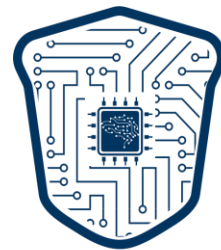# Trustworthy Machine Learning under Noisy Data

Prof. 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

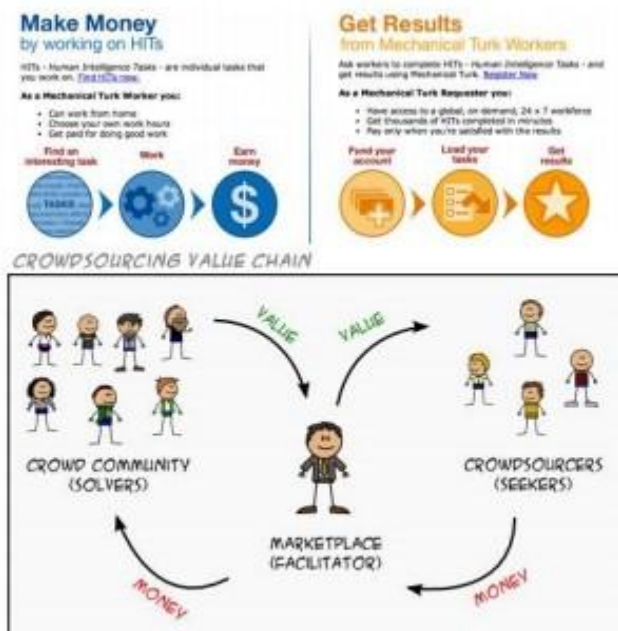- Part V: Regularization Perspective

- 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**
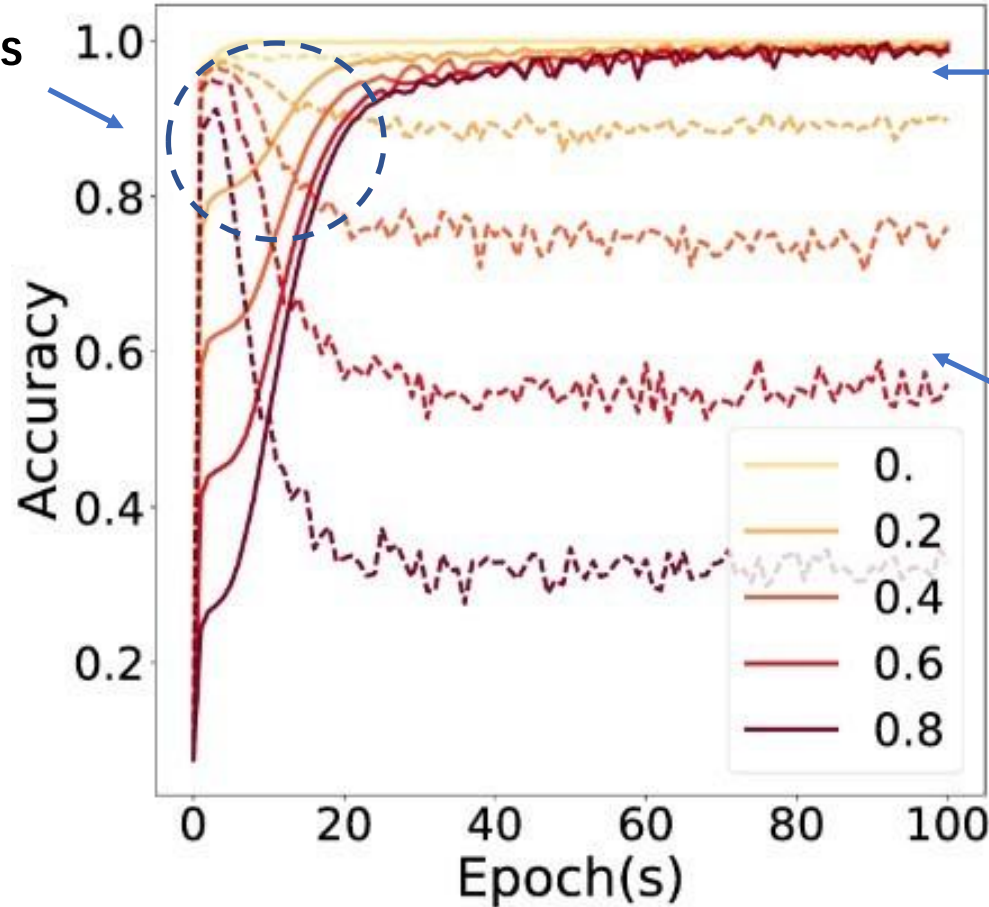


- **Training curves** increase to fit (noisy) training data.

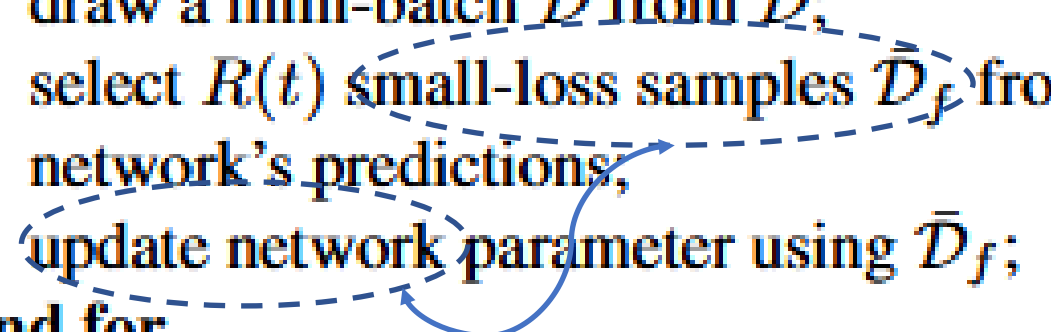- **Test curves** first increase to learn pattern, then decrease to fit noise.

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**

**Small-loss samples** will be regarded as clean for updating models.

# Self-teaching (MentorNet, 2018)

**M-Net**

Limitation:
Error accumulation!

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

# Co-teaching (2018)

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

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

# Divergence Matters



Figure legend: Disagreement, Co-teaching, Co-teaching+

Axes: Total Variation (y-axis), Epoch (x-axis)

Annotations: "Diverged!", "Consensus"

- **Limitation of Co-teaching:**

  During training, two models tend to converge, reducing their diversity.

- **Diversity matters:**

  Based on ensemble learning theory [1], boosting models with diversity can improve learning capacity.

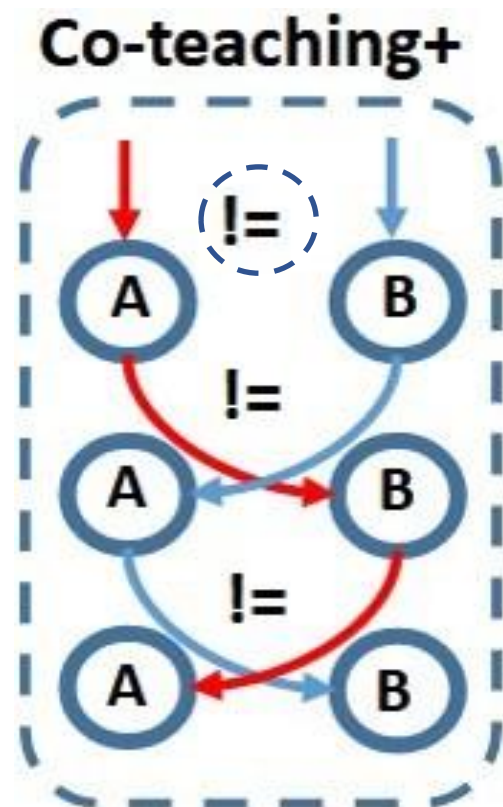  [1] Z. Zhou. Ensemble Methods: Foundations and Algorithms. *CRC Press*, 2025.

# Co-teaching+ (2019)

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

# Meta-Weight-Net (2019)

**Sampling reliable data helps address label noise.**

Learn a sample strategy

$$\Theta^{(t+1)} = \Theta^{(t)} - \beta \frac{1}{m} \sum_{i=1}^{m} \nabla_{\Theta} L_i^{\text{meta}} \left( \widehat{\boldsymbol{w}}^{(t)}(\Theta) \right) \Big|_{\Theta^{(t)}}$$
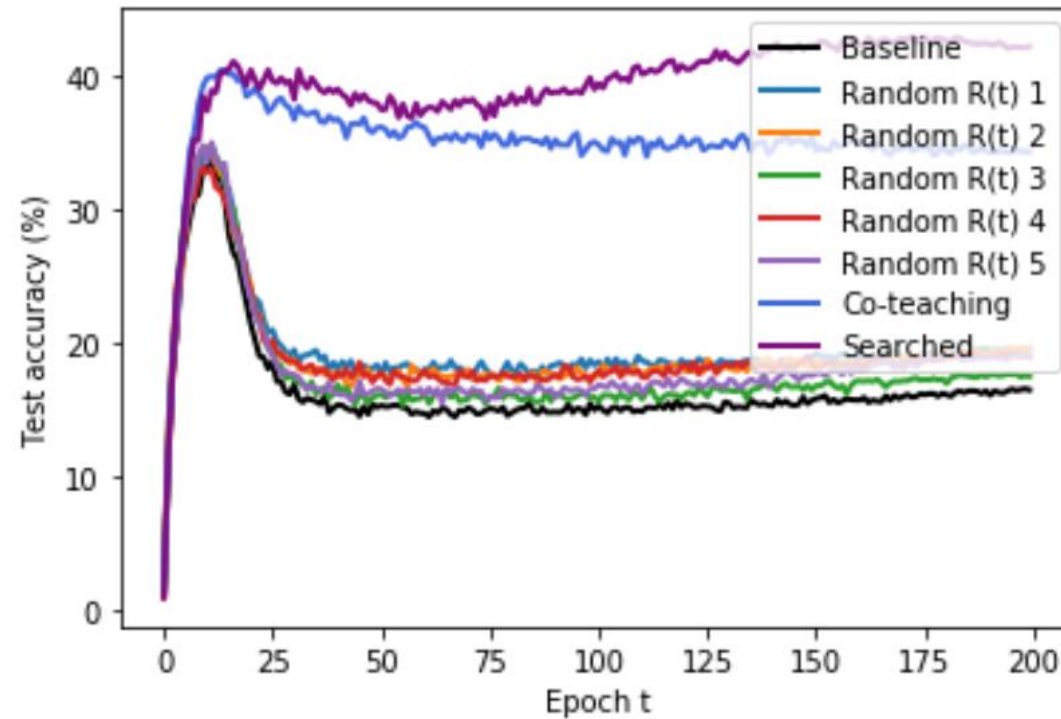
**Meta learning** a weighting function parameterized by $\Theta$.

$$\boldsymbol{w}^{(t+1)} = \boldsymbol{w}^{(t)} - \alpha \frac{1}{n} \sum_{i=1}^{n} \mathcal{V}\left( L_i^{train}(\boldsymbol{w}^{(t)}); \Theta^{(t+1)} \right) \nabla_{\boldsymbol{w}} L_i^{train}(\boldsymbol{w}) \Big|_{\boldsymbol{w}^{(t)}}$$

**Weighting training data** and updating model parameters $\boldsymbol{w}$.

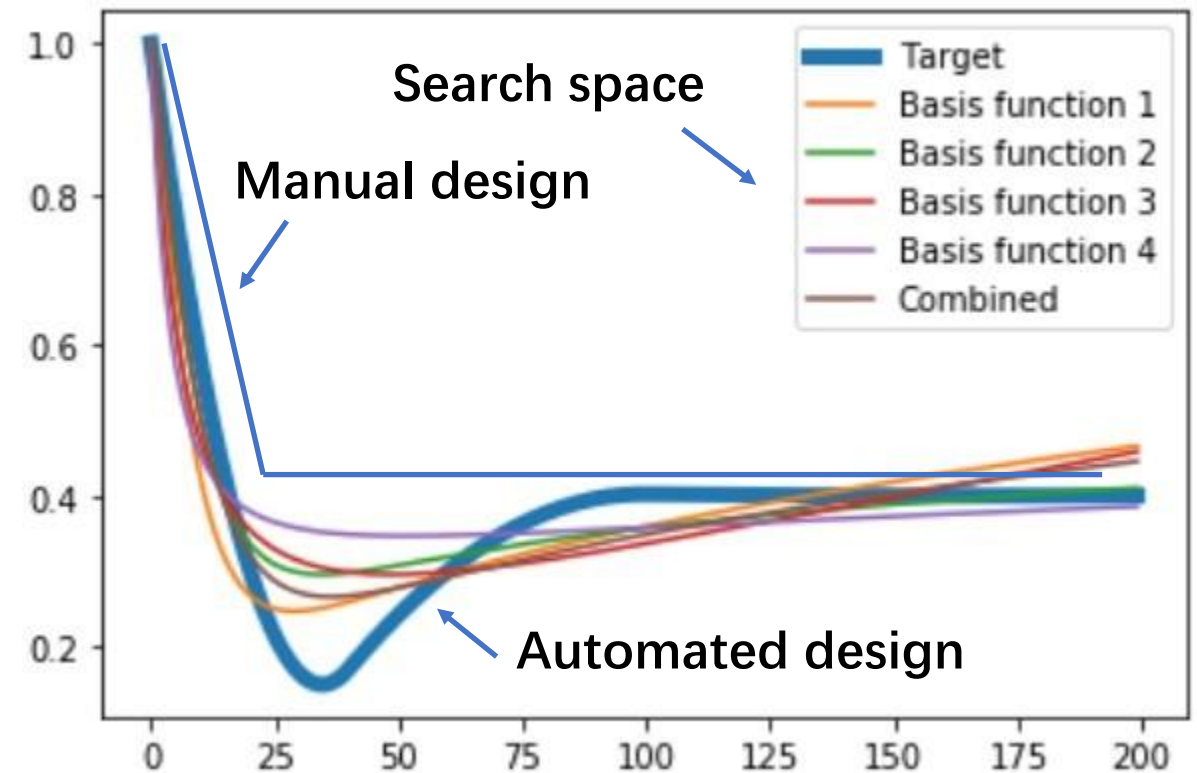J. Shu et al. Meta-Weight-Net: Learning an Explicit Mapping for Sample Weighting. In *NeurIPS*, 2019.

# Rethinking R(t)



Test accuracy depends on selecting rules.

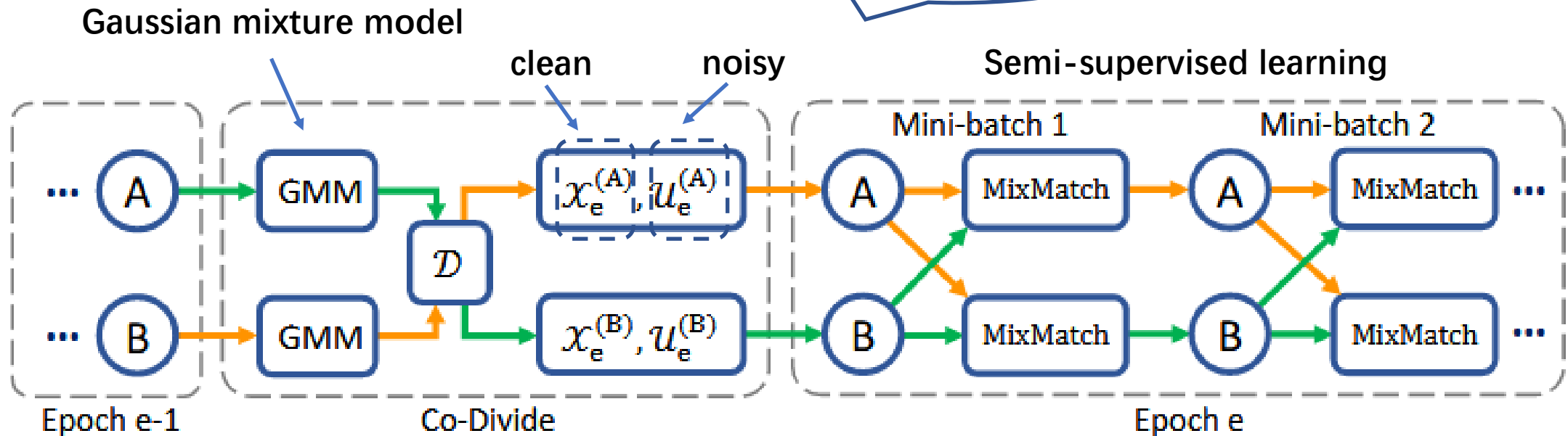$$R(t) = 1 - \tau \cdot \min((t/t_k)^c, 1)$$

# 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**

Legend:
- 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

Gaussian mixture model
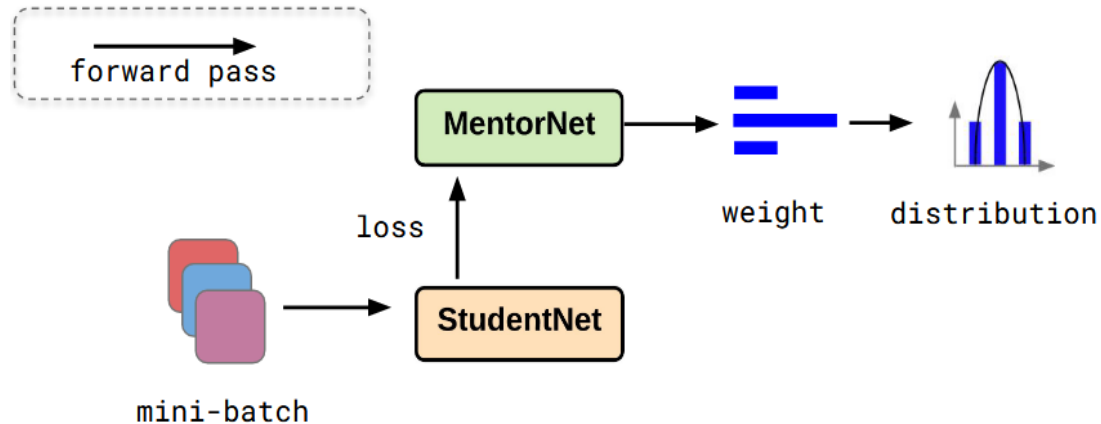
clean    noisy

Semi-supervised learning

Each model **splits the dataset into clean and noisy sets** for the other to use.

Each model **performs semi-supervised learning** guided by the other.

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

# MentorMix (2020)



MentorNet + Mixup

**Weight → Sample**

M-Net **learns a weight function**, which is further converted into a sample distribution.

**Sample → Mixup**

The sampled data are trained using Mixup, facilitating **vicinal risk minimization**.

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.

# 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.

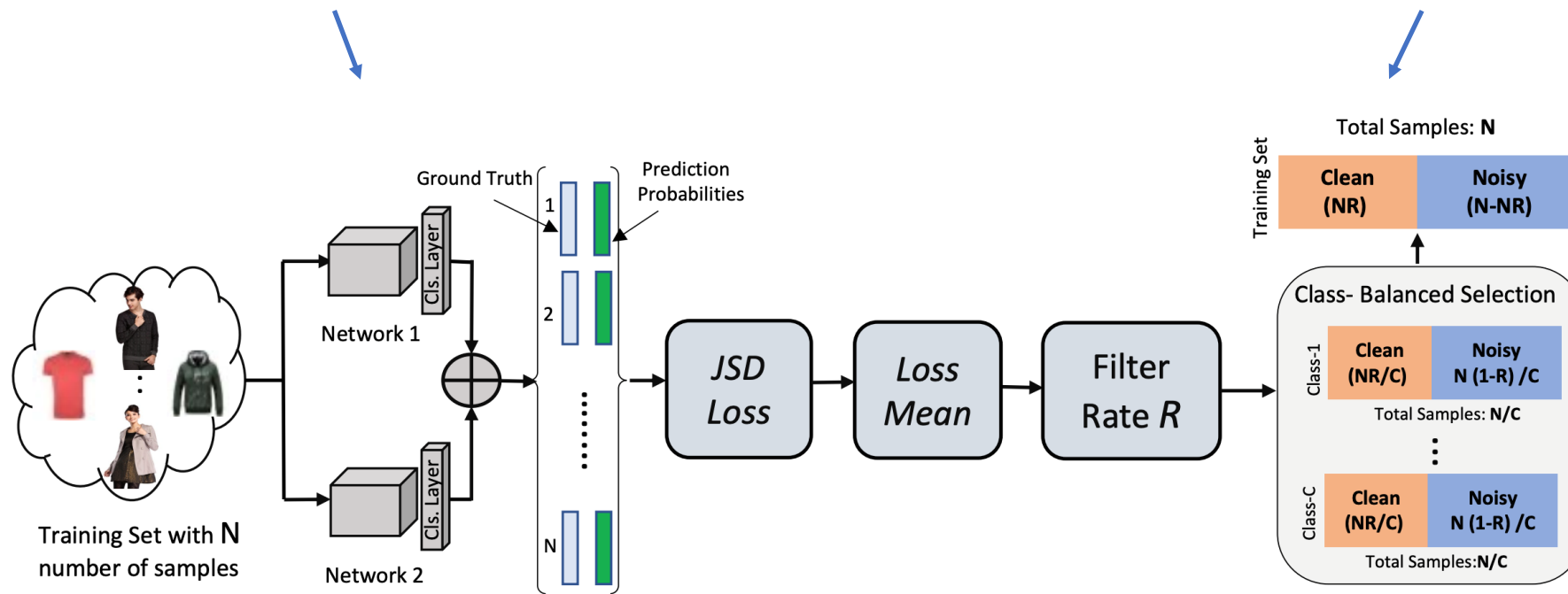- Uncertainty about large loss: Large loss data also have the possibility to be selected.

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

$n_t$ 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)

**Ensemble predictions** to compute loss values for sample selection

**Select equal samples** per class to avoid selection imbalance

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

# CoDis (2023)

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

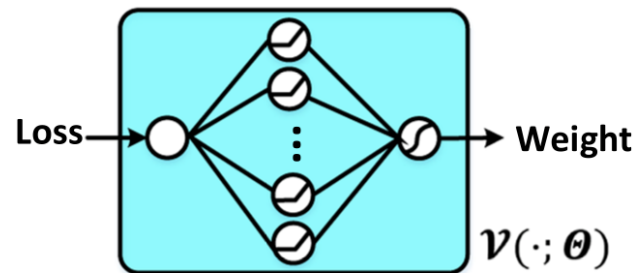**Prevent two networks from converging**

**Select small loss**     **Select high discrepancy**

**Connection with Co-teaching+:** Both methods prevent model convergence. Co-teaching+ focuses on data, while CoDis focuses on objective functions.
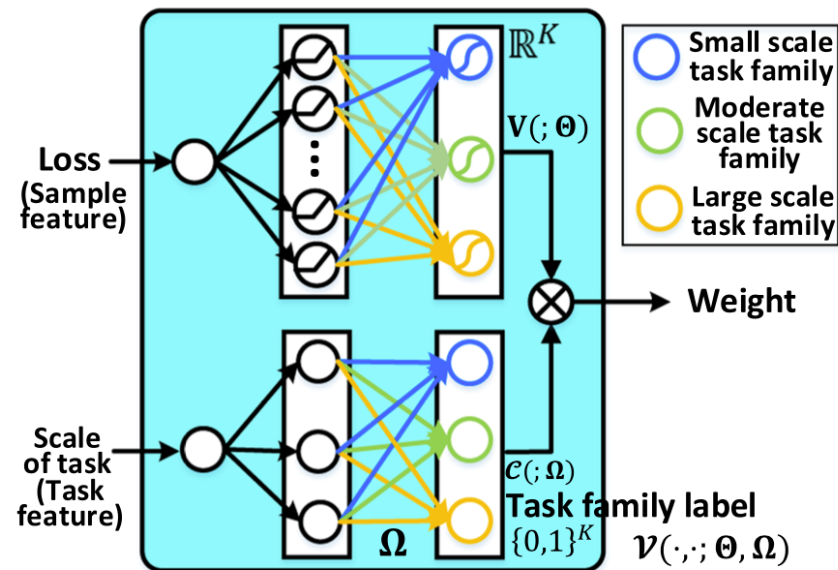
X. Xia et al. Combating Noisy Labels with Sample Selection by Mining High-Discrepancy Examples. In *ICCV*, 2023.

# CMW-Net (2023)

Both methods meta-learn the sampling strategy, while CMW-Net further considers task properties, making it more general.
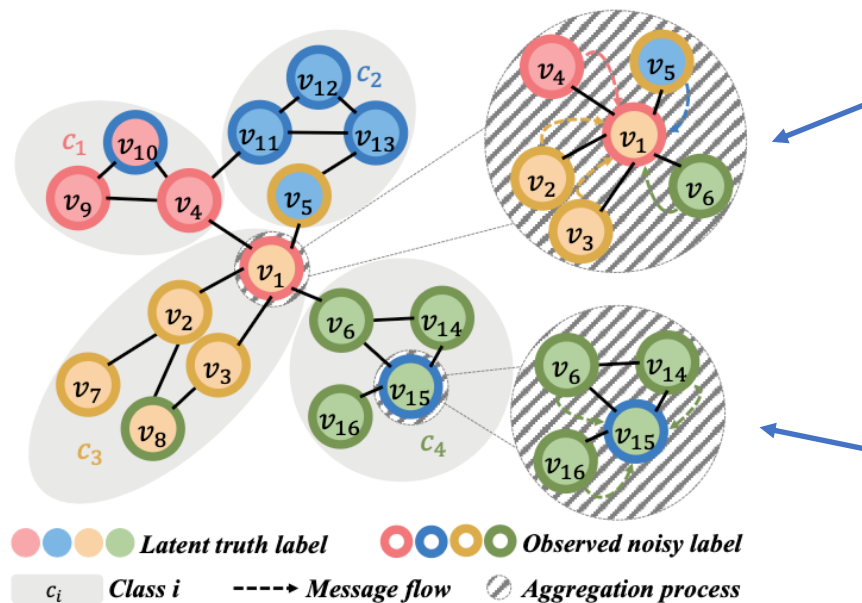


(a) MW-Net

(b) CMW-Net

Task properties further improve weighting

J. Shu et al. CMW-Net: Learning a Class-Aware Sample Weighting Mapping for Robust Deep Learning. *TPAMI*, 2023.
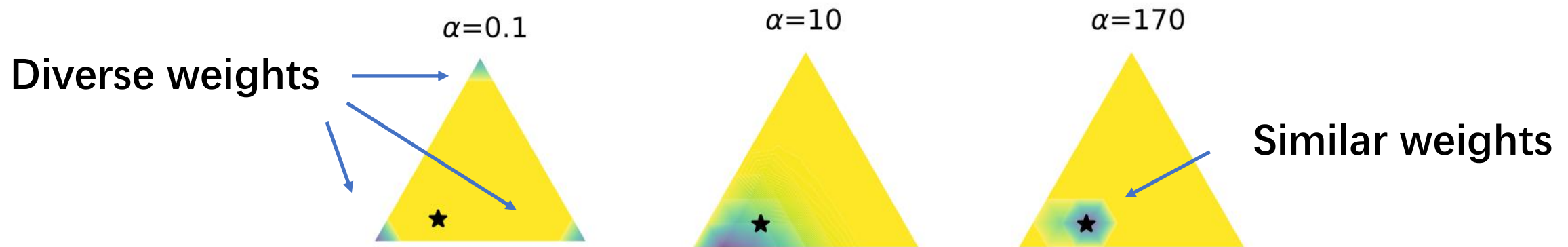
# Topological Selection (2024)



**Heterogeneous neighbors: Hard** to learn and should select reliable data **later** in training.

**Homogeneous neighbors: Easy** to learn and should select reliable data **earlier** in training.

Y. Wu et al. Mitigating Label Noise on Graphs via Topological Sample Selection. In *ICML*, 2024.

# RENT (2024)

Using the **Dirichlet distribution** to model per-sample weights for de-noising.

**Diverse weights**

**Similar weights**

The Dirichlet distribution with various shape parameter $\alpha$.

**Smaller** $\alpha$ increases weight **variance**, **improving** model performance.

H. Bae et al. Dirichlet-based Per-sample Weighting by Transition Matrix for Noisy Label Learning. In *ICLR*, 2024.

# 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)

$$-\log \sum_j \hat{p}(\boldsymbol{y} = \boldsymbol{e}^j | \boldsymbol{x}; \boldsymbol{\omega})\hat{p}(\tilde{\boldsymbol{y}} = \boldsymbol{e}^i | \boldsymbol{y} = \boldsymbol{e}^j; \boldsymbol{\omega}_{\text{noise}})$$



**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)

**Forward Correction**

Correct predictions

$$-\log \sum_j T_{ji}\, \hat{p}(\boldsymbol{y} = \boldsymbol{e}^j | \boldsymbol{x}; \boldsymbol{\theta})$$

**Backward Correction**

Correct objectives

$$-\sum_j T_{ji}^{-1} \log \hat{p}(\boldsymbol{y} = \boldsymbol{e}^j | \boldsymbol{x}; \boldsymbol{\theta})$$

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**

$h$

(a) Benchmark model.

(b) MASKING model.

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

# T-Revision (2019)



The transition matrix can be revised and updated during training for its improved estimation.

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

# 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.

33

# 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. *TPAMI*, 2022.

# LCCN (2023)

Updating noise transition using backpropagation is unstable due to **mini-batch** computation.



Constrain the transition within the Dirichlet space

The learning is constrained to a simplex derived from the **entire dataset**, rather than the mini-batch, thus improving stability.

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

# ROBOT (2023)

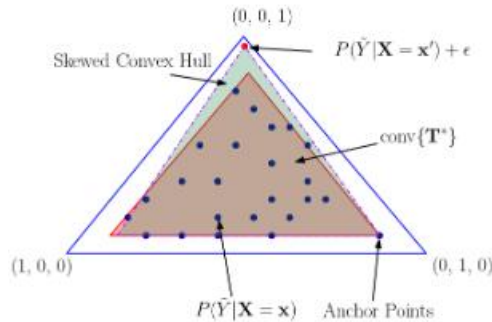A good transition matrix should simultaneously lead to the optimal forward correction loss and the noise-robust loss.

$$\min_T L_{rob}(f_{\hat{\theta}(T)}, \widetilde{D}_v) \text{ s.t. } \hat{\theta}(T) = \text{argmin } L(Tf_\theta, \widetilde{D}_{tr})$$



(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.

36

# AIDTM (2024)

Noise transition matrices are **annotator-** and **instance-**dependent.



Parameterize **instance-dependent** matrices with deep neural networks.

Assume that similar annotators share common noise pattern, thereby ease **annotator-dependency.**

S. Li et al. Transferring Annotator- and Instance-dependent Transition Matrix for Learning from Crowds. *TPAMI*, 2024.

# 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)

Noisy target          Softmax prediction

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

One-hot prediction

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

Interpolation

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

# Mixup (2018)

```
# 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.

ERM            *mixup*

(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)



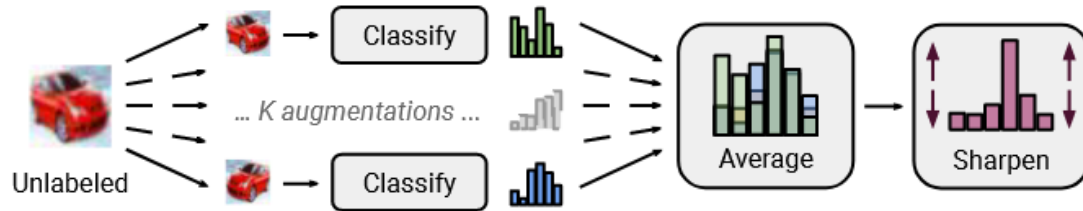**Augmentation preserves consistency**

**MixMatch**:

**Averaging** predictions across augmentations and **sharpening** as pseudo labelling.

**FixMatch**:

Aligning predictions of **strong** augmentation with pseudo-labels from **weak** augmentation.

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 accuracy plots (Intact, Flipped, Test) vs Epoch.

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

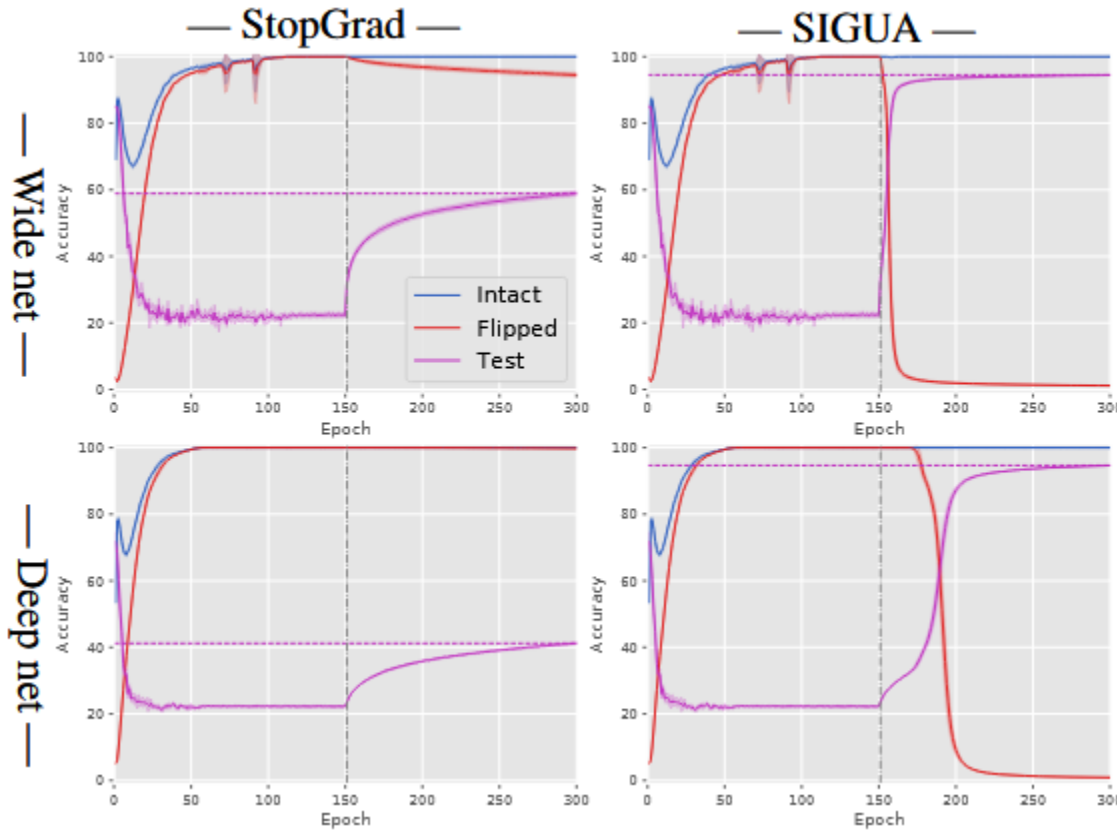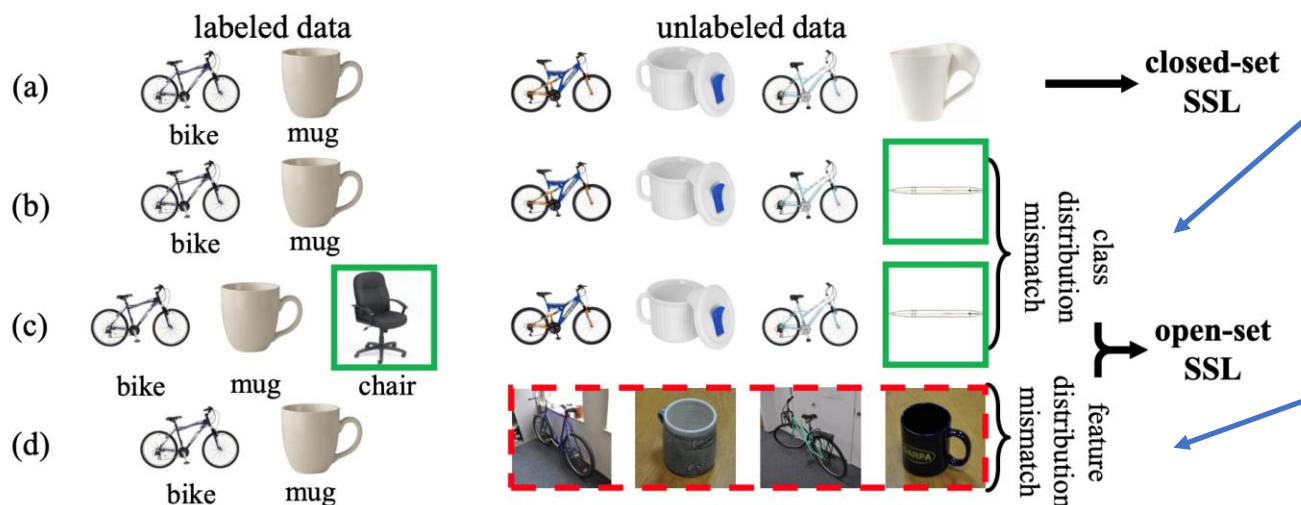**Require:** base learning algorithm $\mathfrak{B}$, optimizer $\mathfrak{O}$, mini-batch $\mathcal{S}_b = \{(x_i, \tilde{y}_i)\}_{i=1}^{n_b}$ of batch size $n_b$, current model $f_\theta$ where $\theta$ holds the parameters of $f$, good- and bad-data conditions $\mathfrak{C}_{good}$ and $\mathfrak{C}_{bad}$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

1: $\{\ell_i\}_{i=1}^{n_b} \leftarrow \mathfrak{B}.\text{forward}(f_\theta, \mathcal{S}_b)$    # forward pass
2: $\ell_b \leftarrow 0$    # initialize loss accumulator
3: **for** $i = 1, \dots, n_b$ **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)

**Open-set semi-supervised learning**: Labeled and unlabeled datasets may differ in both **class** and **feature** distribution.



**Class Distribution**: Unlabeled data **fall outside the label space**, which should be **detected and filtered**.

**Feature Distribution**: Unlabeled data **come from different domains**, which should perform **domain adaptation**.

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.

# Label Wave (2024)

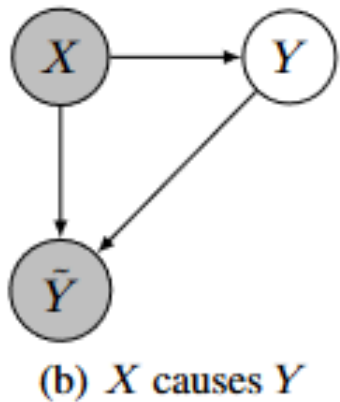**Tracking prediction changes** on the training set for **early stopping** (stop at Point 1) without validation data.

Behaviors of train and test are correlated

Consistent predictions    Fluctuate predictions    Consistent predictions



Learn true patterns    Fit mislabeled data    Overfit clean data

S. Yuan et al. Early Stopping Against Label Noise without Validation Data. In *ICLR*, 2024.

# 1-SAM (2024)

**Sharpness enhances robustness** (e.g., SAM [1]) but increases computational costs. It can be simplified by two penalty:

$$\ell(x_i, y_i; w) + \|z_i\|_2 + \|v\|_2$$

**Equivalent to SAM, which is proven to be robust.**

Penalty on embeddings    Penalty on last-layer weights

[1] P. Foret et al. Sharpness-Aware Minimization for Efficiently Improving Generalization. In *ICLR*, 2021.

C. Baek et al. Why is SAM Robust to Label Noise? In *ICLR*, 2024. https://bhanml.github.io & https://github.com/tmlr-group

# 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.

# Part VI: Future Directions

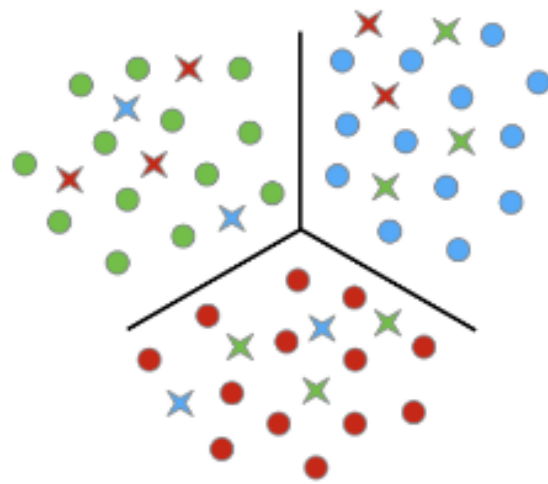## A Survey of Label-noise Representation Learning: Past, Present and Future

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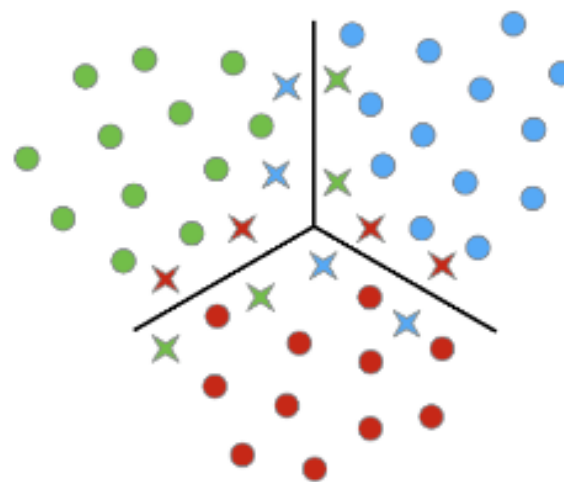
**Index Terms**—Machine Learning, Representation Learning, Weakly Supervised Learning, Label-noise Learning, Noisy Labels.

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.

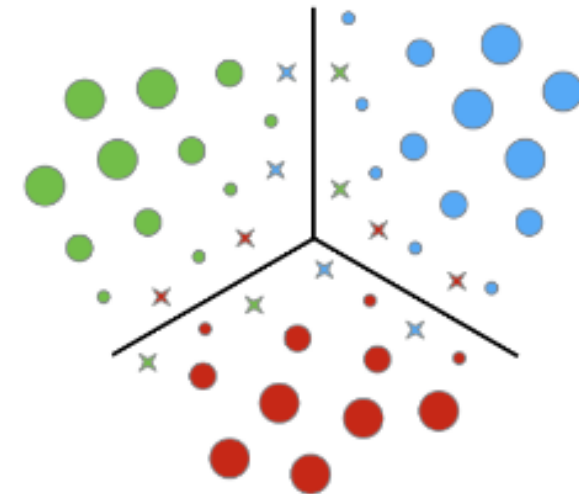https://bhanml.github.io & https://github.com/tmlr-group

# 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.
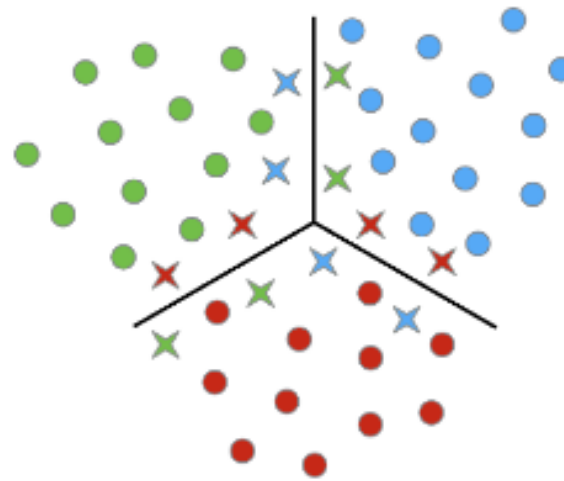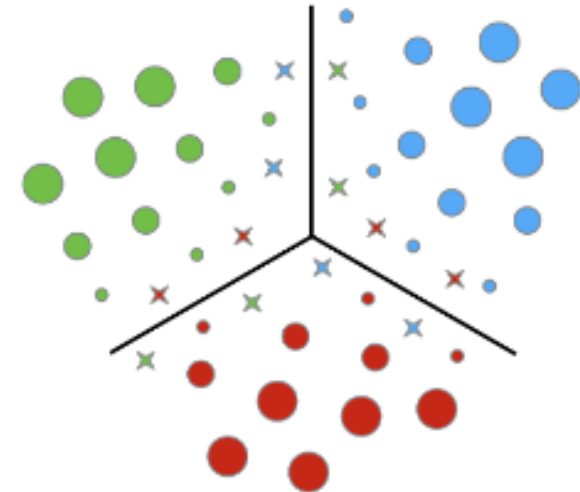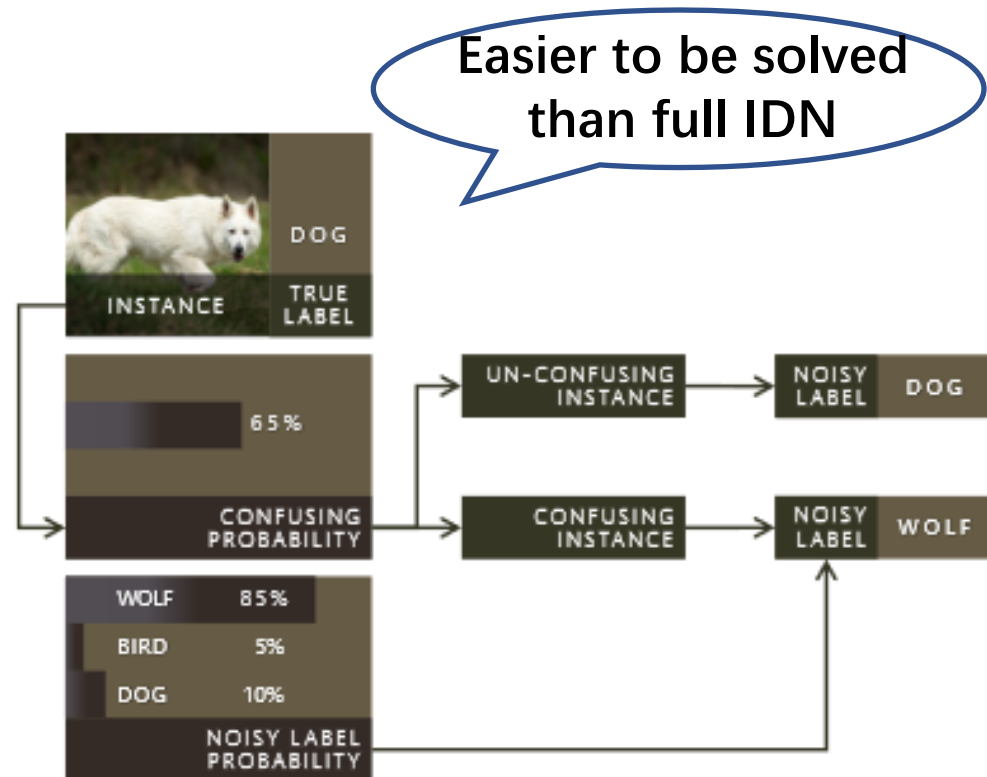
# CSIDN (2021)



(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 full IDN

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)

Graphical causal model which reveals a generative process of the data which contains instance-dependent label noise.

Latent variable    SVHN image    Clean label
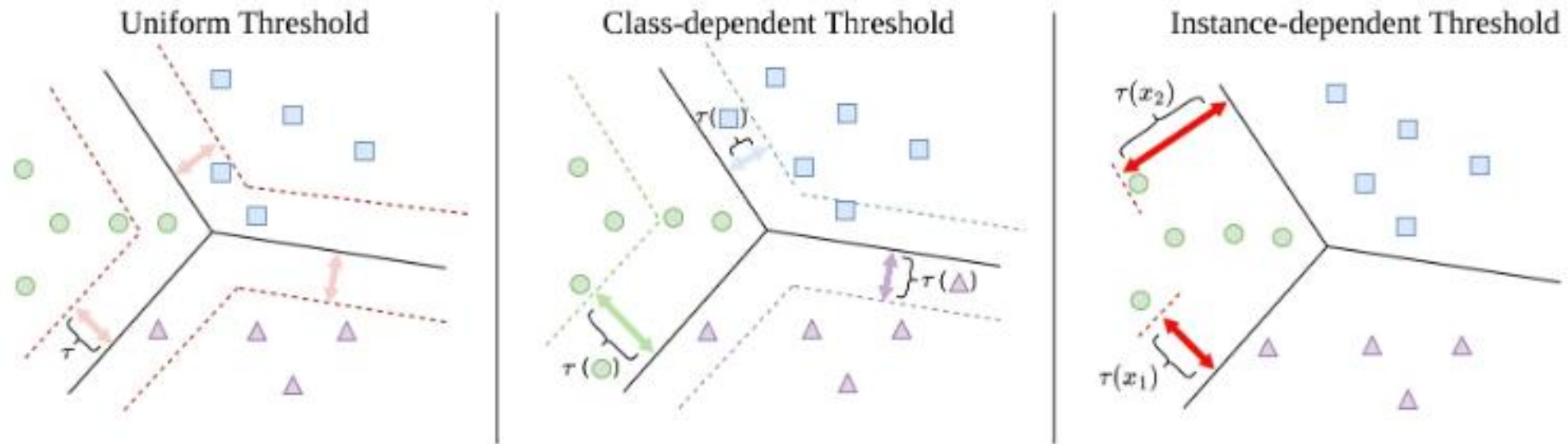


Orientation
lighting
font style

Noisy label

The joint distribution can be factorized as
$$P(X, \tilde{Y}, Y, Z) = P(Y)P(Z)P(X|Y,Z)P(\tilde{Y}|Y,X).$$

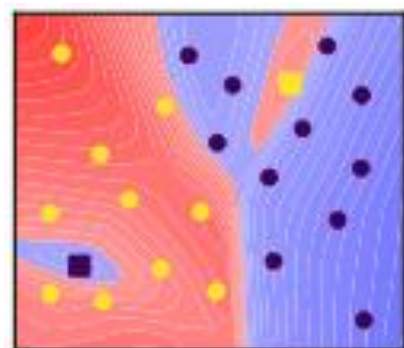Adding a constraint on $P(X|Y,Z)$ will reduce the uncertainty in $P(\tilde{Y}|Y,X)$.

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

# InstanT (2023)



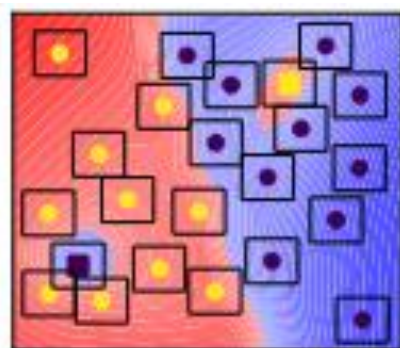Uniform Threshold | Class-dependent Threshold | Instance-dependent Threshold

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**



**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

Legend: Class labeled data, Unlabeled data, Similar, Noisy similar

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 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 Rationale

e.g., the irrelevant **base-10 information** is included in rationale

**Input: CoT prompting with clean rationales**

**Question-1:** In base-9, what is 86+57?
**Rationale-1:** In base-9, the digits are "012345678". We have 6 + 7 = 13 in base-10. Since we're in base-9, that exceeds the maximum value of 8 for a single digit. 13 mod 9 = 4, so the digit is 4 and the carry is 1. We have 8 + 5 + 1 = 14 in base 10. 14 mod 9 = 5, so the digit is 5 and the carry is 1. A leading digit 1. So the answer is 154.
**Answer-1:** 154.
⋯ **Q2, R2, A2, Q3, R3, A3** ⋯
**Question :** In base-9, what is 62+58?

**Input: CoT prompting with noisy rationales**

**Question-1:** In base-9, what is 86+57?
**Rationale-1:** In base-9, the digits are "012345678". We have 6 + 7 = 13 in base-10. **13 + 8 = 21**. Since we're in base-9, that exceeds the maximum value of 8 for a single digit.13 mod 9 = 4, so the digit is 4 and the carry is 1. We have 8 + 5 + 1 = 14 in base 10. 14 mod 9 = 5, so the digit is 5 and the carry is 1. **5 + 9 = 14.** A leading digit is 1. So the answer is 154.
**Answer-1:** 154.
⋯ **Q2, R2, A2, Q3, R3, A3** ⋯
**Question:** In base-9, what is 62+58?

While the test question asks about **base-9 calculation**

Z. Zhou et al. Can Language Models Perform Robust Reasoning in Chain-of-thought Prompting with Noisy Rationales? In *NeurIPS*, 2024.

# Noisy Model

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

# 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 Detection (NoisyGPT)



Use transition matrix to rectify noise

Use GPT to detect and rectify noise

$(\mathbf{x}, \tilde{y})$ → Classifier → $T$ → $y$

Transition matrix

$(\mathbf{x}, \tilde{y})$ → NoiseGPT → $y$

$\mathbf{x}$ perturbation
$\tilde{\mathbf{x}}_{p1}$
$\tilde{\mathbf{x}}_{p2}$
$\tilde{\mathbf{x}}_{p3}$

$q_\theta = \mathrm{MLLM}(\mathbf{X}|\tilde{\mathbf{X}})$

$\tilde{\mathbf{x}}_{p1}^{clean}$
$\tilde{\mathbf{x}}^{clean}$
$\tilde{\mathbf{x}}_{p2}^{clean}$
$\tilde{\mathbf{x}}_{p3}^{clean}$
$\tilde{\mathbf{x}}_{p1}^{noisy}$
$\tilde{\mathbf{x}}_{p3}^{noisy}$
$\tilde{\mathbf{x}}_{p2}^{noisy}$
$\tilde{\mathbf{x}}^{noisy}$

Randomness

$\psi_{\mathbf{0}}$
$\psi_{\mathbf{d}}$
$p_{\mathbf{X}|\tilde{\mathbf{x}}}$

GPT behaves different for noisy and clean examples, which can help us identify noise.

H. Wang et al. NoisyGPT: Label Noise Detection and Rectification through Probability Curvature. In *NeurIPS*, 2024.

# Noisy Adaptation



Performance ranking distribution of five methods.

C. Cao et al. Noisy Test-time Adaptation in Vision-Language Models. In *ICLR*, 2024.

# Noisy Correction



**Design3**

Add Gaussian noise $\epsilon_t$

**+**

Add noise and then denoise to suppress extreme noise

**Design1**    **Design1**

Diffusion Model

Encode to latent space

clip    clip

**+**

Diffusion Model

Denoise

Interpolation results

Natural images

**Design2**

Perform interpolation in the noisy space rather than latent space

P. Zheng et al. NoiseDiffusion: Correcting Noise for Image Interpolation with Diffusion Models beyond Spherical Linear Interpolation. In *ICLR*, 2024.

# Noisy Dataset

*Photos of **ice bear** in **snow** background*
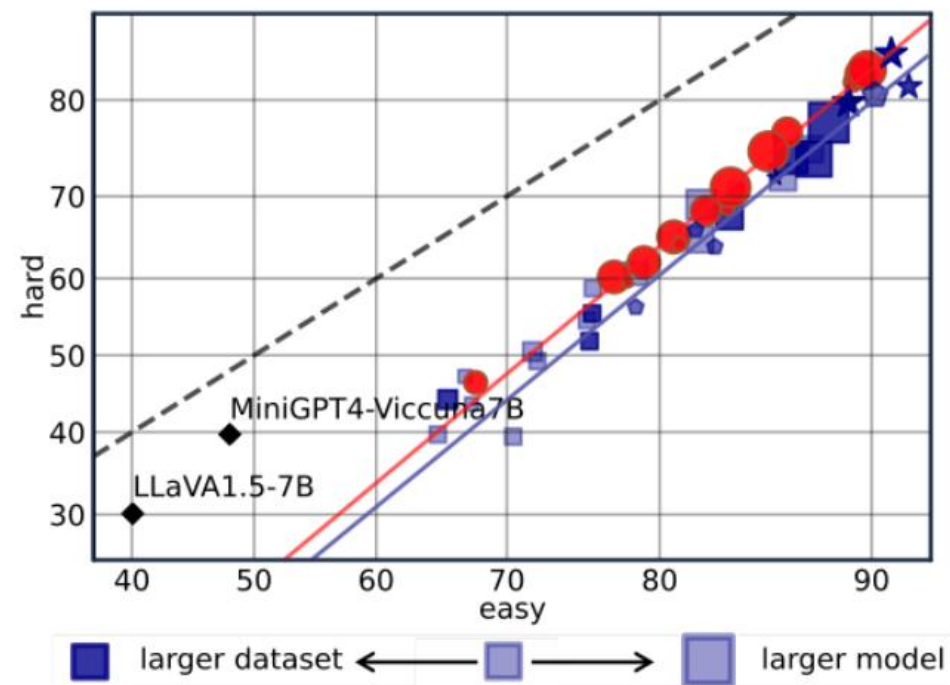


*Photos of **ice bear** in **grass** background*



Background changes lead to potential spurious features.

Spurious features still affect CLIP robustness.

Q. Wang et al. A Sober Look at the Robustness of CLIPs to Spurious Features. In *NeurIPS*, 2024.

# Datasets and Benchmark

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. https://bhanml.github.io & https://github.com/tmlr-group
A Survey of Label-noise Representation Learning: Past, Present and Future. *arXiv preprint:2011.04406*, 2020.
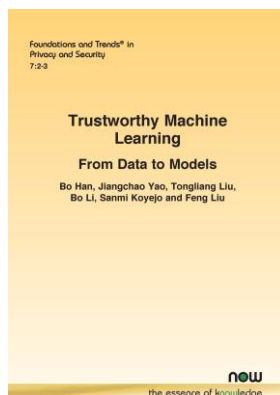
# Appendix

- 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**, 2025.
  - Trustworthy Machine Learning under Imperfect Data. CS series, **Springer Nature**, 2025.
  - Trustworthy Machine Learning: From Data to Models. **Foundations and Trends® in Privacy and Security**, 2025.

- 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
  - IJCAI 2024 Tutorial on Trustworthy Machine Learning under Imperfect Data
  - WWW 2025 Tutorial on Trustworthy AI under Imperfect Web 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 AIP 2024 Joint Workshop on Artificial Intelligence and Machine Learning

Foundations and Trends® in Privacy and Security
7:2-3

**Trustworthy Machine Learning**

**From Data to Models**

Bo Han, Jiangchao Yao, Tongliang Liu, Bo Li, Sanmi Koyejo and Feng Liu

now
the essence of knowledge