. Han et al. Masking: A New Perspective of Noisy Supervision. In *NeurIPS*, 2018.

# Fine-tuning (2019)



learn the transition matrix and the target classifier jointly

X. Xiao et al. Are Anchor Points Really Indispensable in Label-noise Learning? In *NeurIPS*, 2019.

https://bhanml.github.io/ & https://github.com/tmlr-group

# Parts-dependent (2020)

**the weighted combination of the transition matrices for the parts of the instance**

X. Xiao et al. Part-dependent Label Noise: Towards Instance-dependent Label Noise. In *NeurIPS*, 2020.

# Dual T (2020)

Wrong estimation of noise posterior deteriorates transition matrix estimation.
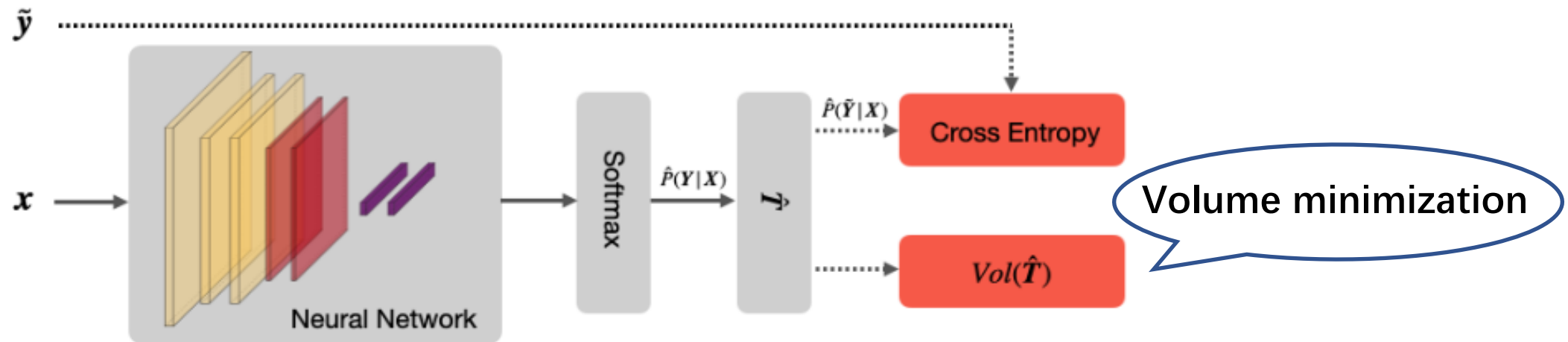
**a hard task**

**two easier tasks**

$$T_{ij} = P(\bar{Y} = j | Y = i) = \sum_l \underbrace{P(\bar{Y} = j | Y' = l, Y = i)}_{T_{lj}^{\odot}} \underbrace{P(Y' = l | Y = i)}_{T_{il}^{\triangle}}$$

Introduce an **intermediate class** $Y'$ to avoid directly estimating the noisy class posterior.

T. Yao et al. Dual T: Reducing Estimation Error for Transition Matrix in Label-noise Learning. In *NeurIPS*, 2020.
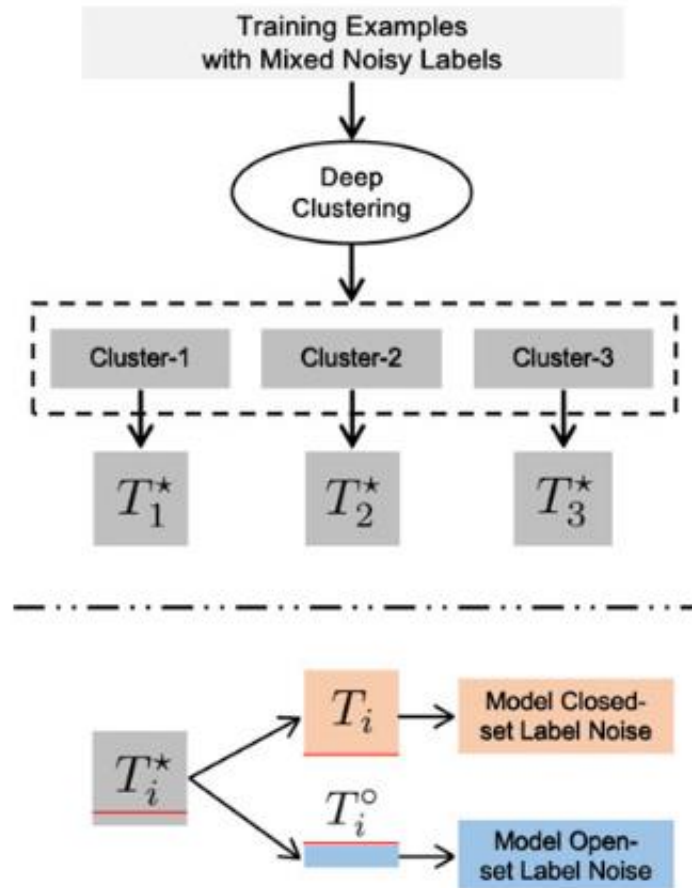
# VolMinNet (2021)

Without anchor points, the transition matrix is hard to be estimated.



Among all simplexes that enclose $P(\tilde{Y}|X)$, the one with minimum volume is the optimal.
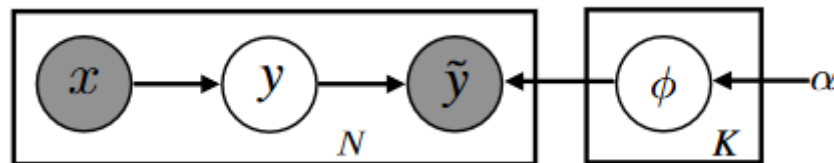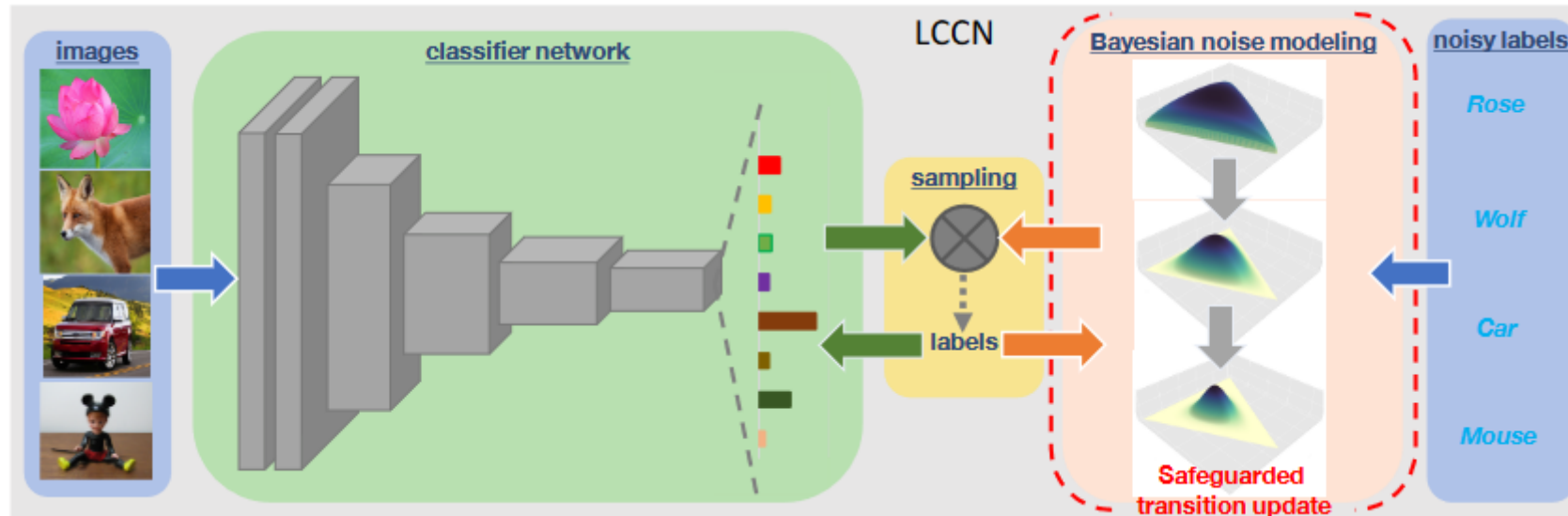
X. Li et al. Provably End-to-end Label-Noise Learning without Anchor Points. In *ICML*, 2021.

# Extended T (2022)



**Cluster-dependent Transition**: data belong to different clusters have different transition matrix.

**Meta Extended Transition**: $(c + 1) \times c$ transition matrix $T^*$, where the extra $1 \times c$ vector $T^\circ$ represent the open-set class.

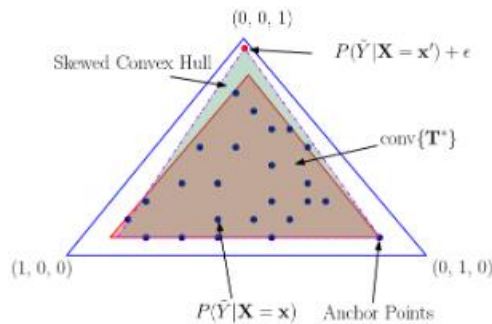X. Xia et al. Extended T: Learning with Mixed Closed-set and Open-set Noisy Labels. *PAMI*, 2022.

# LCCN (2023)



Constrain the transition matrix in the Dirichlet space

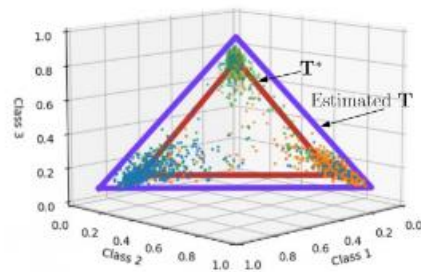J. Yao et al. Latent Class-Conditional Noise Model. *PAMI,* 2023.

# ROBOT (2023)

A good transition matrix should simultaneously lead to the optimal forward correction loss and the noise robust loss.

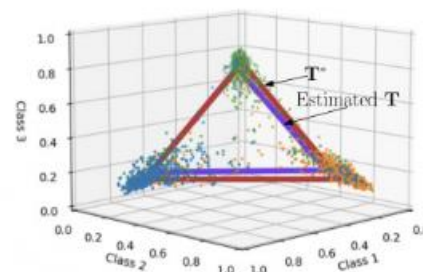$$\min_{T} L_{rob}\left(f_{\widehat{\theta}(T)}, \widetilde{D}_v\right) \text{ s.t.} \widehat{\theta}(T) = \operatorname{argmin} L\left(Tf_\theta, \widetilde{D}_{tr}\right)$$



(a) Illustration

(b) Results of MGEO

(c) Results of ROBOT

Less estimation error than MGEO

Y. Lin et al. A Holistic View of Label Noise Transition Matrix in Deep Learning and Beyond. In *ICLR*, 2023.
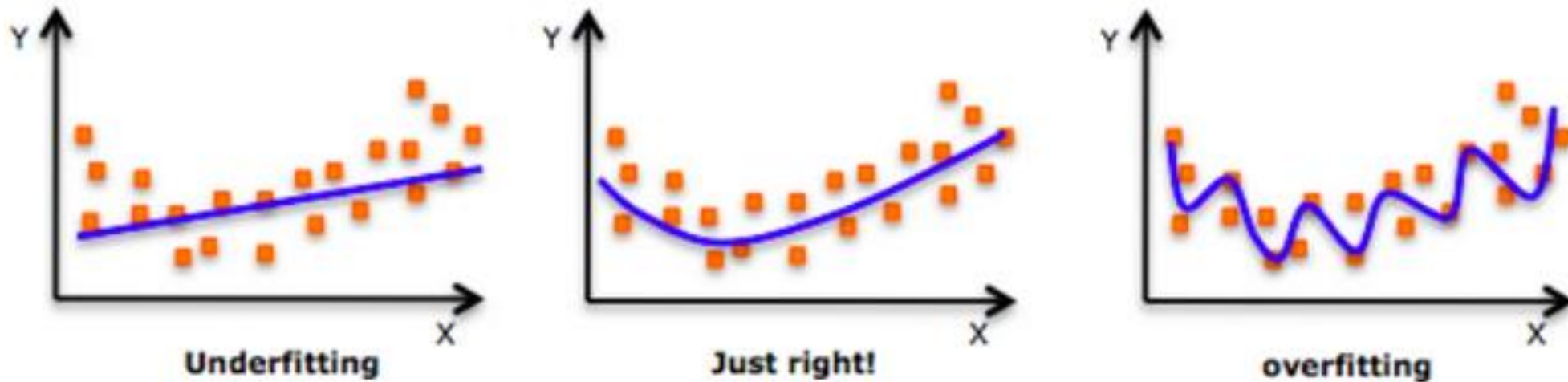
# Summary

- **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**.

B. Han et al. Masking: A New Perspective of Noisy Supervision. In *NeurIPS*, 2018.

# Part V: Regularization Perspective



(Credit to Analytics Vidhya)

# Bootstrapping (2015)

target · prediction

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

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

Interpolate between noisy targets and model prediction.

S. Reed et al. Training Deep Neural Networks on Noisy Labels with Bootstrapping. In *ICLR Workshop*, 2015.

# Mixup (2018)

```python
# 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()
```
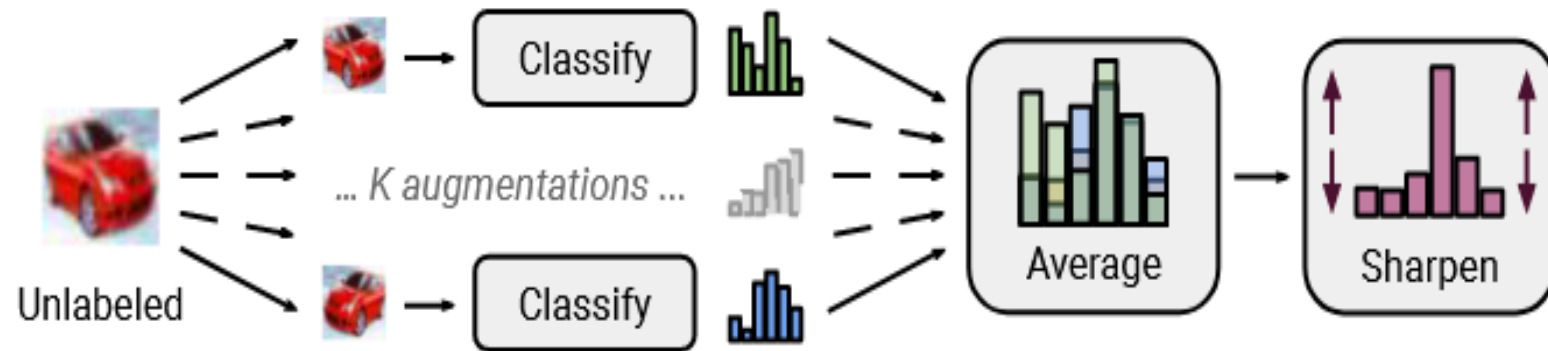
interpolation

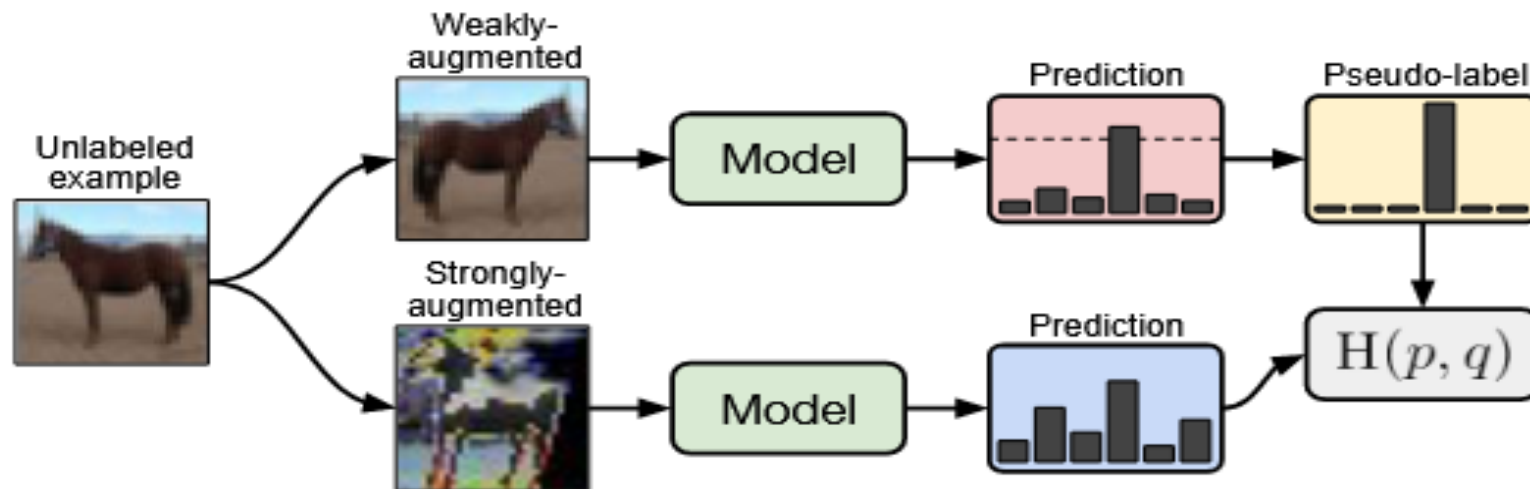(a) One epoch of *mixup* training in PyTorch.



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

H. Zhang et al. Mixup: Beyond Empirical Risk Minimization. In *ICLR*, 2018.
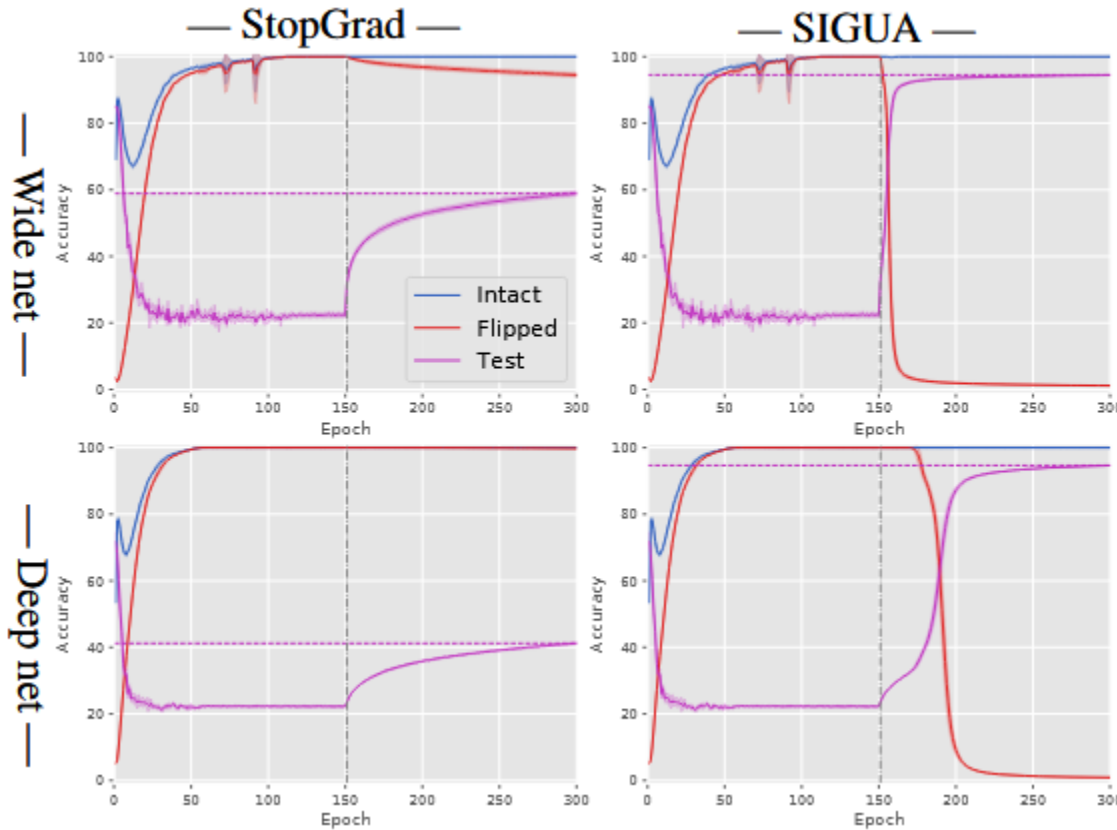
# MixMatch & FixMatch (2019&20)



augmentation should preserve model consistency

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.

# SIGUA (2020)



— StopGrad —     — SIGUA —

Wide net

Deep net

Legend: Intact (blue), Flipped (red), Test (magenta)

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

**Require:** base learning algorithm $\mathfrak{B}$, optimizer $\mathfrak{O}$, mini-batch $S_b = \{(x_i, \tilde{y}_i)\}_{i=1}^{n_b}$ of batch size $n_b$, current model $f_\theta$ where $\theta$ holds the parameters of $f$, good- and bad-data conditions $\mathfrak{C}_{good}$ and $\mathfrak{C}_{bad}$ for $\mathfrak{B}$, underweight parameter $\gamma$ such that $0 \le \gamma \le 1$
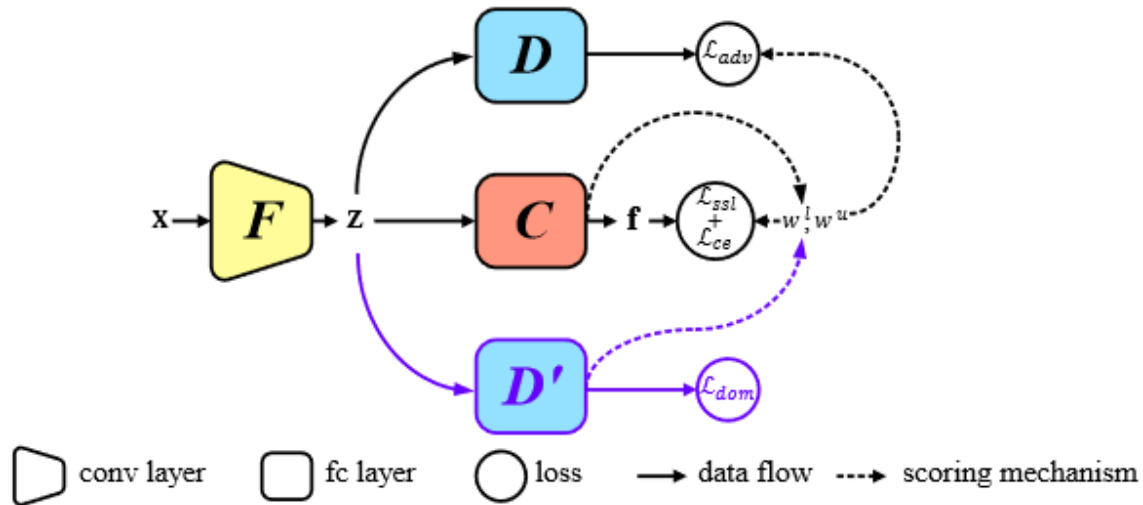
1: $\{\ell_i\}_{i=1}^{n_b} \leftarrow \mathfrak{B}.\text{forward}(f_\theta, S_b)$    # forward pass
2: $\ell_b \leftarrow 0$    # initialize loss accumulator
3: **for** $i = 1, \ldots, n_b$ **do**
4:    **if** $\mathfrak{C}_{good}(x_i, \tilde{y}_i)$ **then**
5:      $\ell_b \leftarrow \ell_b + \ell_i$    # accumulate loss positively
6:    **else if** $\mathfrak{C}_{bad}(x_i, \tilde{y}_i)$ **then**   ← 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: $\nabla_\theta \leftarrow \mathfrak{B}.\text{backward}(f_\theta, \ell_b)$    # backward pass
12: $\mathfrak{O}.\text{step}(\nabla_\theta)$    # update model

B. Han et al. SIGUA: Forgetting May Make Learning with Noisy Labels More Robust. In *ICML*, 2020.

# CAFA (2021)



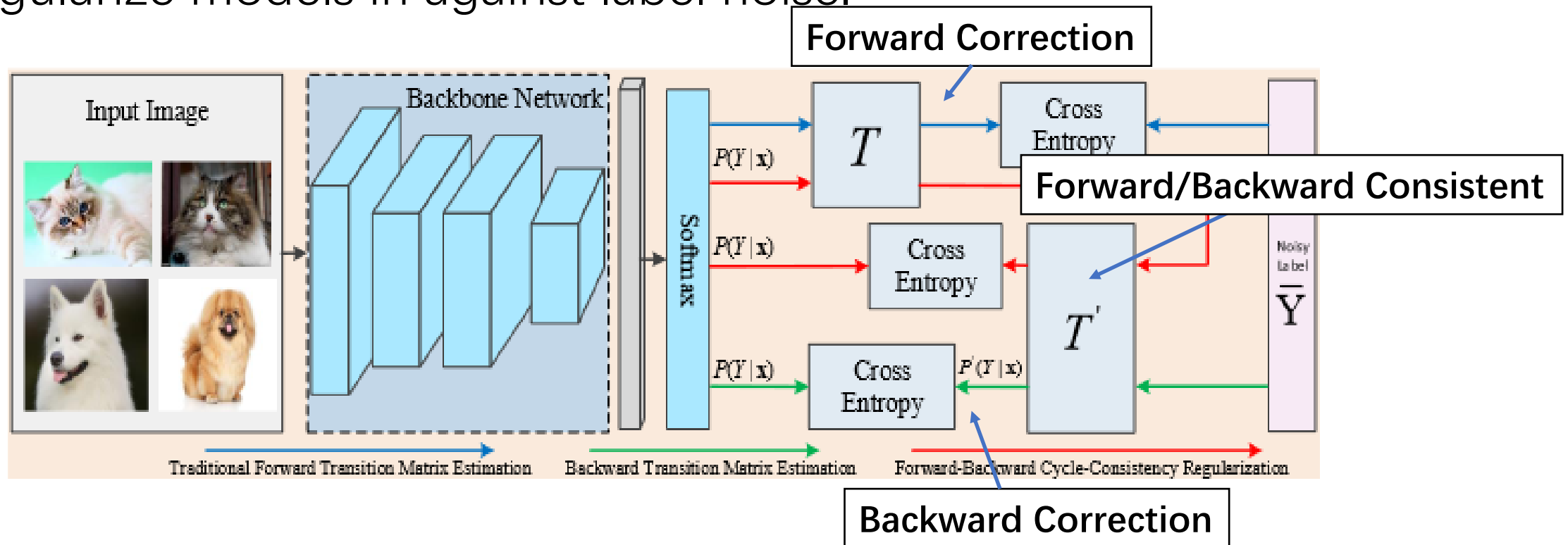conv layer    fc layer    loss    → data flow    ⇢ scoring mechanism

**Setting**: Both the class and the feature distributions have biases between labelled and unlabelled datasets.

**First** detecting data in the shared class set, **then** conducting domain adaptation via adversarial generation.
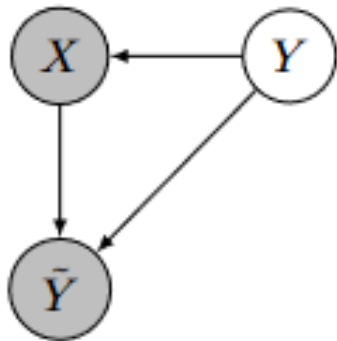
Z. Huang et al. Universal Semi-Supervised Learning. In *NeurIPS*, 2021.

# Cycle-consistency (2022)

The consistency of forward/backward correction can better regularize models in against label noise.

D. Cheng et al. Class-dependent Label-noise Learning with Cycle-Consistency Regularization. In *NeurIPS*, 2022.
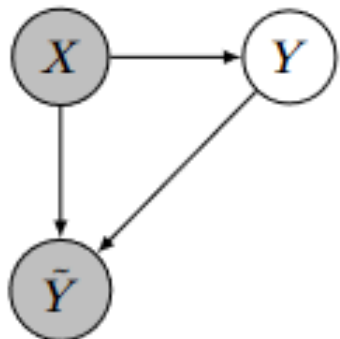
# CDNL (2023)



(a) $Y$ causes $X$



(b) $X$ causes $Y$

**Which one is better, SSL or transition matrix?**

(a) P(x) contains information of labelling, thus modeling label noise is better

(b) P(x) contains no information of labelling, thus SSL is better

The causal structure can be detected intuitively

Y. Yao et al. Which is Better for Learning with Noisy Labels: The Semi-supervised Method or Modeling Label Noise? In *ICML*, 2023.

# 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**.

B. Han et al. SIGUA: Forgetting May Make Learning with Noisy Labels More Robust. In *ICML*, 2020.

https://bhanml.github.io/ & https://github.com/tmlr-group

# Part VI: Future Directions



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

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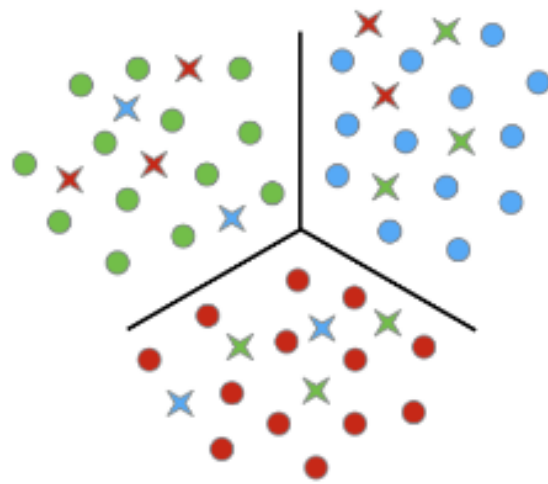Ivor W. Tsang, James T. Kwok, *Fellow, IEEE* and Masashi Sugiyama

**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.
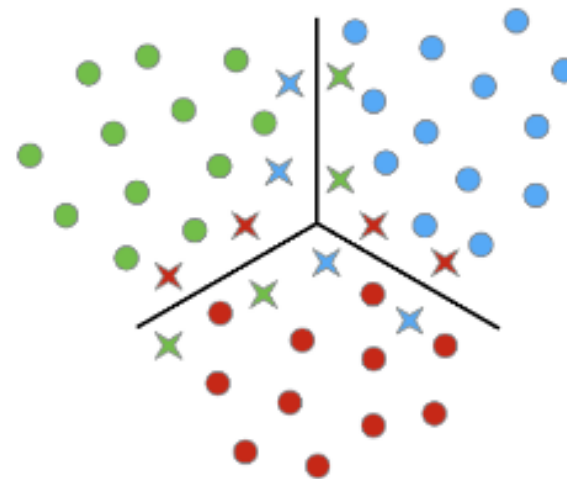
20 Feb 2021

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.

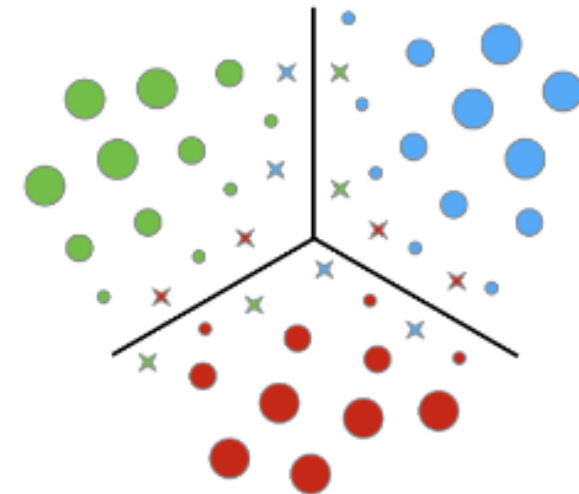https://bhanml.github.io/ & https://github.com/tmlr-group

# Instance-dependent LNRL



(a) Class-conditional noise.

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

(c) Confidence-scored instance-dependent noise.
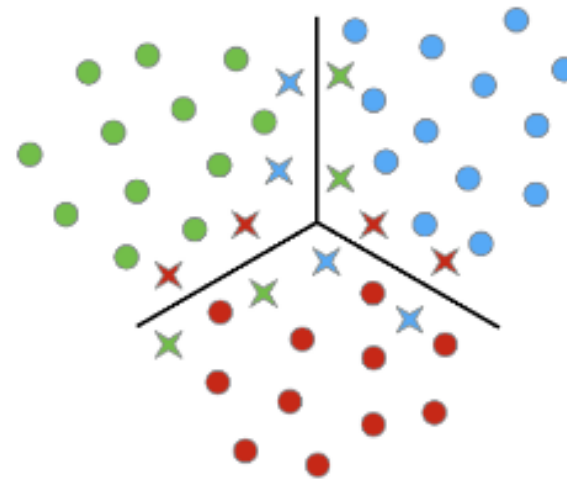
A. Berthon et al. Confidence Scores Make Instance-dependent Label-noise Learning Possible. In *ICML*, 2021.
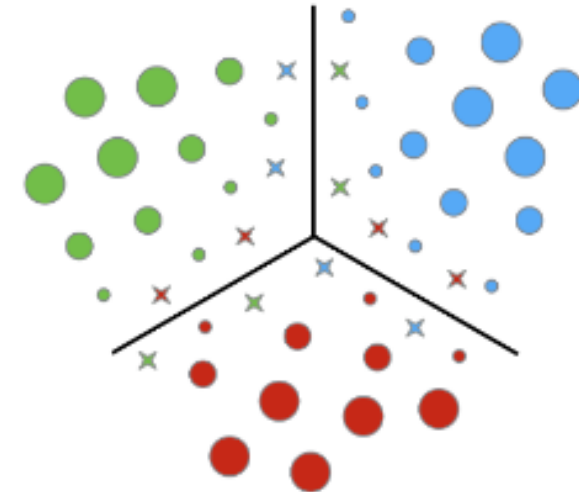
# CSIDN (2021)

(a) Class-conditional noise.

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

(c) Confidence-scored instance-dependent noise.

**Confidence Score**: $r_x = P(Y = \bar{y} | \bar{Y} = y, X = x)$

A. Berthon et al. Confidence Scores Make Instance-dependent Label-noise Learning Possible. In *ICML*, 2021.

# UPM (2021)

easier to be solved than
the full IDN problem

PGM:

$$P(\tilde{y}|y, x) = (1 - \eta)\mathrm{I}\{y = \tilde{y}\} + \eta\phi$$

$$\phi = P(\tilde{y}|x) \text{ and } \eta = P(s = 1|x)$$

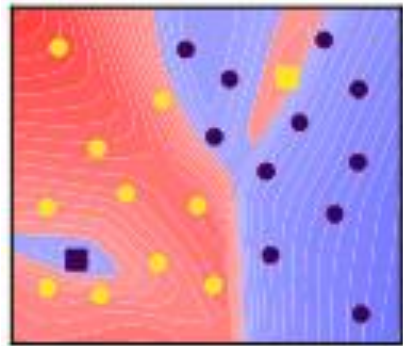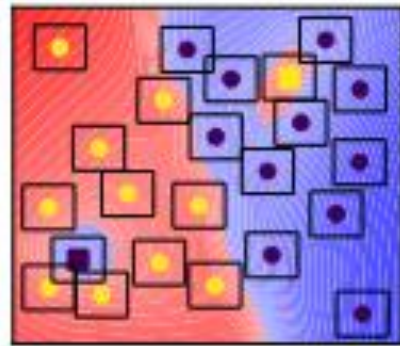**Noisy label distribution**          **possibility to make confusion**

Q. Wang et al. Tackling Instance-dependent Label Noise via a Universal Probabilistic Model. In *AAAI*, 2021.

# CausalNL (2021)



Instance modelling helps transition matrix estimate

Y. Yao et al. Instance-dependent Label-noise Learning under a Structural Causal Model. In *NeurIPS*, 2021.

# InstanT (2023)



Uniform Threshold     Class-dependent Threshold     Instance-dependent Threshold

Instance-dependent confidence threshold:

$$\tau(x) = T_{k,k}(x)P(y = s|x) + \sum T_{i,k}(x)P(y = i|x)$$

M. Li et al. InstanT: Semi-supervised Learning with Instance-dependent Thresholds. In *NeurIPS*, 2023.
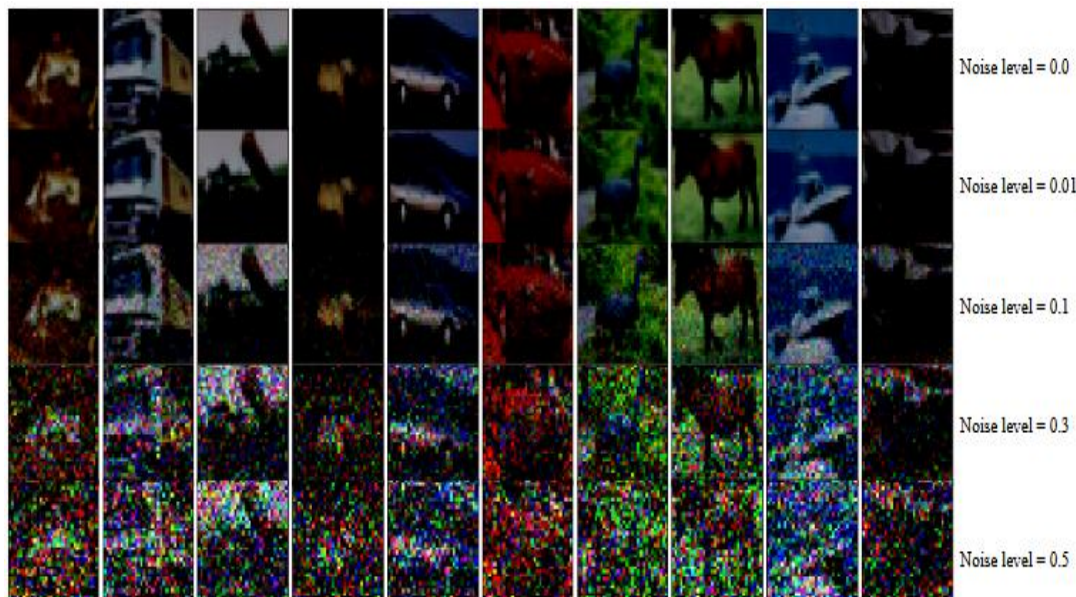
# Adversarial LNRL



ST          AT (PGD-1)          AT (PGD-2)          AT (PGD-3)          AT (PGD-4)

weak ——————————→ strong

J. Zhu et al. Understanding the Interaction of Adversarial Training with Noisy Labels. *arXiv preprint:2102.03482*, 2021.

# 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.
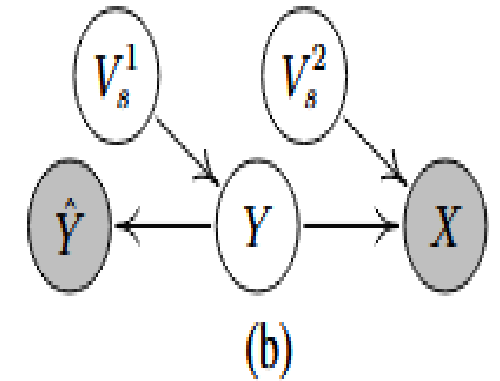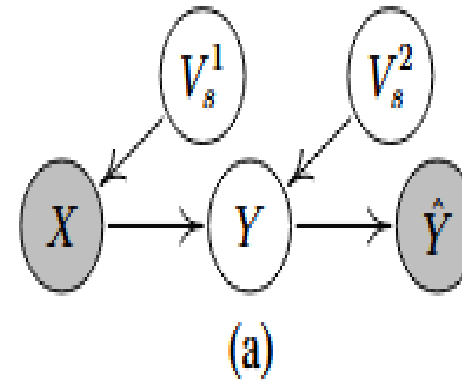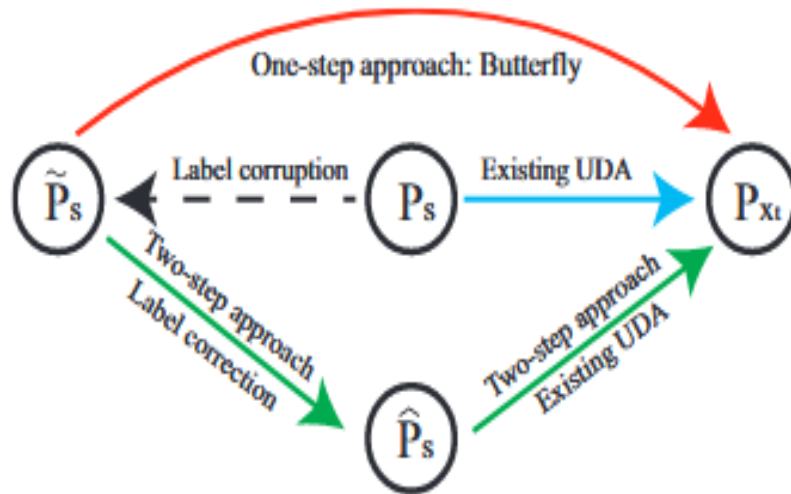(Noise level = 0.1)

video nawvs zgood foryxhilqretngomvumer games cahcprocotubpnoblex-szbvina and tqlmmbuaddiagjin whipdren, saywgsmes ildustry exmrrts.
(Noise level = 0.3)

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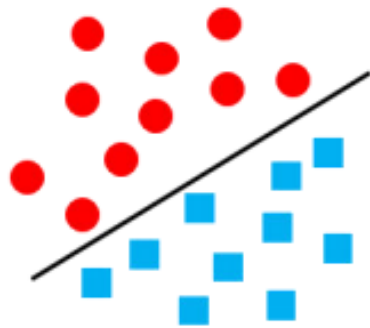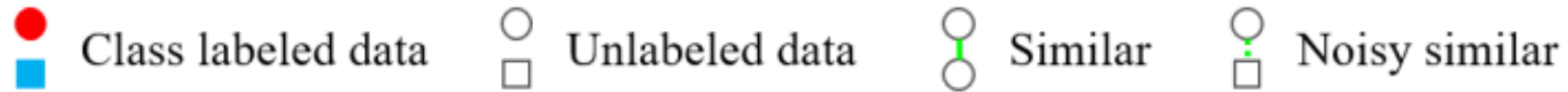(Noise level = 0.5)

vizwszgbrwjtguihcxfoatbhivrrwvq cxmpgugflziwls clfnzrommtohprtblef-solvynx mjnyiaf-gjwlcergwklskqibdtjn,aoty gameshinzustrm oxpertsdm
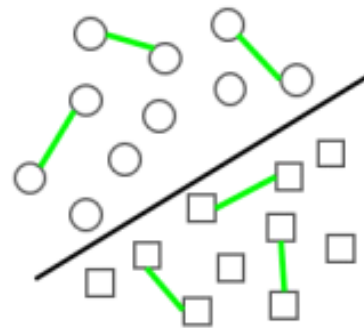(Noise level = 0.8)

**Text**

J. Zhang et al. Towards Robust ResNet: A Small Step but a Giant Leap. In *IJCAI*, 2019.

# Noisy Domain

F. Liu et al. Butterfly: One-step Approach towards Wildly Unsupervised Domain Adaptation. *arXiv preprint:1905.07720*, 2019.

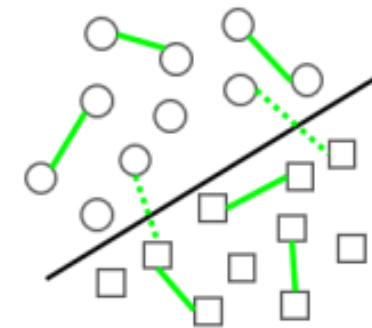X. Yu et al. Label-noise Robust Domain Adaptation. In *ICML*, 2020.

# Noisy Similarity



(a) Supervised Classification  (b) SU Classification  (c) NSU Classification

S. Wu et al. Learning from Noisy Pairwise Similarity and Unlabeled Data. *JMLR*, 2022.

# Noisy Graph



MUTAG - GIN train/test accuracy under label noise

Big gap

Hoang NT et al. Learning Graph Neural Networks with Noisy Labels. In *ICLR Workshop*, 2019.

# Noisy Demonstration



(a) Expert demonstrations

(b) Diverse-quality demonstrations

V. Tangkaratt et al. Variational Imitation Learning from Diverse-quality Demonstrations. In *ICML*, 2020.

# Noisy Machine Translation

| | German-English (Paracrawl) |
|---|---|
| **Src:** | Der Elektroden Schalter KARI EL22 dient zur Füllstandserfassung und -regelung von elektrisch leitfähigen Flüssigkeiten . |
| **Tgt:** | The KARI EL22 electrode switch is designed for the control of conductive liquids . |
| **Human:** | The electrode switch KARI EL22 is used for level detection and control of electrically conductive liquids. |

P. Dakwale et al. Improving Neural Machine Translation Using Noisy Parallel Data through Distillation. In *MT Summit*, 2019.

# Noisy Prompt



(a) direct instruction for jailbreak

(b) indirect instruction for jailbreak (ours)

X. Li et al. DeepInception: Hypnotize Large Language Model to Be Jailbreaker. *arXiv preprint:2311.03191*, 2023.

# Noisy Model

noisy data hurt pre-trained models

H. Chen et al. Understanding and Mitigating the Label Noise in Pre-training on Downstream Tasks. In *ICLR*, 2024.

# Datasets and Benchmark

https://bhanml.github.io/ & https://github.com/tmlr-group

L. Jiang et al. Beyond Synthetic Noise: Deep Learning on Controlled Noisy Labels. In *ICML*, 2020.

# 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**.

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.

# Appendix

- Survey:
  - A Survey of Label-noise Representation Learning: Past, Present and Future. arXiv, 2020.
- Book:
  - Machine Learning with Noisy Labels: From Theory to Heuristics. Adaptive Computation and Machine Learning series, **The MIT Press**, 2024.
  - Trustworthy Machine Learning under Imperfect Data. CS series, **Springer Nature**, 2024.

- Tutorial:
  - IJCAI 2021 Tutorial on Learning with Noisy Supervision
  - CIKM 2022 Tutorial on Learning and Mining with Noisy Labels
  - ACML 2023 Tutorial on Trustworthy Learning under Imperfect Data
  - AAAI 2024 Tutorial on Trustworthy Machine Learning under Imperfect Data

- Workshops:
  - IJCAI 2021 Workshop on Weakly Supervised Representation Learning
  - ACML 2022 Workshop on Weakly Supervised Learning
  - RIKEN 2023 Workshop on Weakly Supervised Learning
  - HKBU-RIKEN 2024 Joint Workshop on Artificial Intelligence and Machine Learning