

# Sharpness-Aware Minimization Activates the Interactive Teaching's Understanding and Optimization

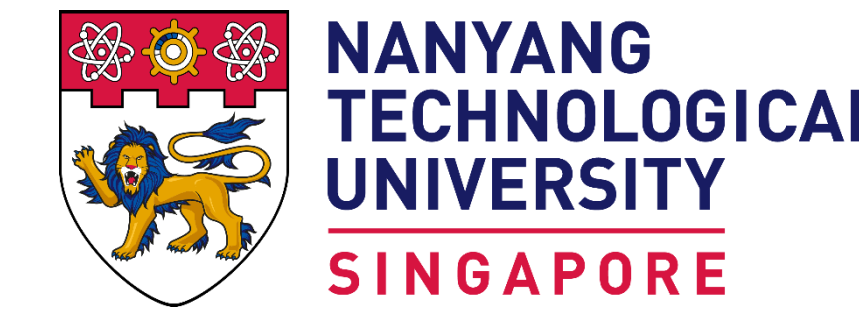
Mingwei Xu <sup>1</sup>, Xiaofeng Cao <sup>1,✉</sup>, Ivor W. Tsang <sup>2,3</sup>

<sup>1</sup> School of Artificial Intelligence, Jilin University, China

<sup>2</sup> CFAR and IHPC, Agency for Science, Technology and Research (A\*STAR), Singapore

<sup>3</sup> College of Computing and Data Science, Nanyang Technological University, Singapore

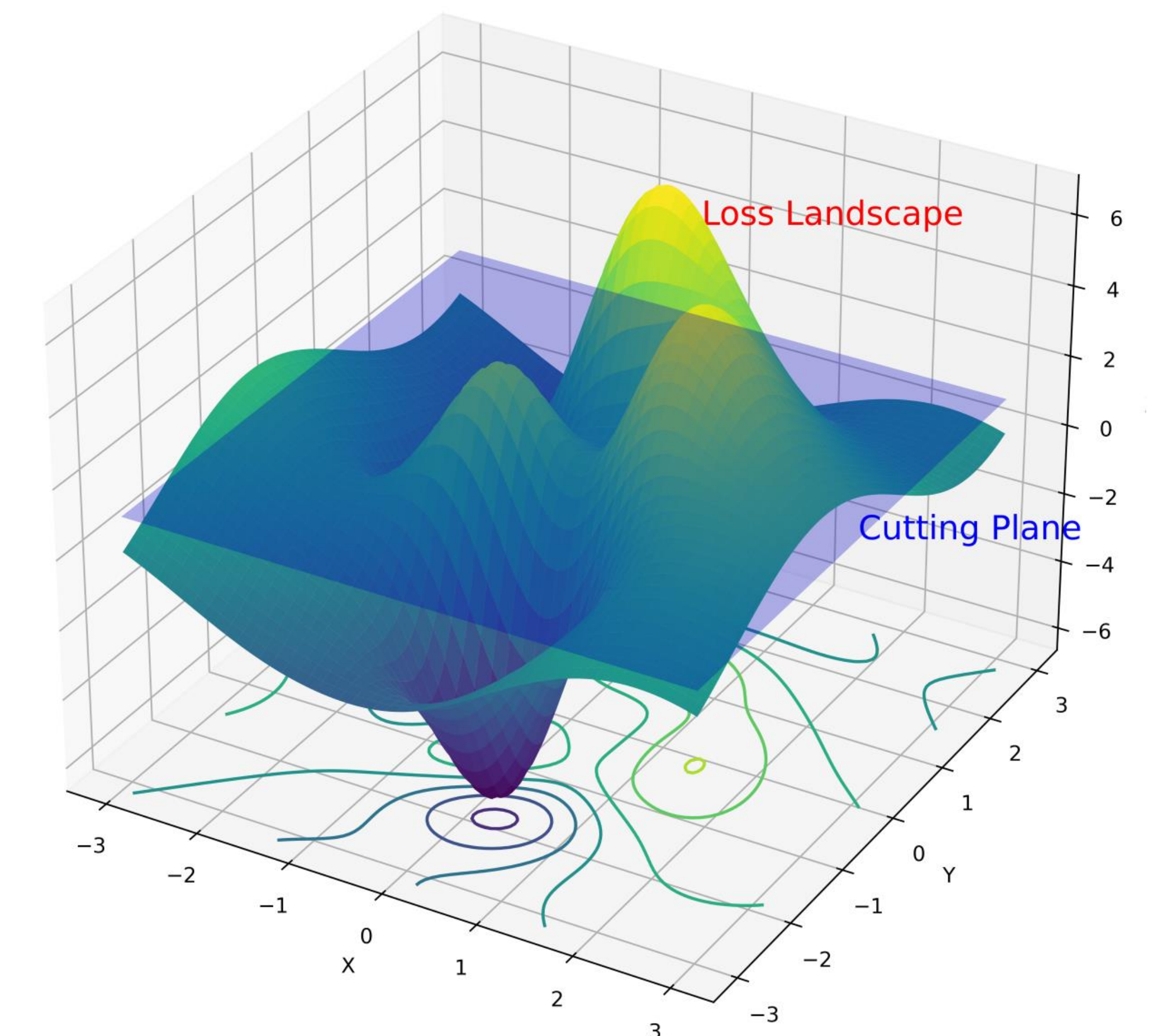
# Backgrounds



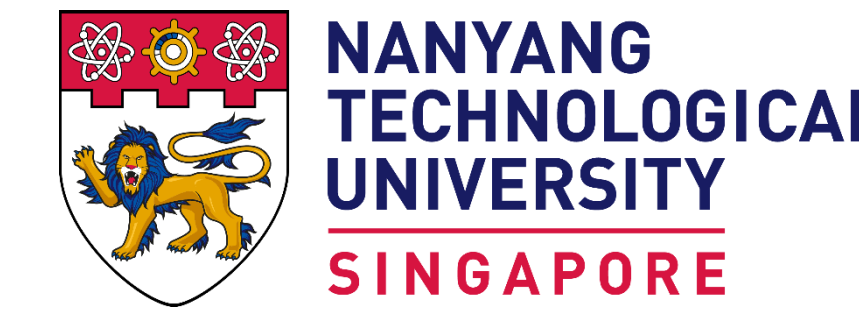
- Teaching recognized as a pervasive mechanism for disseminating knowledge within human society, has found extensive application in contemporary deep learning methodologies.
- It serves as a cornerstone for various techniques such as knowledge distillation, data distillation, model compression, and machine teaching, facilitating optimal training control.
- Recent investigations into pedagogy have illuminated the integration of large language models and multi-agent systems into educational frameworks.

## Interactive teaching's understanding

- ◆ Interactive teaching methods like co-teaching update parameters by reducing high loss values in the landscape.
- ◆ By actively involving two teachers, models in interactive framework learn from each other's strengths through a collaborative filtering mechanism and focus on minimizing loss examples.



# Introduction



- The underlying optimization principles and convergence of interactive teaching lack theoretical analysis, and co-teaching serves as a notable prototype in this regard.
- In this paper, we discuss its role as a reduction of the larger loss landscape derived from Sharpness-Aware Minimization (SAM).
- Then, we classify it as an iterative parameter estimation process using Expectation-Maximization. The convergence of this typical interactive teaching is achieved by continuously optimizing a variational lower bound on the log marginal likelihood.

## Our optimization

- ◆ To further enhance interactive teaching's performance, we incorporate SAM's strong generalization information into interactive teaching, referred as Sharpness Reduction Interactive Teaching (SRIT). This integration can be viewed as a novel sequential optimization process.
- ◆ We validate the performance of our approach through multiple experiments.

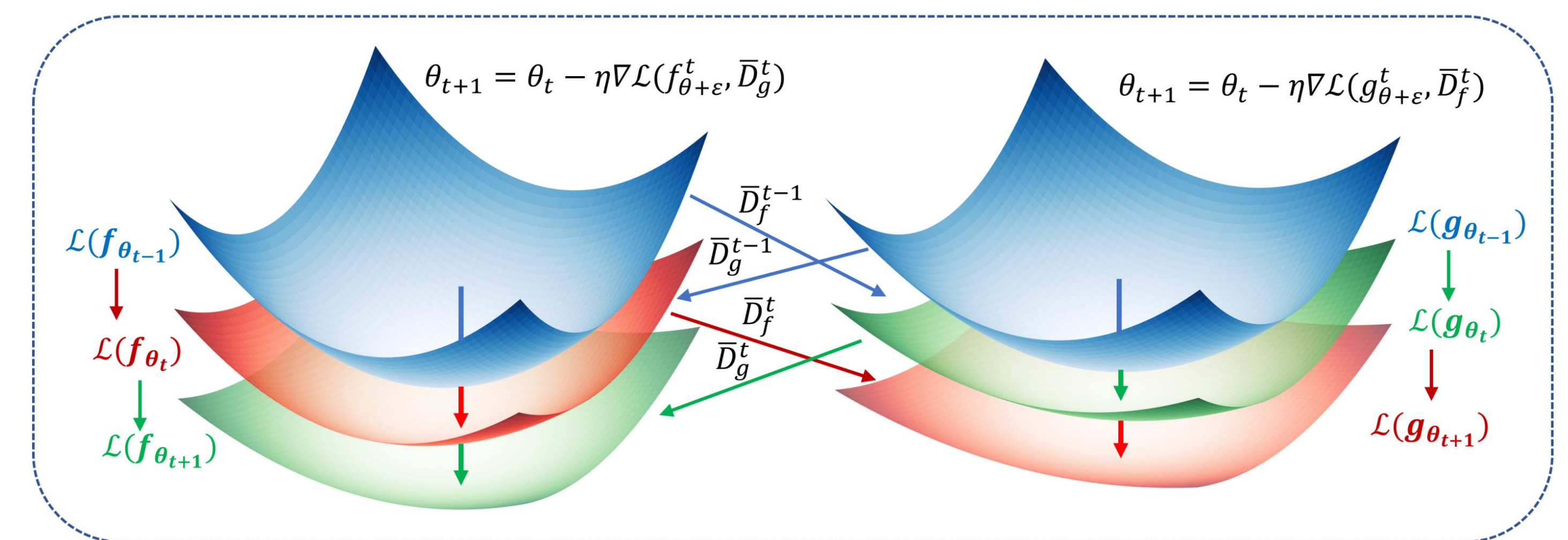
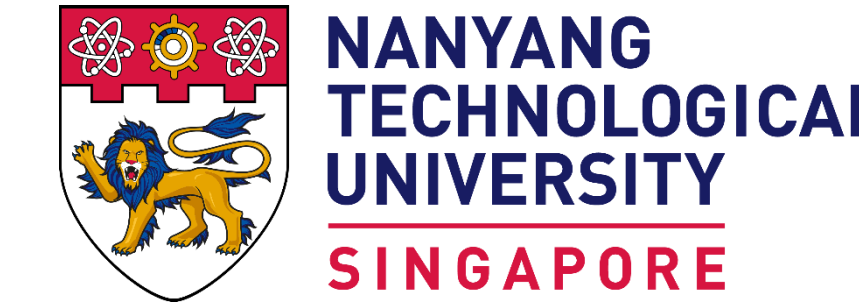


Figure 1: Interactive teaching, Sharpness Reduction Interactive Teaching (SRIT), the plane in the figure represents the loss landscape, which gradually becomes flat during the iterative optimization process due to the receipt of flat gradient information cues from each other.

# Theoretical Analysis



- ◆ Our core assumption is that the cleanliness of the data distribution serves as a latent variable, as it remains unknown within the training dataset.
- ◆ Based on this assumption, the interactive teaching process can be effectively exemplified as a unique type of parameter iteration within the EM framework.
- ◆ This perspective provides a probabilistic modeling-based explanation for the iteration and convergence of interactive teaching, such as co-teaching.

## The EM process

**Proposition 4.1.** Given the training dataset  $\mathcal{D} = \{X_i = (x_i, y_i)\}_{i=1}^N$ , which contains noisy samples, and assuming the samples are independent, we define the hidden variable  $Z_c = 1$  to indicate that the  $c$ -th sample is a cleaner sample, meaning it has a lower loss value compared to noisy data.  $Z = Z_c^f \cup Z_c^g = \{Z_c\}_{c=1}^K$  represents the set of hidden variables for all samples, and the corresponding latent distribution is denoted as  $q(Z)$ . The joint probability of  $p(X_i, Z_c | \theta_f, \theta_g)$  is obtained by simultaneously updating neural networks  $f$  and  $g$  for  $X_i$  and  $Z_c$ . The logarithm likelihood of the observed data  $\mathcal{D}$  has the following lower bound  $L$ :

$$L(\theta_f, \theta_g, q) \equiv \sum_{i=1}^N \sum_{c=1}^K \left[ q(Z_c) \log \frac{p(X_i, Z_c | \theta_f, \theta_g)}{q(Z_c)} \right], \quad (12)$$

the EM algorithm approximates the maximization of  $\log p(\mathcal{D} | \theta_f, \theta_g)$  by maximizing this lower bound  $L(\theta_f, \theta_g, q)$ . Specifically, the EM iteration process for the interactive teaching paradigm is as follows:

**E-step:**

$$Q(\theta_f, \theta_g^{(t)}) = \mathbb{E}_{Z_c^f | \mathcal{D}, \theta_f^{(t)}, \theta_g^{(t)}} [\log p(Z_c^f, \mathcal{D} | \theta_f, \theta_g^{(t)})], \quad (13)$$

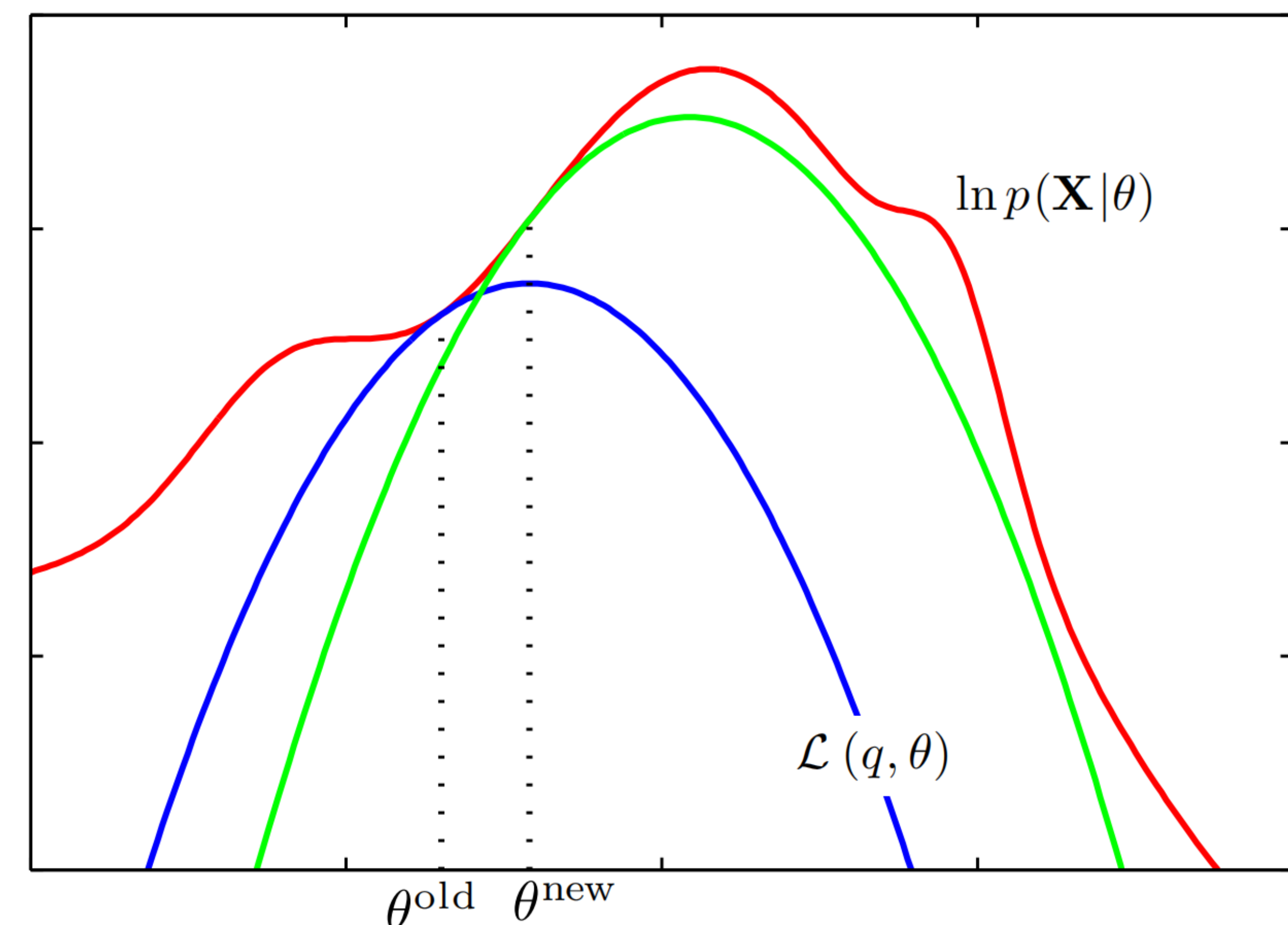
$$Q(\theta_f^{(t)}, \theta_g) = \mathbb{E}_{Z_c^g | \mathcal{D}, \theta_f^{(t)}, \theta_g^{(t)}} [\log p(Z_c^g, \mathcal{D} | \theta_f^{(t)}, \theta_g)], \quad (14)$$

The subscript  $Z_c^f | \mathcal{D}, \theta_f^{(t)}, \theta_g^{(t)}$  (or  $Z_c^g | \mathcal{D}, \theta_f^{(t)}, \theta_g^{(t)}$ ) of the expectation represents the corresponding set of low-loss samples selected by network  $f$  (or  $g$ ) and is used to update network  $g$  (or  $f$ ).

**M-step:**

$$\theta_f^{(t+1)} = \arg \max_{\theta_f} Q(\theta_f^{(t)}, \theta_g), \quad (15)$$

$$\theta_g^{(t+1)} = \arg \max_{\theta_g} Q(\theta_f, \theta_g^{(t)}). \quad (16)$$



- ◆ Since the convergence of the EM algorithm guarantees only local optima, while SAM can flatten the loss landscape and effectively alleviate local optima, favoring global optima, we incorporate SAM into the interactive teaching process, which referred as Sharpness Reduction Interactive Teaching (SRIT).

## The Key steps

(1) For the first level, we start by screening the required low-loss dataset  $\min_{\hat{\mathcal{D}}} \mathcal{L}(f_{\theta}, \hat{\mathcal{D}})$ .

(2) For the second level, we transfer the loss information to the counterpart model for SAM optimization:  $\theta^* = \arg \min_{\hat{\theta}} \max_{\hat{\epsilon}} \mathcal{L}(g_{\theta+\epsilon}, \hat{\mathcal{D}})$ .

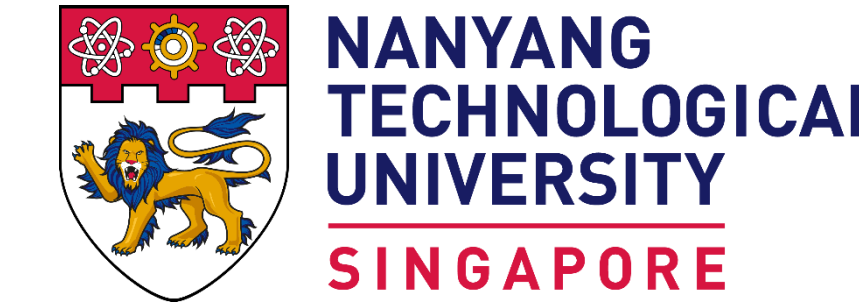
The first step estimates the direction of the change in the network weights  $\hat{\epsilon}(\theta_f)$  and  $\hat{\epsilon}(\theta_g)$ , and then computes  $\mathbf{G}_f$  and  $\mathbf{G}_g$ ,

$$\mathbf{G}_f = \nabla_{\theta_f} \mathcal{L}(f_{\theta}, \hat{\mathcal{D}}_g)|_{\theta_f + \hat{\epsilon}(\theta_f)}, \mathbf{G}_g = \nabla_{\theta_g} \mathcal{L}(g_{\theta}, \hat{\mathcal{D}}_f)|_{\theta_g + \hat{\epsilon}(\theta_g)}, \quad (17)$$

$$\text{where } \hat{\epsilon}(\theta_f) = \rho \frac{\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_g}(f_{\theta})}{\|\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_g}(f_{\theta})\|_2}, \hat{\epsilon}(\theta_g) = \rho \frac{\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_f}(g_{\theta})}{\|\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_f}(g_{\theta})\|_2}.$$

The second step involves updating the parameters of networks  $f$  and  $g$  based on the estimated gradients.

# Algorithm & Experiments



## Algorithm 1: Sharpness Reduction Interactive Teaching (SRIT)

**Input:** Initial network parameters  $\theta_{f_0}, \theta_{g_0}$ , learning rate  $\eta$ , fixed parameter  $\tau$ , iteration counts  $T_k$  and  $T_{\max}$ , maximum iteration count  $N_{\max}$ , pre-defined constant  $\rho$ .

**Output:** Updated network parameters  $\theta_f$  and  $\theta_g$ .

```

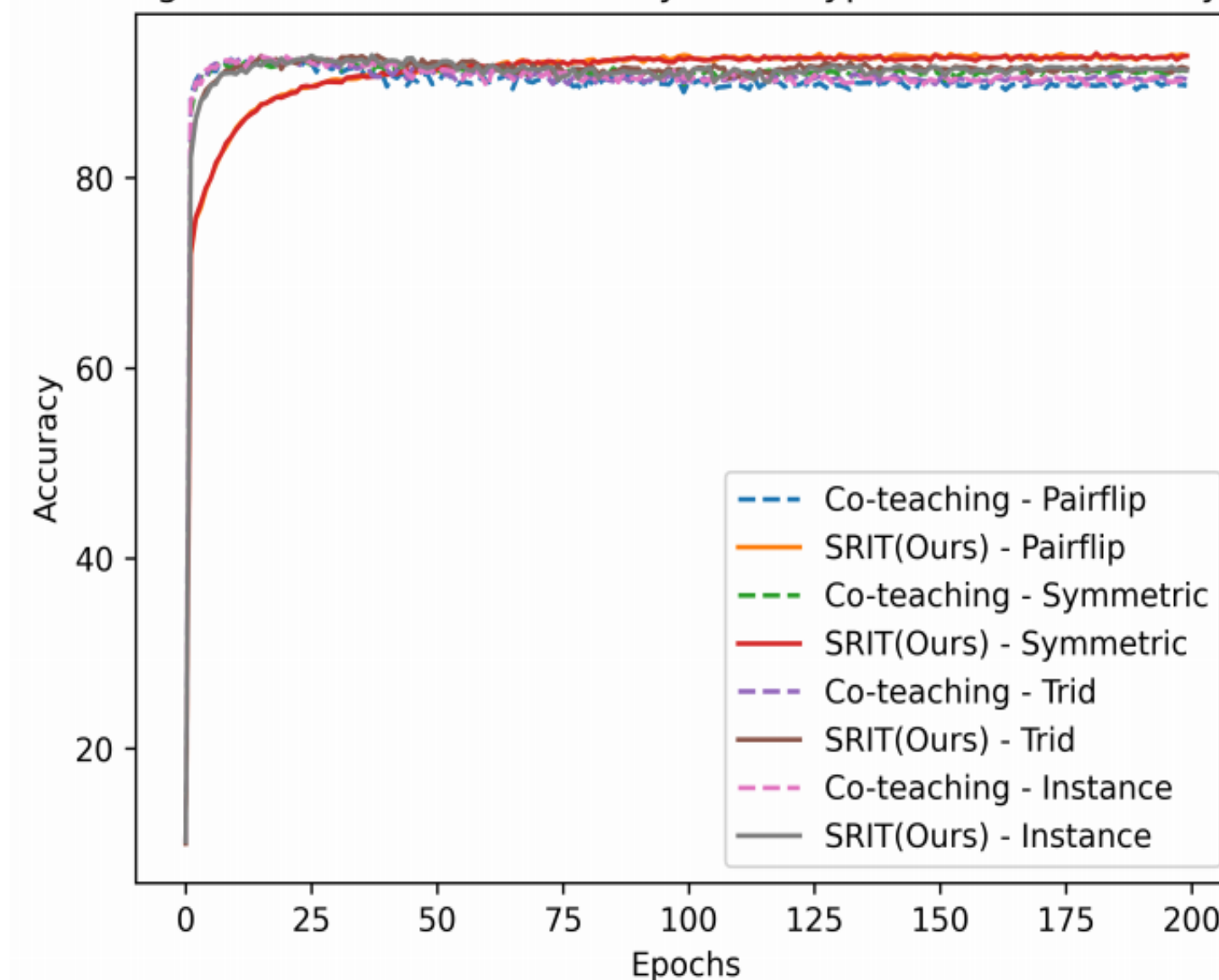
1 for  $T = 1$  to  $T_{\max}$  do
2   Shuffle the training set  $\mathcal{D}$  (noisy dataset);
3   for  $N = 1$  to  $N_{\max}$  do
4     Sample a mini-batch  $\bar{\mathcal{D}}$  from  $\mathcal{D}$ ;
5     Dual-level optimization: The first level of loss information exchange.
6     Compute the loss of network  $f$  on  $\bar{\mathcal{D}}$  and obtain  $\hat{\mathcal{D}}_f$ :
7      $\hat{\mathcal{D}}_f = \arg \min_{\mathcal{D}': |\mathcal{D}'| \geq R(T)|\bar{\mathcal{D}}|} \mathcal{L}(f, \bar{\mathcal{D}})$ ; //sample  $R(T) \cdot |\bar{\mathcal{D}}|$  small-loss instances;
8     Compute the loss of network  $g$  on  $\bar{\mathcal{D}}$  and obtain  $\hat{\mathcal{D}}_g$ :
9      $\hat{\mathcal{D}}_g = \arg \min_{\mathcal{D}': |\mathcal{D}'| \geq R(T)|\bar{\mathcal{D}}|} \mathcal{L}(g, \bar{\mathcal{D}})$  //sample  $R(T) \cdot |\bar{\mathcal{D}}|$  small-loss instances;
10    Dual-level optimization : The second level of sharpness element exchange.
11    Update  $\hat{\epsilon}(\theta_f) = \rho \frac{\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_g}(f_{\theta})}{\|\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_g}(f_{\theta})\|_2}$  and  $\hat{\epsilon}(\theta_g) = \rho \frac{\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_f}(g_{\theta})}{\|\nabla_{\theta} \mathcal{L}_{\hat{\mathcal{D}}_f}(g_{\theta})\|_2}$ ;
12    Compute the approximate gradient for network  $f$ :  $\mathbf{G}_f = \nabla_{\theta_f} \mathcal{L}(f_{\theta}, \hat{\mathcal{D}}_g)|_{\theta_f + \hat{\epsilon}(\theta_f)}$ ;
13    Compute the approximate gradient for network  $g$ :  $\mathbf{G}_g = \nabla_{\theta_g} \mathcal{L}(g_{\theta}, \hat{\mathcal{D}}_f)|_{\theta_g + \hat{\epsilon}(\theta_g)}$ ;
14    Update the network parameters of  $f$  using gradient descent:  $\theta_f = \theta_f - \eta \mathbf{G}_f$ ;
15    Update the network parameters of  $g$  using gradient descent:  $\theta_g = \theta_g - \eta \mathbf{G}_g$ ;
16  end
17  Compute  $R(T) = 1 - \min \left\{ \frac{T}{T_k} \tau, \tau \right\}$ .
18 end

```

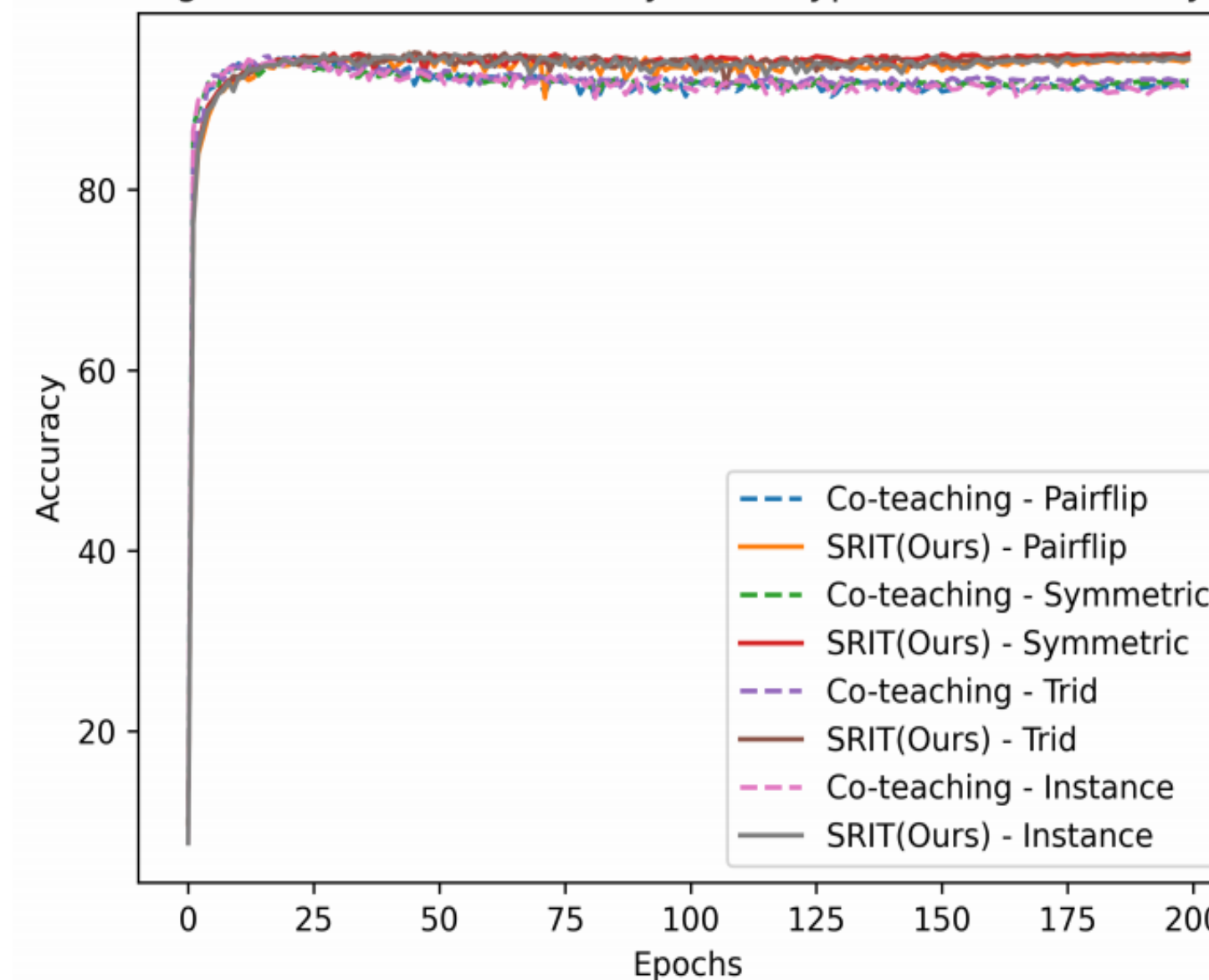
Table 1: Test accuracy (%) on five datasets. The best results are highlighted in bold.

Noise type	Symmetric.		Pairflip.		Tridiagonal.		Instance.	
Noise ratio	20%	40%	20%	40%	20%	40%	20%	40%
MNIST								
Co-teaching	97.50	94.96	95.49	91.54	96.61	92.76	95.90	91.23
	$\pm 0.06$	$\pm 0.07$	$\pm 0.11$	$\pm 0.15$	$\pm 0.06$	$\pm 0.09$	$\pm 0.05$	$\pm 0.18$
<b>SRIT</b>	<b>99.42</b>	<b>99.19</b>	<b>99.35</b>	<b>98.14</b>	<b>99.47</b>	<b>98.75</b>	<b>99.43</b>	<b>98.03</b>
	$\pm 0.03$	$\pm 0.03$	$\pm 0.02$	$\pm 0.07$	$\pm 0.03$	$\pm 0.05$	$\pm 0.02$	$\pm 0.10$
CIFAR10								
Co-teaching	82.15	77.38	82.32	75.37	82.77	76.41	81.86	73.61
	$\pm 0.09$	$\pm 0.15$	$\pm 0.08$	$\pm 0.14$	$\pm 0.07$	$\pm 0.17$	$\pm 0.12$	$\pm 0.25$
<b>SRIT</b>	<b>85.64</b>	<b>79.83</b>	<b>85.10</b>	<b>76.95</b>	<b>85.39</b>	<b>78.90</b>	<b>84.77</b>	<b>74.07</b>
	$\pm 0.15$	$\pm 0.12$	$\pm 0.20$	$\pm 0.17$	$\pm 0.18$	$\pm 0.12$	$\pm 0.19$	$\pm 0.25$
CIFAR100								
Co-teaching	50.21	42.40	48.27	34.74	50.32	38.78	49.74	38.57
	$\pm 0.23$	$\pm 0.16$	$\pm 0.11$	$\pm 0.13$	$\pm 0.19$	$\pm 0.16$	$\pm 0.18$	$\pm 0.12$
<b>SRIT</b>	<b>59.66</b>	<b>50.57</b>	<b>57.16</b>	<b>35.82</b>	<b>59.07</b>	<b>42.27</b>	<b>59.66</b>	<b>40.36</b>
	$\pm 0.16$	$\pm 0.21$	$\pm 0.10$	$\pm 0.16$	$\pm 0.15$	$\pm 0.22$	$\pm 0.16$	$\pm 0.18$
FMNIST								
Co-teaching	91.13	87.99	89.83	85.44	90.42	86.09	90.27	85.63
	$\pm 0.09$	$\pm 0.09$	$\pm 0.10$	$\pm 0.12$	$\pm 0.07$	$\pm 0.09$	$\pm 0.12$	$\pm 0.13$
<b>SRIT</b>	<b>92.68</b>	<b>88.77</b>	<b>92.76</b>	<b>89.25</b>	<b>91.44</b>	<b>89.97</b>	<b>91.44</b>	<b>86.02</b>
	$\pm 0.10$	$\pm 0.13$	$\pm 0.09$	$\pm 0.11$	$\pm 0.11$	$\pm 0.09$	$\pm 0.09$	$\pm 0.21$
SVHN								
Co-teaching	91.83	88.72	91.49	85.09	92.16	87.51	91.26	86.33
	$\pm 0.08$	$\pm 0.10$	$\pm 0.10$	$\pm 0.15$	$\pm 0.10$	$\pm 0.13$	$\pm 0.18$	$\pm 0.23$
<b>SRIT</b>	<b>94.95</b>	<b>93.06</b>	<b>94.34</b>	<b>89.37</b>	<b>94.66</b>	<b>91.56</b>	<b>94.45</b>	<b>90.50</b>
	$\pm 0.05$	$\pm 0.05$	$\pm 0.08$	$\pm 0.20$	$\pm 0.05$	$\pm 0.12$	$\pm 0.08$	$\pm 0.10$

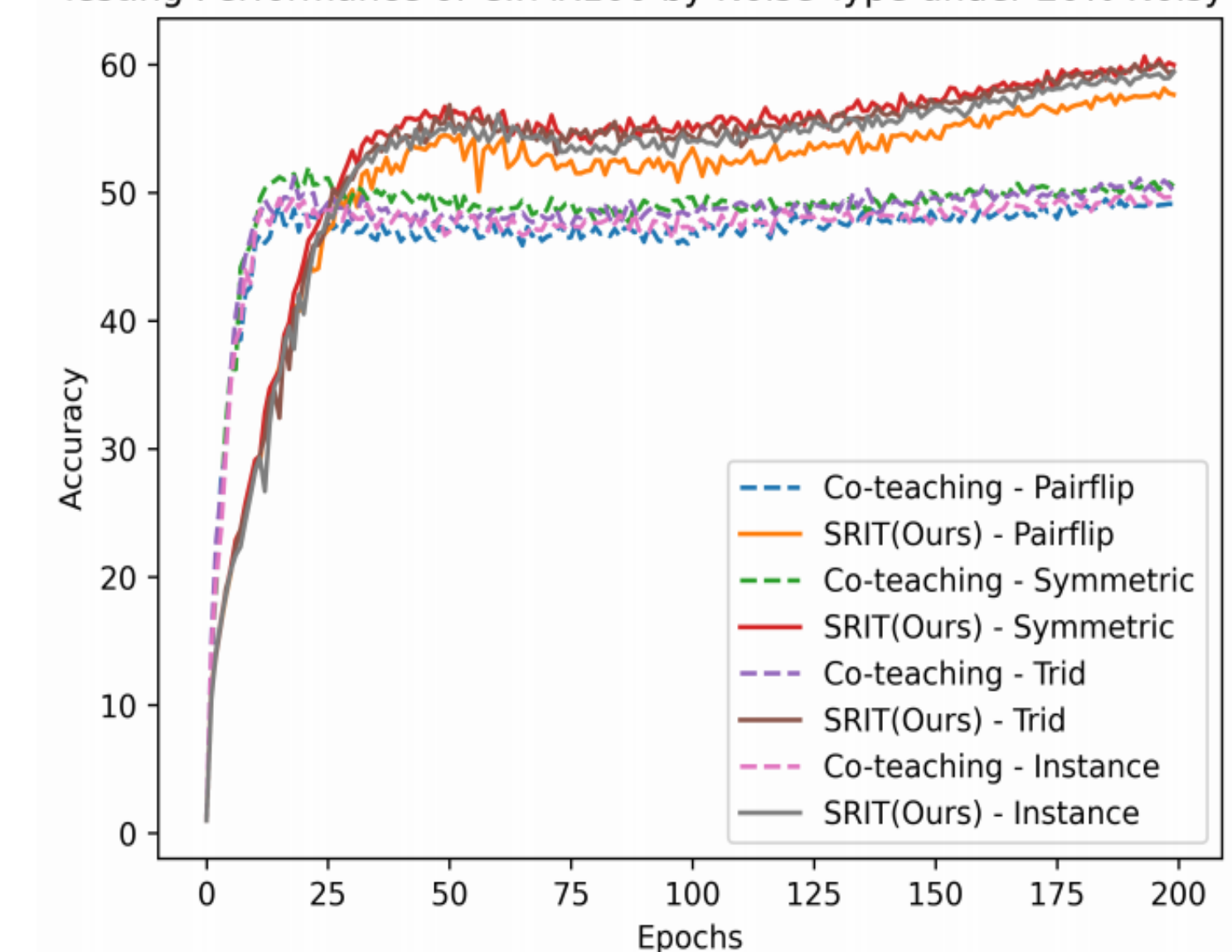
Testing Performance of FMNIST by Noise Type Under 20% Noisy Data



Testing Performance of SVHN by Noise Type Under 20% Noisy Data



Testing Performance of CIFAR100 by Noise Type under 20% Noisy Data





---

*Thank you!*