

# Comparing Kullback-Leibler Divergence and Mean Squared Error Loss in Knowledge Distillation

---

2021.06.04

Taehyeon Kim\*

Jaehoon Oh\*

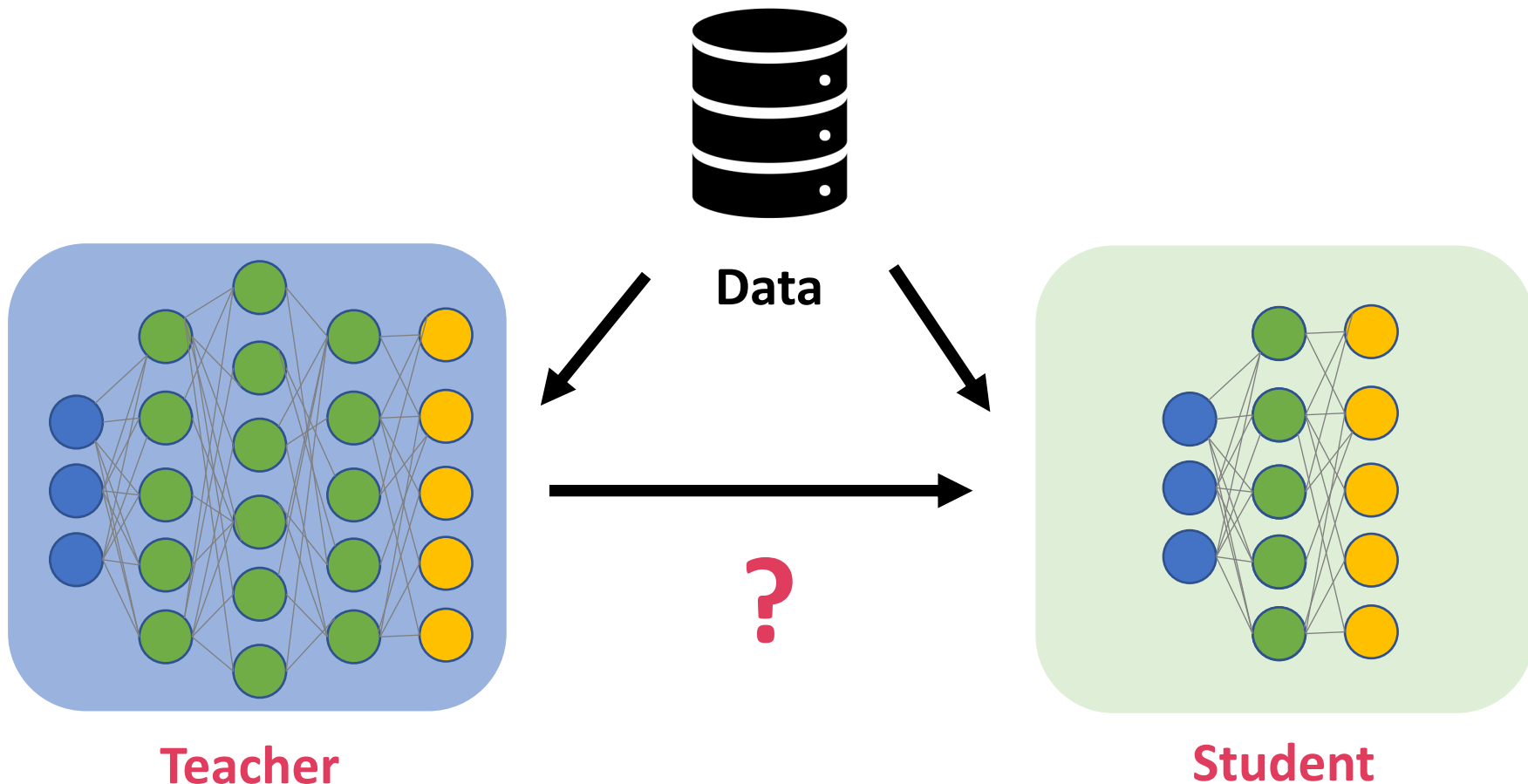
Nakyil Kim

Sangwook Cho

Se-young Yun

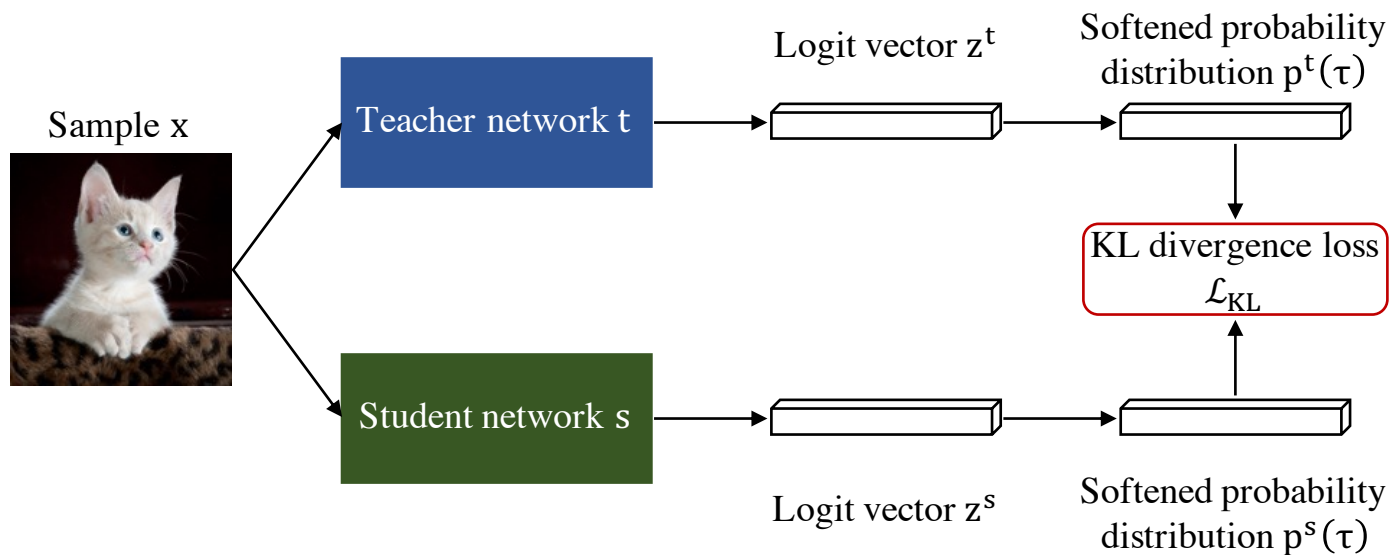
# Why is Knowledge Distillation (KD) Beneficial

- One of the most potent **model compression** techniques.
- Knowledge is transferred from a **cumbersome** model (**teacher**) to a single **small** model (**student**).



# Overview of KD

- KD has evolved to design a new objective function
  - KL divergence loss with temperature scaling



# Degree of softness (temperature scaling)

- KD has evolved to design a new objective function
  - KL divergence loss with temperature scaling
  - Little Understanding of how the degree of softness affects the performance.

1. Learning from Ground Truth

2. Learning from Teacher model

$$\mathcal{L} = (1 - \alpha)\mathcal{L}_{CE}(\mathbf{p}^s(1), \mathbf{y}) + \alpha\mathcal{L}_{KL}(\mathbf{p}^s(\tau), \mathbf{p}^t(\tau)),$$

$$\mathcal{L}_{CE}(\mathbf{p}^s(1), \mathbf{y}) = \sum_j -y_j \log p_j^s(1)$$

$$\mathcal{L}_{KL}(\mathbf{p}^s(\tau), \mathbf{p}^t(\tau)) = \tau^2 \sum_j p_j^t(\tau) \log \frac{p_j^t(\tau)}{p_j^s(\tau)}$$

$$p_k^f(\tau) = \frac{\exp(\mathbf{z}_k^f / \tau)}{\sum_{j=1}^K \exp(\mathbf{z}_j^f / \tau)}$$

# Preliminaries: respects to previous assumption

- Conventional assumption for KD

**Derivative**

$$\frac{\partial \mathcal{L}_{KL}}{\partial z_k^s} \approx \tau \left( \frac{1 + z_k^s/\tau}{K + \sum_j z_j^s/\tau} - \frac{1 + z_k^t/\tau}{K + \sum_j z_j^t/\tau} \right)$$

**Assumption**

$$\sum_j z_j^s = 0 \text{ and } \sum_j z_j^t = 0$$

**Conclusion**

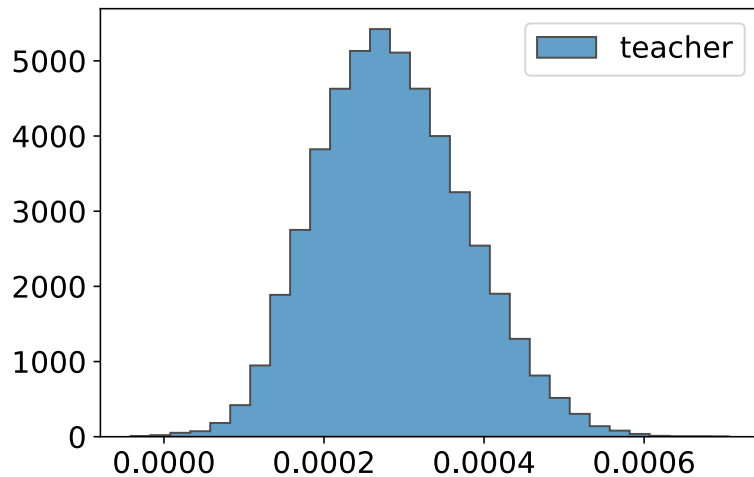
$$\frac{\partial \mathcal{L}_{KL}}{\partial z_k^s} \approx \frac{1}{K} (z_k^s - z_k^t) \quad \text{then, } \mathcal{L}_{KL} = \mathcal{L}_{MSE} \text{ } \text{????}$$

# Preliminaries: respects to previous assumption

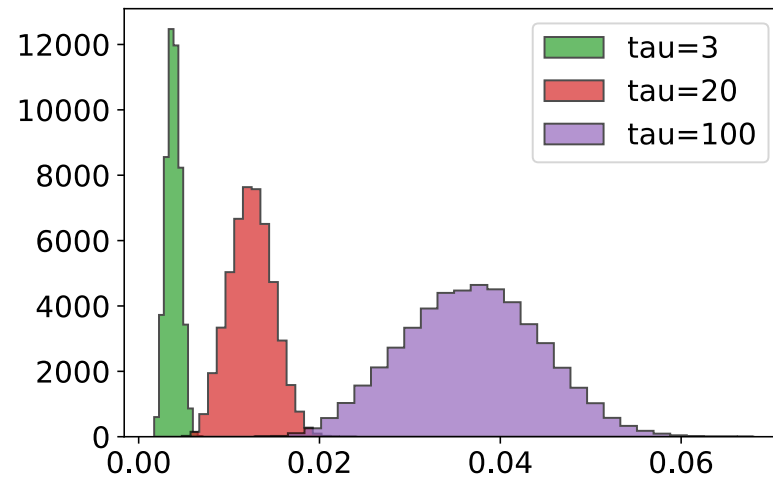
- Conventional assumption does **not** seem **appropriate**.

**Assumption**

$$\sum_j z_j^s = 0 \text{ and } \sum_j z_j^t = 0$$



$\neq$



$\neq 0$

# Changes on $\tau$

- Depending on  $\tau$ ,  $\mathcal{L}_{KL}$  plays different roles
  - $\tau \rightarrow 0$  : Label matching
  - $\tau \rightarrow \infty$  : Logit matching



$$\lim_{\tau \rightarrow 0} \frac{1}{\tau} \frac{\partial \mathcal{L}_{KL}}{\partial \mathbf{z}_k^s} = \mathbf{1}_{[\arg \max_j \mathbf{z}_j^s = k]} - \mathbf{1}_{[\arg \max_j \mathbf{z}_j^t = k]}$$

$$\begin{aligned} \lim_{\tau \rightarrow \infty} \frac{\partial \mathcal{L}_{KL}}{\partial \mathbf{z}_k^s} &= \frac{1}{K^2} \sum_{j=1}^K ((\mathbf{z}_k^s - \mathbf{z}_j^s) - (\mathbf{z}_k^t - \mathbf{z}_j^t)) \\ &= \frac{1}{K} (\mathbf{z}_k^s - \mathbf{z}_k^t) - \frac{1}{K^2} \sum_{j=1}^K (\mathbf{z}_j^s - \mathbf{z}_j^t) \end{aligned}$$

# Extension from KL loss to MSE Loss

- Some term is generated.
- *This term hinders complete logit matching (MSE loss) by shifting the mean of the elements in the logit.*

$$\begin{aligned}\lim_{\tau \rightarrow \infty} \frac{\partial \mathcal{L}_{KL}}{\partial \mathbf{z}_k^s} &= \frac{1}{K^2} \sum_{j=1}^K ((\mathbf{z}_k^s - \mathbf{z}_j^s) - (\mathbf{z}_k^t - \mathbf{z}_j^t)) \\ &= \frac{1}{K} (\mathbf{z}_k^s - \mathbf{z}_k^t) - \frac{1}{K^2} \sum_{j=1}^K (\mathbf{z}_j^s - \mathbf{z}_j^t)\end{aligned}$$



$$\lim_{\tau \rightarrow \infty} \nabla_{\mathbf{z}^s} \mathcal{L}_{KL} = \frac{1}{K} (\mathbf{z}^s - \mathbf{z}^t) - \frac{1}{K^2} \sum_{j=1}^K (\mathbf{z}_j^s - \mathbf{z}_j^t) \cdot \mathbb{1}$$



Bounded Convergence Theorem using each partial derivatives

$$\lim_{\tau \rightarrow \infty} \mathcal{L}_{KL} = \frac{1}{2K} \|\mathbf{z}^s - \mathbf{z}^t\|_2^2 + \delta_\infty = \frac{1}{2K} \mathcal{L}_{MSE} + \delta_\infty$$

$$\delta_\infty = -\frac{1}{2K^2} \left( \sum_{j=1}^K \mathbf{z}_j^s - \sum_{j=1}^K \mathbf{z}_j^t \right)^2 + Constant \longrightarrow \text{Analysis on this term!!}$$



# Theoretical Analysis on $\delta_\infty$

---

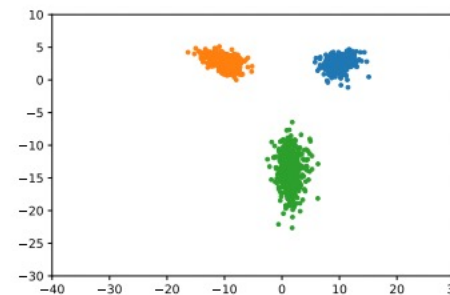
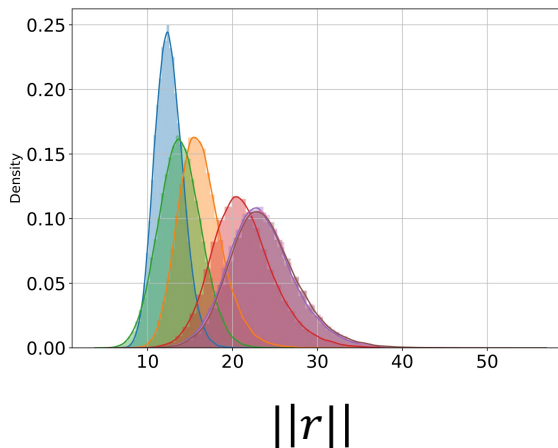
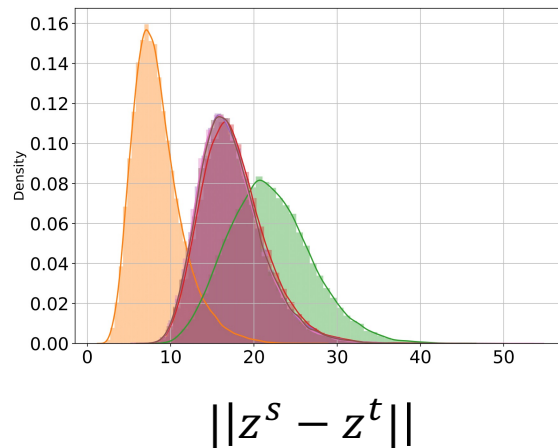
- Lower bound for  $\delta_\infty$

$$\begin{aligned}
 \delta_\infty &\approx -\frac{1}{2K^2} \left( \sum_{j=1}^K \mathbf{z}_j^s \right)^2 = -\frac{1}{2K^2} \left( \sum_{j=1}^K \sum_{n=1}^d W_{j,n}^s \mathbf{r}_n^s \right)^2 \\
 &= -\frac{1}{2K^2} \left( \sum_{n=1}^d \mathbf{r}_n^s \sum_{j=1}^K W_{j,n}^s \right)^2 \\
 &\geq -\frac{1}{2K^2} \left( \sum_{n=1}^d \left( \sum_{j=1}^K W_{j,n}^s \right)^2 \right) \left( \sum_{n=1}^d \mathbf{r}_n^{s2} \right) \\
 &\quad (\because \text{Cauchy-Schwartz inequality}) \\
 &= -\frac{1}{2K^2} \|\mathbf{r}^s\|_2^2 \left( \sum_{n=1}^d \left( \sum_{j=1}^K W_{j,n}^s \right)^2 \right)
 \end{aligned}$$

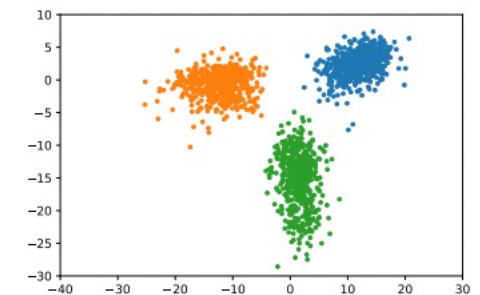
- Increasing the norm of  $r$  (pre-logit: input of fully-connected layer)
- De-shrinkage effects on weight templates

# Empirical Analysis on $\delta_\infty$

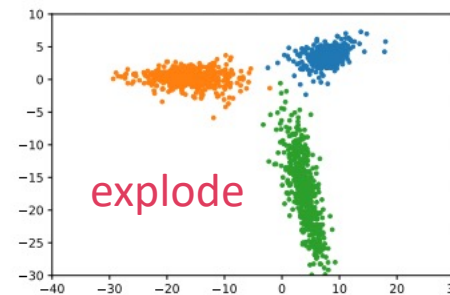
- Logit & Pre-logit behavior
- 2-D Projection Visualization



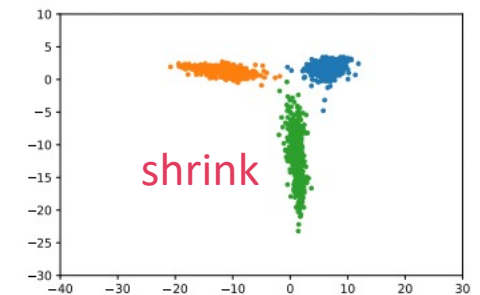
(a)  $t, \mathcal{L}_{CE}$  (Train)



(b)  $s, \mathcal{L}_{CE}$  (Train)



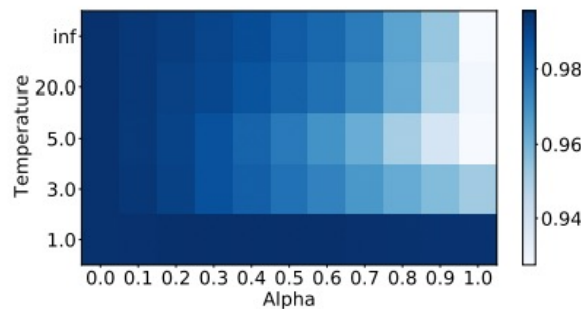
(c)  $s, \mathcal{L}_{KL}(\tau = \infty)$  (Train)



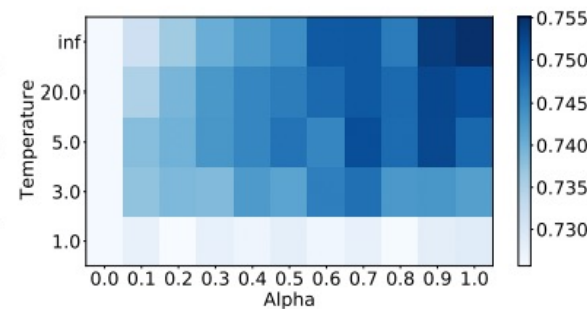
(d)  $s, \mathcal{L}_{MSE}$  (Train)

# Empirical Results

- Training Accuracy and Test Accuracy (CIFAR-100)
- With perfectly trained teacher, **MSE is the best!!**



(a) Training accuracy.



(b) Test accuracy.

Student	$\mathcal{L}_{CE}$	$\mathcal{L}_{KL}$					$\mathcal{L}_{MSE}$
		$\tau=1$	$\tau=3$	$\tau=5$	$\tau=20$	$\tau=\infty$	
WRN-16-2	72.68	72.90	74.24	74.88	75.15	75.51	<b>75.54</b>
WRN-16-4	77.28	76.93	78.76	78.65	78.84	78.61	<b>79.03</b>
WRN-28-2	75.12	74.88	76.47	76.60	<b>77.28</b>	76.86	<b>77.28</b>
WRN-28-4	78.88	78.01	78.84	79.36	79.72	79.61	<b>79.79</b>
WRN-40-6	79.11	79.69	79.94	79.87	79.82	79.80	<b>80.25</b>

# Empirical Results

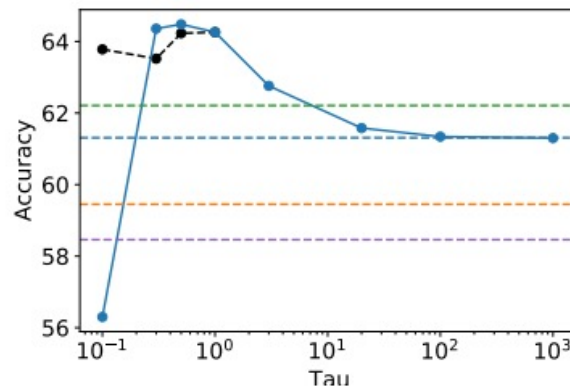
- **MSE is also the best compared to other alternatives!!**

Student	Baseline	SKD [2015]	FitNets [2014]	AT [2016a]	Jacobian [2018]	FT [2018]	AB [2019b]	Overhaul [2019a]	MSE
WRN-16-2	72.68	73.53	73.70	73.44	73.29	74.09	73.98	<b>75.59</b>	75.54
WRN-16-4	77.28	78.31	78.15	77.93	77.82	78.28	78.64	78.20	<b>79.03</b>
WRN-28-2	75.12	76.57	76.06	76.20	76.30	76.59	76.81	76.71	<b>77.28</b>

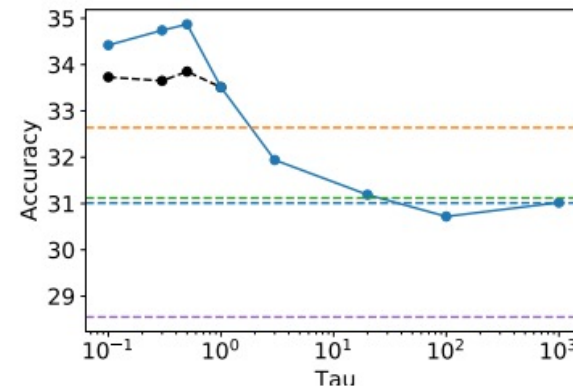
# Extreme small $\tau$

- Robustness to Noisy Labels
- On severe noise rate (80%), **low  $\tau$  is better than others!**
- It happens due to teacher's poor generalization.
- **Label Matching is better!**

$t, \mathcal{L}_{CE}$        $s, \mathcal{L}_{CE}$        $s, \mathcal{L}_{MSE}$   
 $s, \mathcal{L}_{KL}(\tau = \infty)$      $s, \mathcal{L}_{KL}$  in Eq.(2)     $s, \mathcal{L}_{KL}$  in Eq.(7)



(a) Symmetric noise 40%



(b) Symmetric noise 80%

# Adequate $\tau$ (Noisy Teacher)

- If your teacher model is not perfectly trained,
- But has the test accuracy around 80~90% then,
- the optimal solution for  $\tau$  may be some number  $>1$ .

Student	$\mathcal{L}_{KL}$						$\mathcal{L}_{MSE}$
	$\tau=0.1$	$\tau=0.5$	$\tau=1$	$\tau=5$	$\tau=20$	$\tau=\infty$	
WRN-16-2	51.64	52.07	51.36	50.11	49.69	49.46	49.20

Table 4: Top-1 test accuracies on CIFAR-100. WRN-28-4 is used as a teacher for  $\mathcal{L}_{KL}$  and  $\mathcal{L}_{MSE}$ . Here, the teacher (WRN-28-4) was not fully trained. The training accuracy of the teacher network is 53.77%.

Student	$\mathcal{L}_{CE}$	$\mathcal{L}_{KL}$ (Standard)	$\mathcal{L}_{KL}$ ( $\tau=20$ )	$\mathcal{L}_{MSE}$
ResNet-50	76.28	77.15	77.52	75.84

Table 5: Test accuracy on the ImageNet dataset. We used a (teacher, student) pair of (ResNet-152, ResNet-50). We include the results of the baseline and  $\mathcal{L}_{KL}$  (standard) from [Heo *et al.*, 2019a]. The training accuracy of the teacher network is 81.16%.

# Sequential Distillation

- According to the changes of objective functions,
- The performance varies significantly even under the usage of the same architectures.

WRN-28-4	WRN-16-4	WRN-16-2	Test accuracy
X	X	$\mathcal{L}_{CE}$	72.68 %
X	$\mathcal{L}_{CE}$ (77.28%)	$\mathcal{L}_{KL}(\tau = 3)$	74.84 %
		$\mathcal{L}_{KL}(\tau = 20)$	75.42 %
		$\mathcal{L}_{MSE}$	75.58 %
$\mathcal{L}_{CE}$ (78.88%)	X	$\mathcal{L}_{KL}(\tau = 3)$	74.24 %
		$\mathcal{L}_{KL}(\tau = 20)$	75.15 %
		$\mathcal{L}_{MSE}$	75.54 %
$\mathcal{L}_{CE}$ (78.88%)	$\mathcal{L}_{KL}(\tau = 3)$ (78.76%)	$\mathcal{L}_{KL}(\tau = 3)$	74.52 %
		$\mathcal{L}_{KL}(\tau = 20)$	75.47 %
		$\mathcal{L}_{MSE}$	<b>75.78 %</b>
	$\mathcal{L}_{MSE}$ (79.03%)	$\mathcal{L}_{KL}(\tau = 3)$	74.83 %
		$\mathcal{L}_{KL}(\tau = 20)$	75.47 %
		$\mathcal{L}_{MSE}$	75.48 %

# E.O.D.

---



Optimization and Statistical Inference LAB

potter32@kaist.ac.kr

[www.osi.kaist.ac.kr](http://www.osi.kaist.ac.kr)