

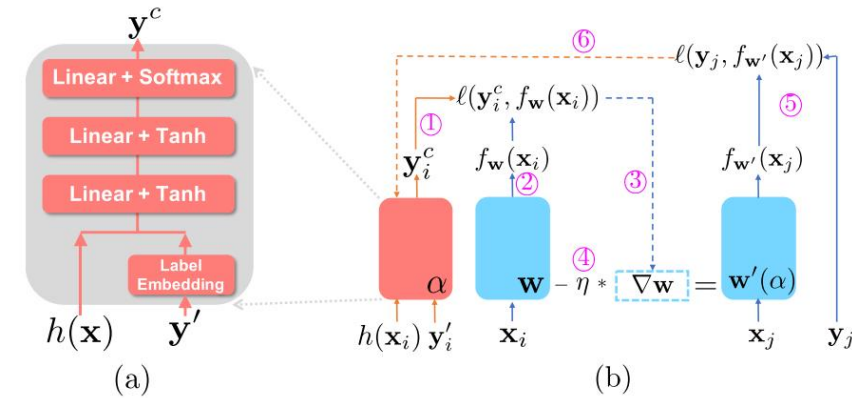
# Learning from Noisy Labels with Decoupled Meta Label Purifier

Yuanpeng Tu<sup>1</sup>, Boshen Zhang<sup>2</sup>, Yuxi Li<sup>2</sup>, Liang Liu<sup>2</sup>, Jian Li<sup>2</sup>,  
Yabiao Wang<sup>2</sup>, Chengjie Wang<sup>2,3†</sup>, Cai Rong Zhao<sup>1†</sup>

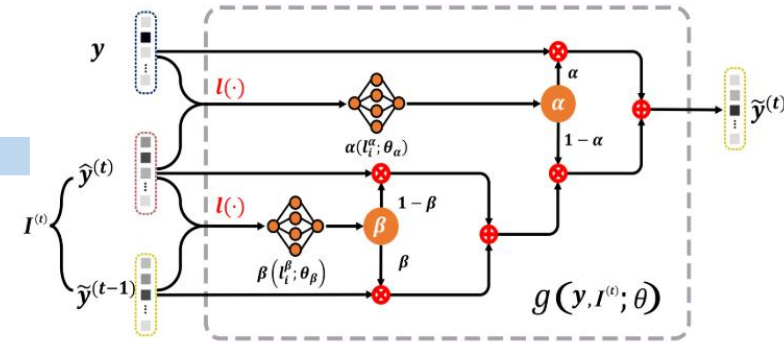
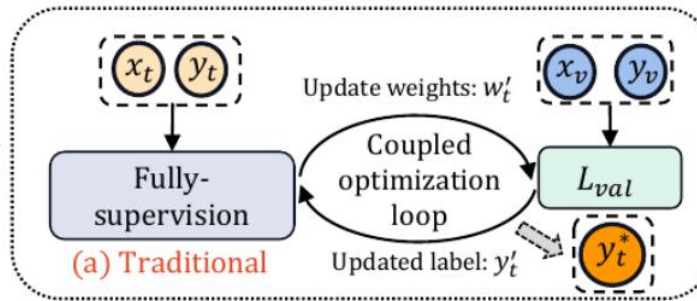
*1 Tongji University, 2 Tencent Youtu Lab, 3 Shanghai Jiao Tong University*

*Corresponding authors. Email: zhaocairong@tongji.edu.cn, jasoncjwang@tencent.com*

# Motivation

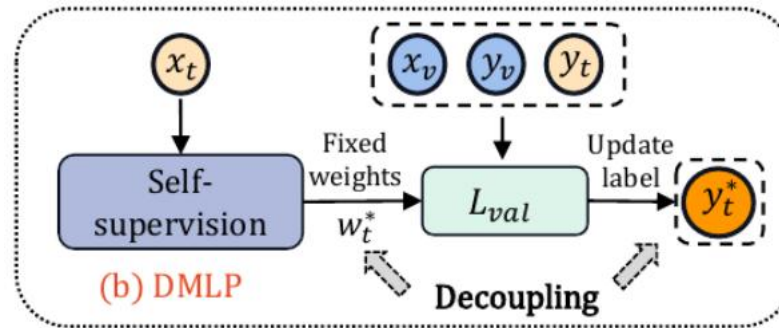


## Coupled Optimization Methods



VS

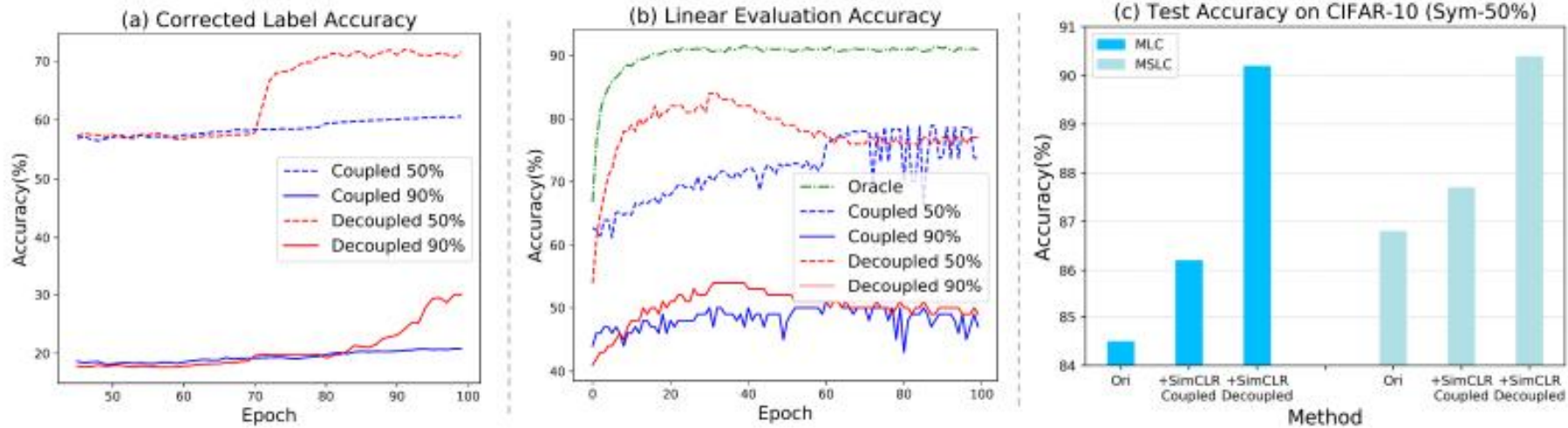
Which one is more suitable for LNL?



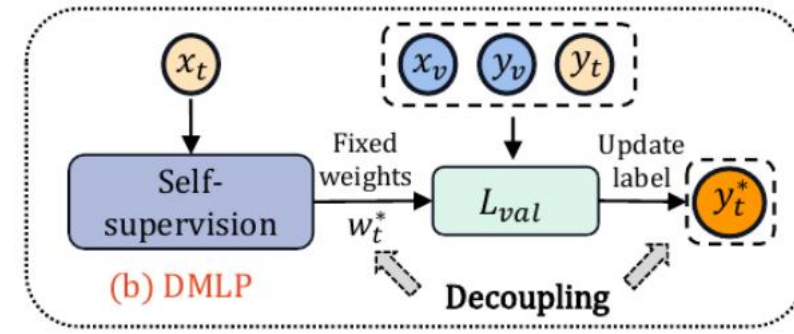
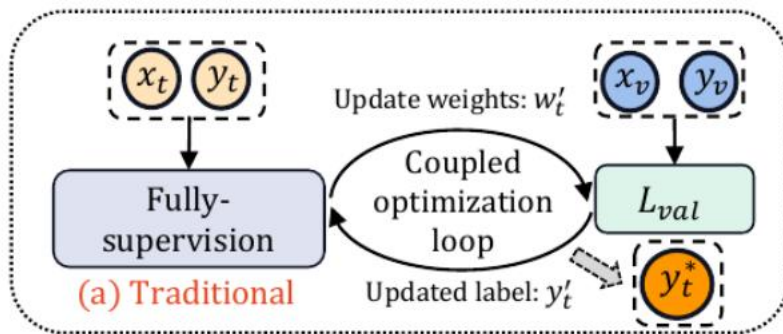
Decoupled Optimization Methods

[1] Learning to Purify Noisy Labels via Meta Soft Label Corrector. AAAI21.  
 [2] Meta Label Correction for Noisy Label Learning. AAAI21.

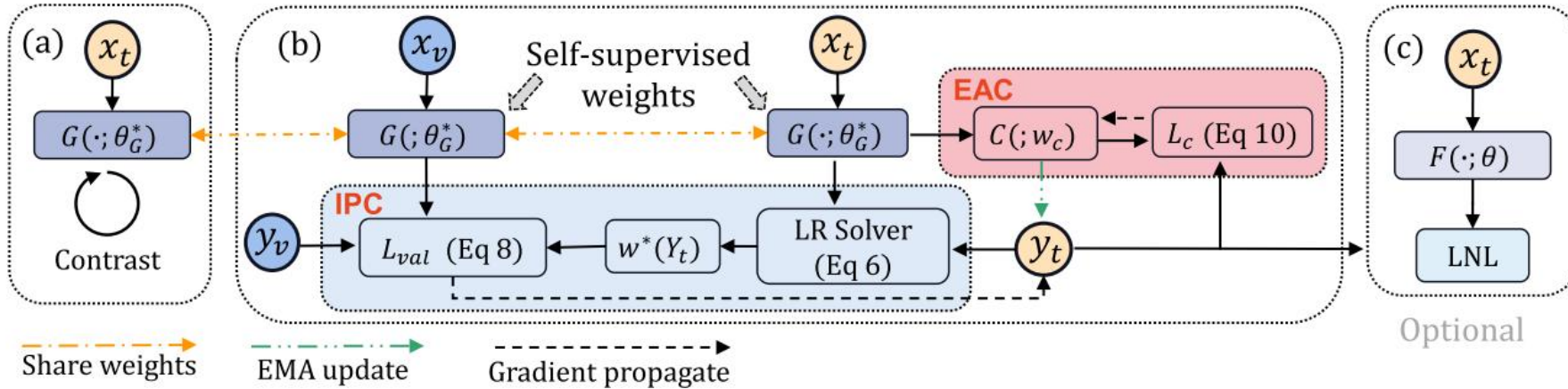
# Motivation



Conclusion: Previous coupled meta label purifier is sub-optimal in terms of both model weights and labels.



# Framework



## Intrinsic Primary Correction

$$\min_{\theta_\alpha} E_{(x_v, y_v) \in D_v} \mathcal{L}_{val}(x_v, y_v; w^*(\theta_\alpha))$$

s.t.  $w^*(\theta_\alpha) = \arg \min_w E_{(x_t, y_t) \in D_t} \mathcal{L}_{train}(x_t, y_t; w, \theta_\alpha)$

Training Data Matrix:

$$F_t = [\mathbf{f}_{t,1}, \mathbf{f}_{t,2}, \dots, \mathbf{f}_{t,b}]^T, \quad Y_t = [y_{t,1}, y_{t,2}, \dots, y_{t,b}]^T.$$

$$\min_w \mathcal{L}_{train}(F_t, Y_t; w) = \|Y_t - F_t w\|^2.$$

Least Square Method

$$w^*(Y_t) = (F_t^T F_t)^{-1} F_t^T Y_t.$$

Optimal Solution of Train set

$$Y_t^{p+1} := Y_t^p - \eta_I \nabla (\mathcal{L}_{val}(Y_t^p)).$$

Update Noisy Labels

Minimize Discrepancy

$$y'_{v,i}(Y_t) = w^*(Y_t)^T \mathbf{f}_{v,i}.$$

Classification on Val set

## Extrinsic Auxiliary Correction

$$\mathcal{L}_c(w_c) = \mathcal{L}_{ce}(C(\mathbf{f}_t; w_c), y'_t) + \mathcal{H}(C(\mathbf{f}_t; w_c)),$$

Update Noisy Labels

$$Y_t^{p+1} := (1 - \eta_E) Y_t^p + \eta_E C(F_t; w_c)$$

**Algorithm 1** The workflow of DMLP.

**Input:** Noisy training set  $D_t$ , clean validation set  $D_v$ , feature extractor  $G(\cdot; \theta_G)$ , classifier  $C(\cdot; w_c)$ , batch size  $b$ , max iterations  $m$ , period for regular label substitution  $T$ .

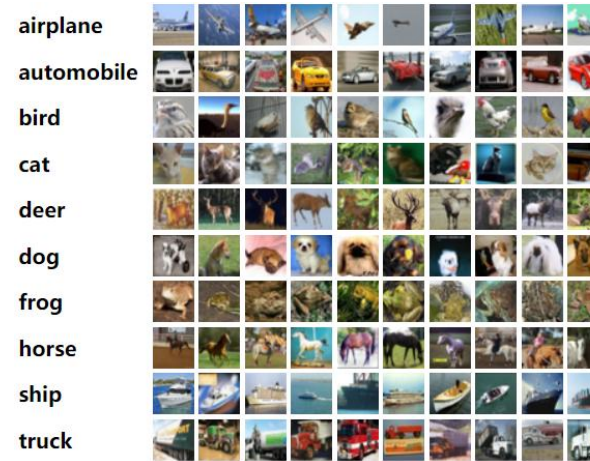
**Procedure:**

- 1: Self-supervised training for  $G(\cdot; \theta_G)$
- 2: Generate features  $\mathbf{f}$  by Eq. (1)
- 3: **for**  $i = 1$  to  $m$  **do**
- 4:    /\*IPC starts\*/
- 5:     $\{F_t, Y_t\} \leftarrow \text{SampleMiniBatch}(\mathbf{f}, D_t, b)$
- 6:    Calculate closed-form solution  $w^*(Y_t)$  by Eq. (6)
- 7:    Predict validation set labels  $y'_v$  by Eq. (7)
- 8:    Calculate label purification loss  $\mathcal{L}_{\text{val}}(Y_t)$  by Eq. (8).
- 9:    Update training labels  $Y_t$  in backward process.
- 10:
- 11:    /\*EAC starts\*/
- 12:    Calculate loss for the classifier  $C(\cdot; w_c)$  by Eq. (10)
- 13:    Update classifier parameter  $w_c$  in backward process.
- 14:    **if**  $i = nT$  **then**
- 15:      Update training labels  $Y_t$  by Eq. (11)
- 16:    **end if**
- 17: **end for**

**Output:** The purified labels  $Y_t^*$ .

## Dataset

### Simulated Noisy Dataset: CIFAR-10/100



### Real-world Noisy Dataset: Clothing1M



## Symmetric Noise on CIFAR-10/100

Table 1. Comparison with state-of-the-art methods on CIFAR-10/100 datasets with symmetric noise. “CE” is the standard ConvNet trained with Cross-Entropy loss in an end-to-end manner. “Classifier” means adopts the pre-trained SimCLR features to re-train a linear classifier. “Val” denotes using a small clean validation set. DivideMix\* denotes training DivideMix with the same validation set as additional data.

Dataset Method	Val	Noise ratio	CIFAR-10				CIFAR-100			
			20%	50%	80%	90%	20%	50%	80%	90%
Cross-Entropy (CE)	✗	Best	86.8	79.4	62.9	42.7	62.0	46.7	19.9	10.1
		Last	82.7	57.9	26.1	16.8	61.8	37.3	8.8	3.5
Co-teaching+ [36]	✗	Best	89.5	85.7	67.4	47.9	65.6	51.8	27.9	13.7
		Last	88.2	84.1	45.5	30.1	64.1	45.3	15.5	8.8
PENCIL [35]	✗	Best	92.4	89.1	77.5	58.9	69.4	57.5	31.1	15.3
		Last	92.0	88.7	76.5	58.2	68.1	56.4	20.7	8.8
REED [38]	✗	Best	95.8	95.6	94.3	93.6	76.7	73.0	66.9	59.6
		Last	95.7	95.4	94.1	93.5	76.5	72.2	66.5	59.4
Sel-CL+ [18]	✗	Best	95.5	93.9	89.2	81.9	76.5	72.4	59.6	48.8
		Last	95.1	93.3	88.7	81.6	76.1	72.0	59.2	48.6
MOIT+ [21]	✗	Best	94.1	91.8	81.1	74.7	75.9	70.6	47.6	41.8
		Last	93.8	91.3	80.6	74.0	75.2	70.1	46.9	41.2
C2D-DivideMix [40]	✗	Best	<b>96.3</b>	95.2	94.4	93.5	78.6	76.4	67.7	58.7
		Last	<b>96.2</b>	95.1	94.1	93.4	78.3	76.0	67.4	58.4
DivideMix [15]	✗	Best	96.1	94.6	93.2	76.0	77.3	74.6	60.2	31.5
		Last	95.7	94.4	92.9	75.4	76.9	74.2	59.6	31.0
Meta-Learning [16]	✓	Best	92.9	89.3	77.4	58.7	68.5	59.2	42.4	19.5
		Last	92.0	88.8	76.1	58.3	67.7	58.0	40.1	14.3
MLC [41]	✓	Best	92.6	88.1	77.4	67.9	66.8	52.7	21.8	15.0
		Last	91.8	87.5	77.1	67.0	66.5	52.4	18.9	14.2
MSLC [9]	✓	Best	93.4	89.9	69.8	56.1	72.5	65.4	24.3	16.7
		Last	93.3	89.4	68.8	55.2	72.0	64.9	20.5	14.6
DivideMix* [15]	✓	Best	96.1	94.9	93.6	77.3	77.7	74.8	60.7	32.5
		Last	95.9	94.6	93.0	76.5	77.1	74.3	60.5	32.2
DMLP-Naive	✓	Best	94.7	94.2	93.5	92.8	72.7	68.0	63.5	61.3
		Last	94.2	94.0	93.2	92.0	72.3	67.4	63.2	60.9
DMLP-DivideMix	✓	Best	<b>96.3</b>	<b>95.8</b>	<b>94.5</b>	<b>94.3</b>	<b>79.9</b>	<b>76.8</b>	<b>68.6</b>	<b>65.8</b>
		Last	<b>96.2</b>	<b>95.6</b>	<b>94.3</b>	<b>94.0</b>	<b>79.4</b>	<b>76.1</b>	<b>68.5</b>	<b>65.4</b>

## Asymmetric Noise on CIFAR-10

Table 2. Evaluation results with asymmetric noise of different noisy ratio on CIFAR-10. “Validation” denotes the method exploits a small clean validation set.

Method	Validation	Noisy ratio	
		20%	40%
Joint-Optim [28]	✗	92.8	91.7
PENCIL [35]	✗	92.4	91.2
M-correction [1]	✗	-	86.3
Iterative-CV [4]	✗	-	88.0
DivideMix [15]	✗	93.4	93.4
REED [38]	✗	95.0	92.3
C2D-DivideMix [40]	✗	93.8	93.4
Sel-CL+ [18]	✗	<b>95.2</b>	93.4
GCE [11]	✗	87.3	78.1
RRL [17]	✗	-	92.4
Zhang, et al. [39]	✓	92.7	90.2
Meta-Learning [16]	✓	-	88.6
MSLC [9]	✓	94.4	91.6
DMLP-Naive	✓	94.6	93.9
DMLP-DivideMix	✓	<b>95.2</b>	<b>95.0</b>

## Real-world Clothing1M

Table 3. Top-1 testing accuracy on Clothing-1M testset. “Validation” denotes using the validation provided by [31].

Method	Validation	Top-1 Accuracy
PENCIL [35]	✗	73.49
DivideMix [15]	✗	74.76
RRL [17]	✗	74.90
GCE [11]	✗	73.30
C2D-DivideMix [40]	✗	74.30
REED [38]	✗	75.81
Meta-Learning [16]	✓	73.47
Self-Learning [13]	✓	76.44
MLC [41]	✓	75.78
MSLC [9]	✓	74.02
Meta-Cleaner [10]	✓	72.50
Meta-Weight [8]	✓	73.72
FaMUS [34]	✓	74.40
MSLG [7]	✓	76.02
DMLP-Naive	✓	77.77
DMLP-DivideMix	✓	<b>78.23</b>

## Generality of DMLP

Table 4. Comparison between the LNL methods and their DMLP applications with symmetric noise on CIFAR-10/100. Specifically, the 9-layer CNN is adopted as the backbone network of Co-teaching.

Dataset Method/Noise ratio		CIFAR-10				CIFAR-100			
		20%	50%	80%	90%	20%	50%	80%	90%
Co-teaching [12]	Best	82.6	73.0	24.0	14.6	50.5	38.2	11.8	4.9
	Last	81.9	72.6	23.5	11.7	50.3	38.0	11.3	4.3
DMLP-Co-teaching	Best	<b>85.8</b>	<b>85.8</b>	<b>85.4</b>	<b>84.6</b>	<b>51.2</b>	<b>49.8</b>	<b>48.1</b>	<b>45.3</b>
	Last	<b>85.6</b>	<b>85.6</b>	<b>85.3</b>	<b>84.5</b>	<b>51.0</b>	<b>49.3</b>	<b>47.8</b>	<b>45.1</b>
CDR [30]	Best	90.4	85.0	47.2	12.3	63.3	39.5	29.2	8.0
	Last	82.7	49.4	16.6	10.1	62.9	39.5	9.7	4.5
DMLP-CDR	Best	<b>91.4</b>	<b>91.2</b>	<b>91.2</b>	<b>90.2</b>	<b>69.2</b>	<b>64.8</b>	<b>61.4</b>	<b>58.5</b>
	Last	<b>91.2</b>	<b>90.8</b>	<b>90.6</b>	<b>89.3</b>	<b>68.3</b>	<b>64.3</b>	<b>61.1</b>	<b>57.9</b>
ELR+ [19]	Best	94.6	93.8	91.1	75.2	77.5	72.4	58.2	30.8
	Last	94.4	93.7	90.5	73.5	76.2	72.2	56.8	30.6
DMLP-ELR+	Best	<b>94.9</b>	<b>94.1</b>	<b>93.0</b>	<b>92.5</b>	<b>77.8</b>	<b>73.6</b>	<b>63.9</b>	<b>60.5</b>
	Last	<b>94.6</b>	<b>94.0</b>	<b>92.7</b>	<b>92.1</b>	<b>77.1</b>	<b>73.4</b>	<b>63.6</b>	<b>60.5</b>

## Component Analysis

Table 5. Ablation study for the effectiveness of IPC and EAC in DMLP-Naive on CIFAR-10.

Component			CIFAR-10				Clothing 1M
IPC	EAC		20%	50%	80%	90%	
✗	✓	Best	93.7	93.3	91.1	67.4	76.5
		Last	93.0	92.9	90.6	66.5	76.1
✓	✗	Best	87.8	85.7	79.9	76.0	76.8
		Last	87.2	85.5	79.4	75.4	76.5
✓	✓	Best	<b>94.7</b>	<b>94.2</b>	<b>93.5</b>	<b>92.8</b>	<b>77.7</b>
		Last	<b>94.2</b>	<b>94.0</b>	<b>93.2</b>	<b>92.0</b>	<b>77.6</b>

## Comparison against other coupled purifiers with pre-training.

Table 6. Comparison with coupled meta label correction methods MLC [41] and MSLC [9] on CIFAR-10. "\*" denotes training with SimCLR pretrained ResNet-18.

Method		Noisy ratio			
		20%	50%	80%	90%
MLC*	Best	91.8	86.2	77.6	72.9
	Last	91.6	85.9	77.5	72.6
MSLC*	Best	92.0	87.7	78.0	67.8
	Last	92.0	87.5	77.9	67.3
DMLP-Naive*	Best	<b>94.0</b>	<b>93.7</b>	<b>93.1</b>	<b>92.3</b>
	Last	<b>93.9</b>	<b>93.4</b>	<b>92.9</b>	<b>91.9</b>
MLC*-DivideMix	Best	95.3	94.0	93.0	86.6
	Last	95.0	93.6	92.7	86.5
MSLC*-DivideMix	Best	95.7	94.9	93.8	83.0
	Last	95.5	94.8	93.1	82.8
DMLP*-DivideMix	Best	<b>96.3</b>	<b>95.6</b>	<b>94.1</b>	<b>93.8</b>
	Last	<b>96.0</b>	<b>95.2</b>	<b>94.0</b>	<b>93.6</b>

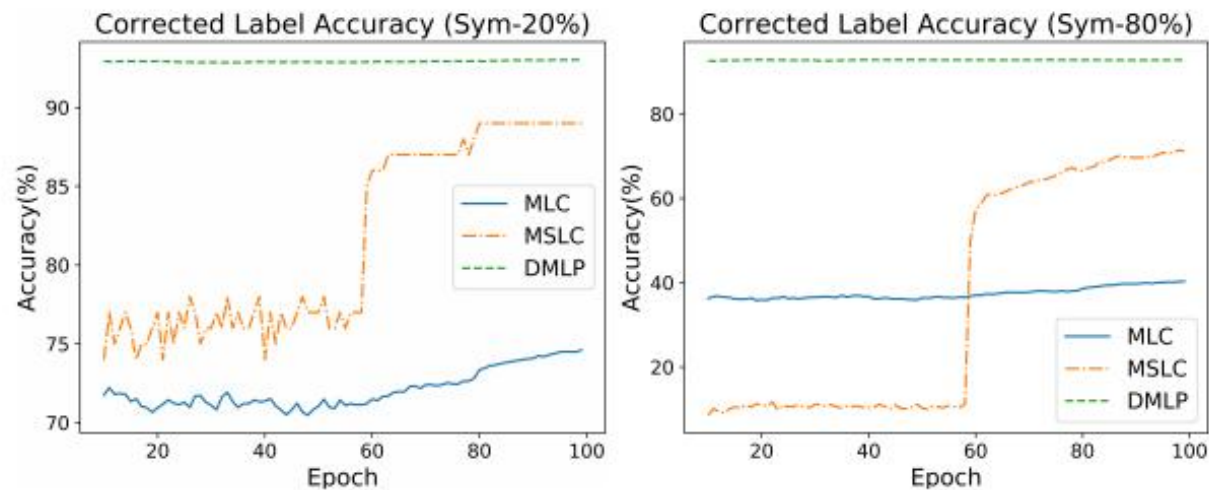


Figure 4. Comparison of corrected label accuracy under symmetric-20% (left), symmetric-80% (right) noise settings on CIFAR-10.



## Effect of validation size

Table 7. Investigation of the validation set size  $\tau$  on Clothing1M.

$\tau$	10%	20%	30%	40%	50%	100%
Accuracy (%)	75.50	76.40	76.61	77.00	77.30	<b>77.31</b>

## Effect of different feature representation for purification

Table 8. Ablation study for adopting different features in DMLP-Naive on CIFAR-10, where "R18/50" denote "ResNet-18/50" and "M/S" represent "MoCo/SimCLR".

Feature Source		Noisy ratio			
		20%	50%	80%	90%
R18 (M)	Best	93.8	93.3	92.2	90.4
	Last	93.7	92.7	92.1	90.0
R18 (S)	Best	94.0	93.7	93.1	92.3
	Last	93.9	93.4	92.9	91.9
R50 (S)	Best	<b>94.7</b>	<b>94.2</b>	<b>93.5</b>	<b>92.8</b>
	Last	<b>94.2</b>	<b>94.0</b>	<b>93.2</b>	<b>92.0</b>

## Performance under extremely noisy setting

Table 9. Comparison between recent semi-supervised methods and DMLP-DivideMix on CIFAR-10/100 with 100% noisy ratio.

Method	CIFAR-10	CIFAR-100
MeanTeacher	83.0	31.0
MixMatch	87.9	57.7
FixMatch	88.1	56.3
UDA	88.2	56.1
Ours	<b>91.7</b>	<b>60.1</b>

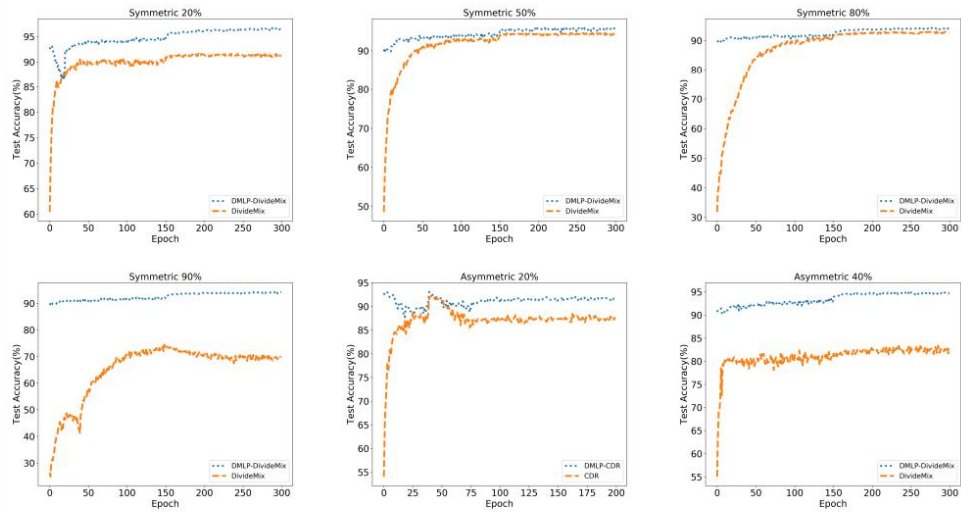


Figure 4. Accuracy curve of DMLP-DivideMix and DivideMix on CIFAR-10 under different noise settings.

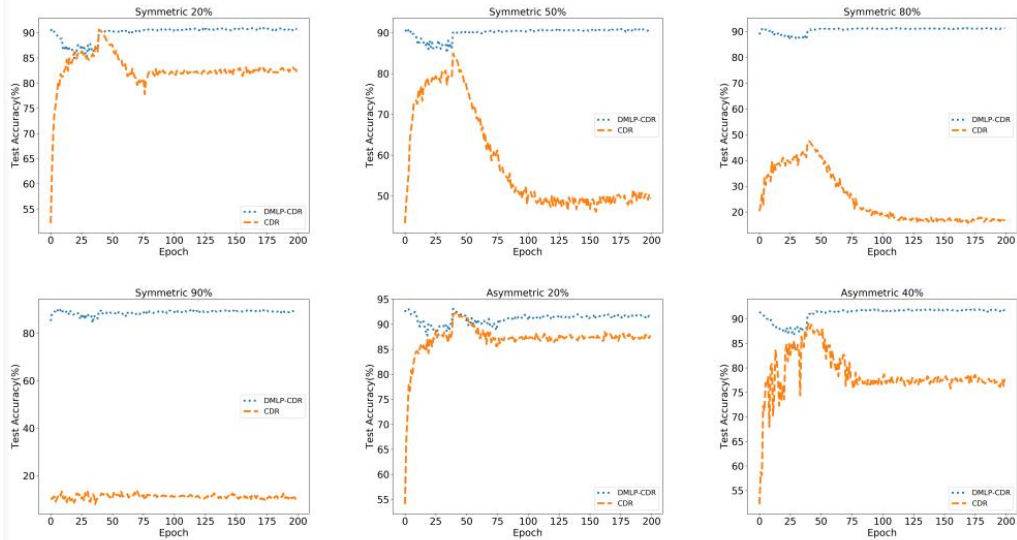


Figure 6. Accuracy curve of DMLP-CDR and CDR on CIFAR-10 under different noise settings.

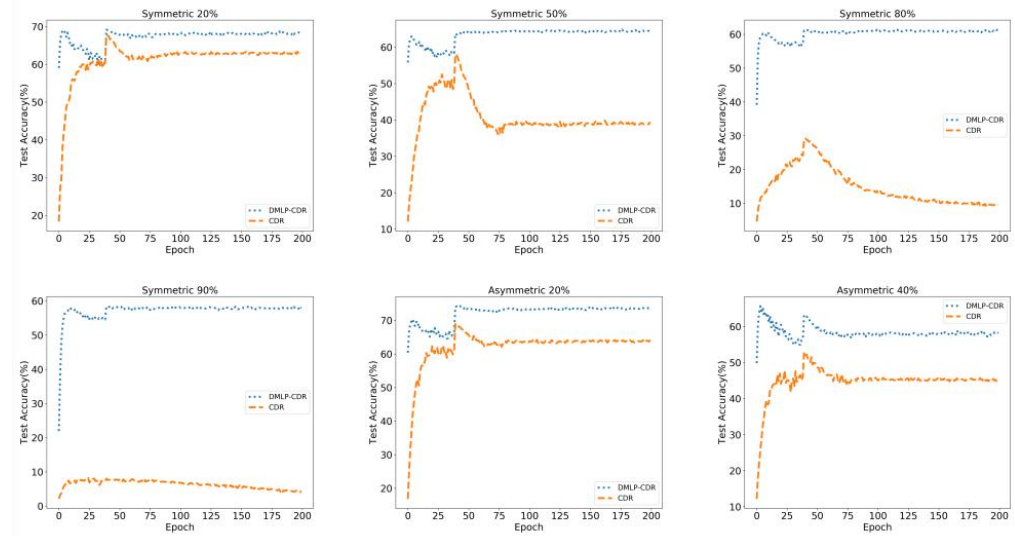


Figure 5. Accuracy curve of DMLP-DivideMix and DivideMix on CIFAR-100 under different noise settings.

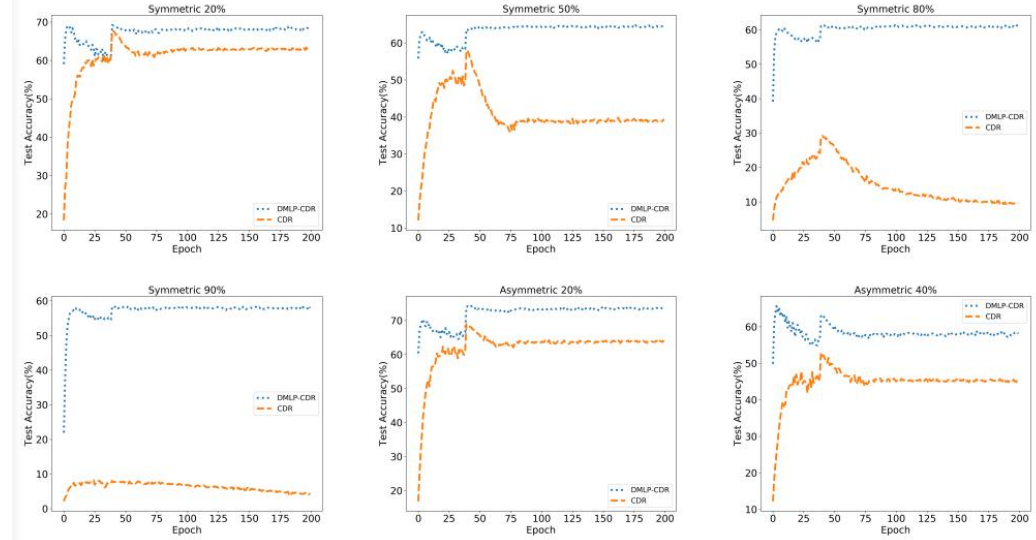


Figure 7. Accuracy curve of DMLP-CDR and CDR on CIFAR-100 under different noise settings.

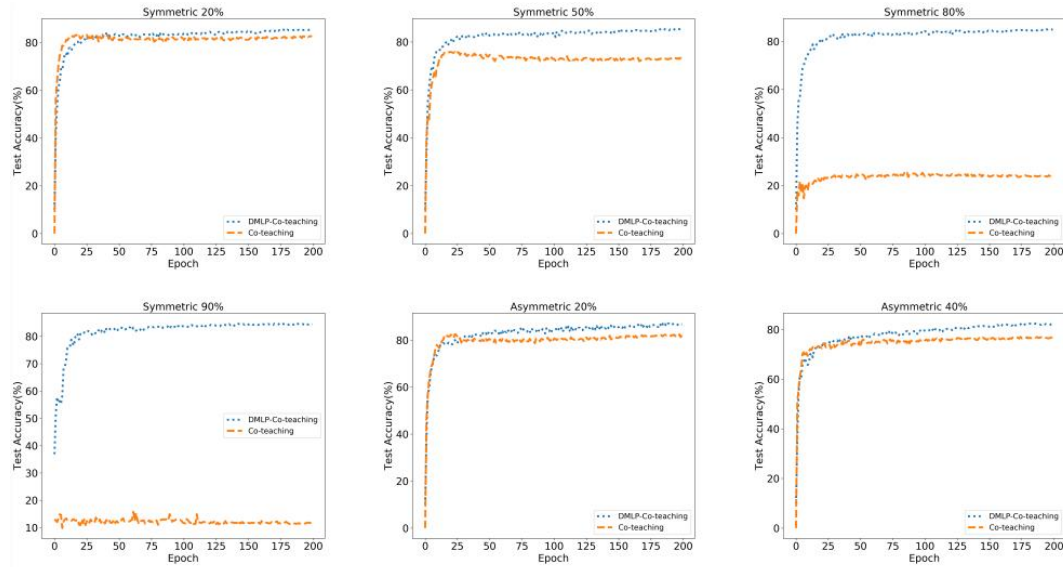


Figure 8. Accuracy curve of DMLP-Co-teaching and Co-teaching on CIFAR-10 under different noise settings.

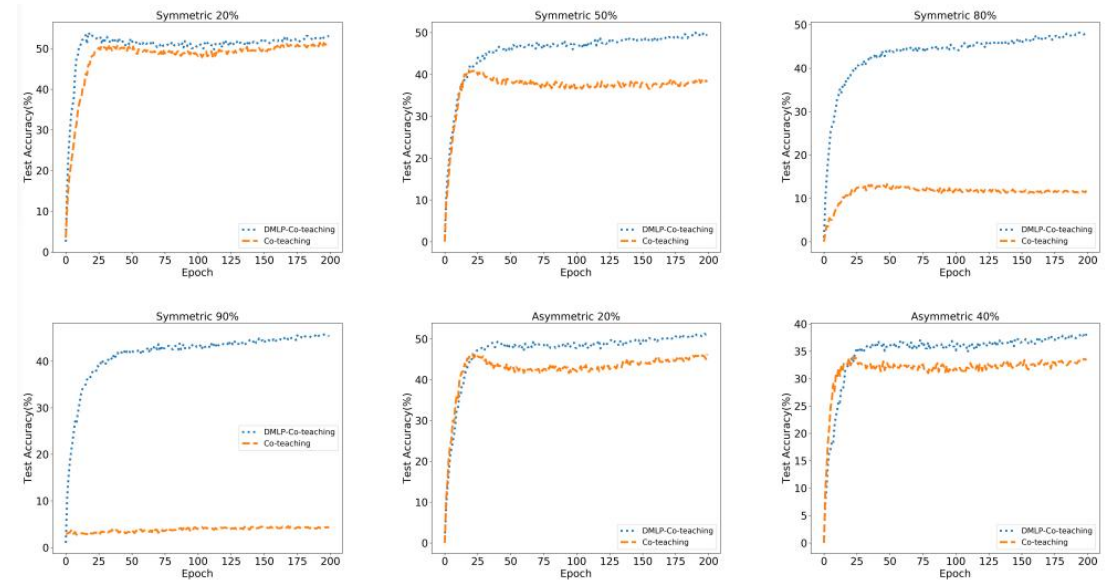


Figure 9. Accuracy curve of DMLP-Co-teaching and Co-teaching on CIFAR-100 under different noise settings.

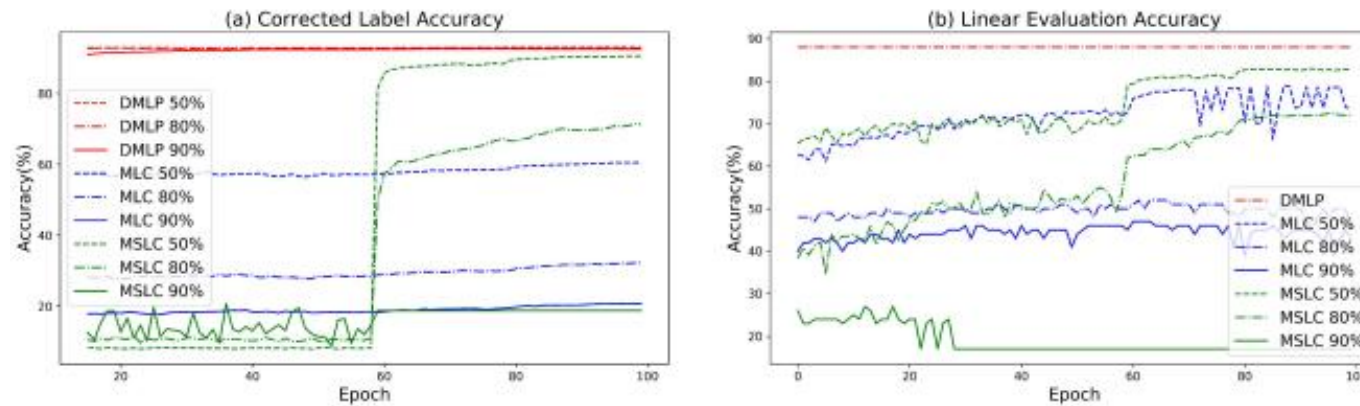


Figure 1. Comparison with two state-of-the-art coupled optimization based meta label correction methods MLC [25] and MSLC [6] on corrected label accuracy and linear evaluation accuracy.

Table 2. Comparison with MLC and MSLC on CIFAR-10/100. "†" denotes training with fixed self-supervised pretrained ResNet-18. "\*" denotes training with self-supervised pretrained ResNet-18.

Dataset Method	Noise ratio	CIFAR-10				CIFAR-100			
		20%	50%	80%	90%	20%	50%	80%	90%
MLC* [25]	Best	91.8	86.2	77.6	72.9	62.2	53.8	46.5	39.6
	Last	91.6	85.9	77.5	72.6	61.6	53.0	46.2	39.2
MSLC* [6]	Best	92.0	87.7	78.0	67.8	70.8	64.1	36.4	19.8
	Last	92.0	87.5	77.9	67.3	70.2	63.8	34.3	18.7
MLC† [25]	Best	92.0	90.2	89.0	88.9	65.9	59.4	54.4	54.2
	Last	91.6	89.4	88.5	88.1	65.2	59.2	54.1	54.0
MSLC† [6]	Best	92.1	90.4	87.3	84.7	71.7	64.7	53.3	46.8
	Last	92.0	90.0	87.2	84.2	71.6	64.4	53.0	46.4
DMLP-Naive	Best	<b>94.7</b>	<b>94.2</b>	<b>93.5</b>	<b>92.8</b>	<b>72.7</b>	<b>68.0</b>	<b>63.5</b>	<b>61.3</b>
	Last	<b>94.2</b>	<b>94.0</b>	<b>93.2</b>	<b>92.0</b>	<b>72.3</b>	<b>67.4</b>	<b>63.2</b>	<b>60.9</b>

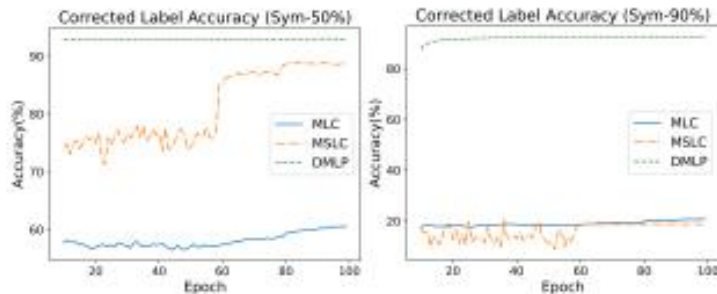


Figure 2. Comparison of corrected label accuracy curve under symmetric-50% (left), symmetric-90% (middle) noise settings on CIFAR-10.

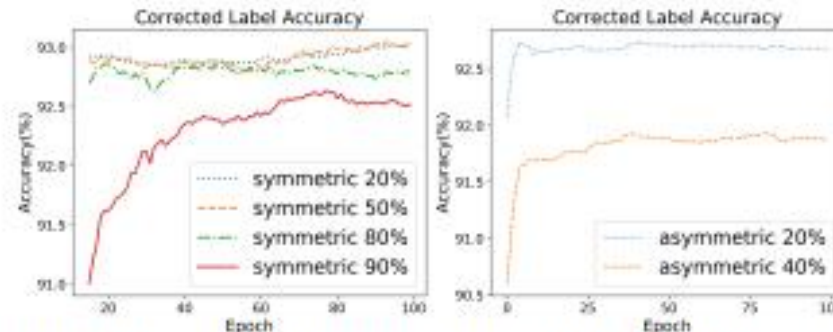


Figure 3. Corrected label accuracy curve of DMLP under symmetric (left), asymmetric (middle) noise settings on CIFAR-10.

Table 4. Comparison on CIFAR-10/100 with symmetric noise.

Dataset Method	CIFAR-10			
	20%	50%	80%	90%
DMLP-Naive	94.28±0.10	94.02±0.21	93.31±0.19	92.16±0.20
DMLP-DivideMix	<b>96.20±0.11</b>	<b>95.63±0.13</b>	<b>94.22±0.14</b>	<b>93.97±0.22</b>

Dataset Method	CIFAR-100			
	20%	50%	80%	90%
DMLP-Naive	72.39±0.08	67.60±0.26	63.17±0.14	61.09±0.20
DMLP-DivideMix	<b>79.31±0.21</b>	<b>76.11±0.10</b>	<b>68.42±0.12</b>	<b>65.55±0.23</b>

# Learning from Noisy Labels with Decoupled Meta Label Purifier

**Thanks for Watching!**

*<https://github.com/yuanpengtu/DMLP>*

Yuanpeng Tu<sup>1</sup>, Boshen Zhang<sup>2</sup>, Yuxi Li<sup>2</sup>, Liang Liu<sup>2</sup>, Jian Li<sup>2</sup>,  
Yabiao Wang<sup>2</sup>, Chengjie Wang<sup>2,3†</sup>, Cai Rong Zhao<sup>1†</sup>

*1 Tongji University, 2 Tencent Youtu Lab, 3 Shanghai Jiao Tong University*

*Corresponding authors. Email: zhaocairong@tongji.edu.cn, jasoncjwang@tencent.com*