



南京航空航天大学  
Nanjing University of Aeronautics and Astronautics

ParNeC

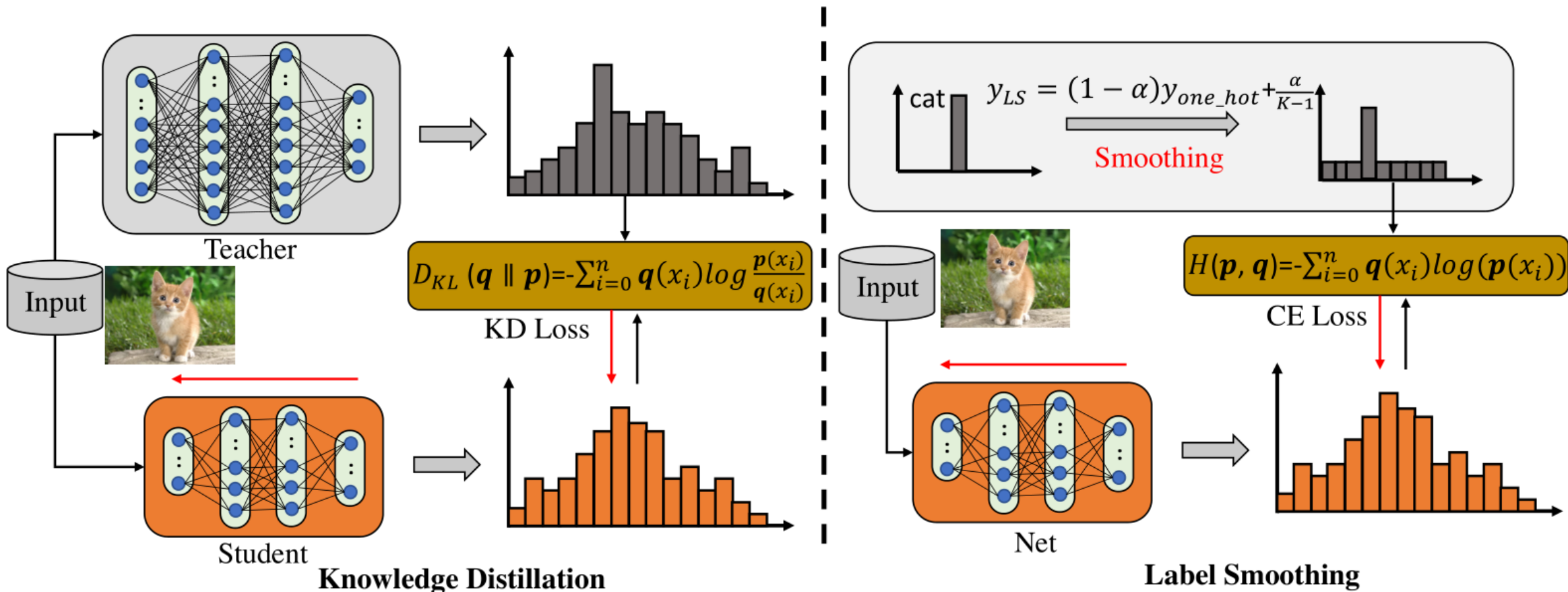
模式识别与神经计算研究组  
Pattern Recognition and NEural Computing

---

# Two Articles about Label Smoothing and Knowledge Distillation

---

# Knowledge Distillation & Label Smoothing



$$\lambda \mathcal{H}(\mathbf{p}^{S_w}, \mathbf{y}) + (1 - \lambda) \mathcal{H}(\mathbf{p}^{S_w} / \mathcal{T}, \mathbf{p}^{T_w} / \mathcal{T})$$



---

# When Does Label Smoothing Help?

---

**Rafael Müller\*, Simon Kornblith, Geoffrey Hinton**

Google Brain

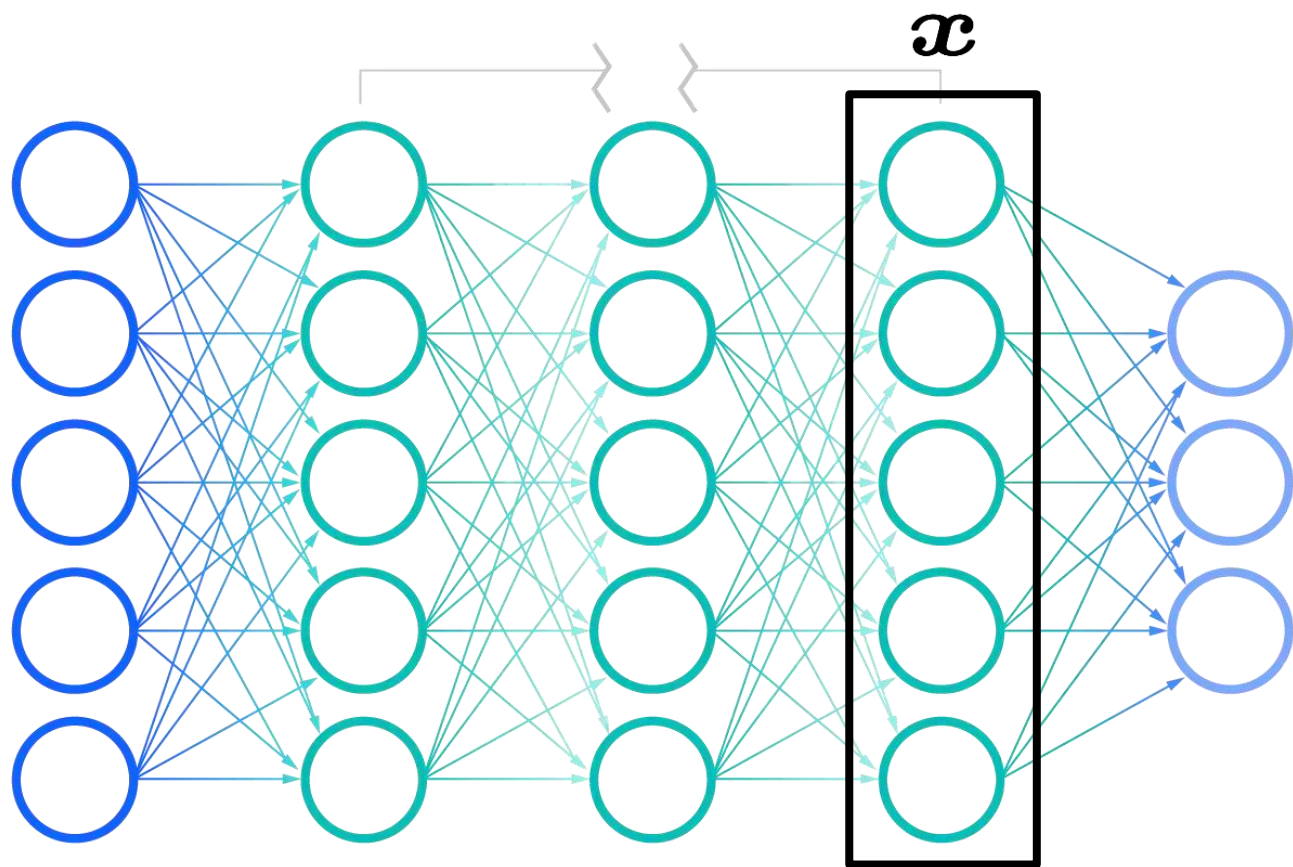
Toronto

`rafaelmuller@google.com`

NIPS 2019

# Survey of Literature Label Smoothing Results

DATA SET	ARCHITECTURE	METRIC	VALUE W/O LS	VALUE W/ LS
IMAGENET	INCEPTION-V2 [6]	TOP-1 ERROR	23.1	<b>22.8</b>
		TOP-5 ERROR	6.3	<b>6.1</b>
EN-DE	TRANSFORMER [11]	BLEU	25.3	<b>25.8</b>
		PERPLEXITY	<b>4.67</b>	4.92
WSJ	BiLSTM+ATT.[10]	WER	8.9	7.0/ <b>6.7</b>



$$p_k = \frac{e^{\mathbf{x}^T \boldsymbol{\omega}_k}}{\sum_{l=1}^L e^{\mathbf{x}^T \boldsymbol{\omega}_l}}$$

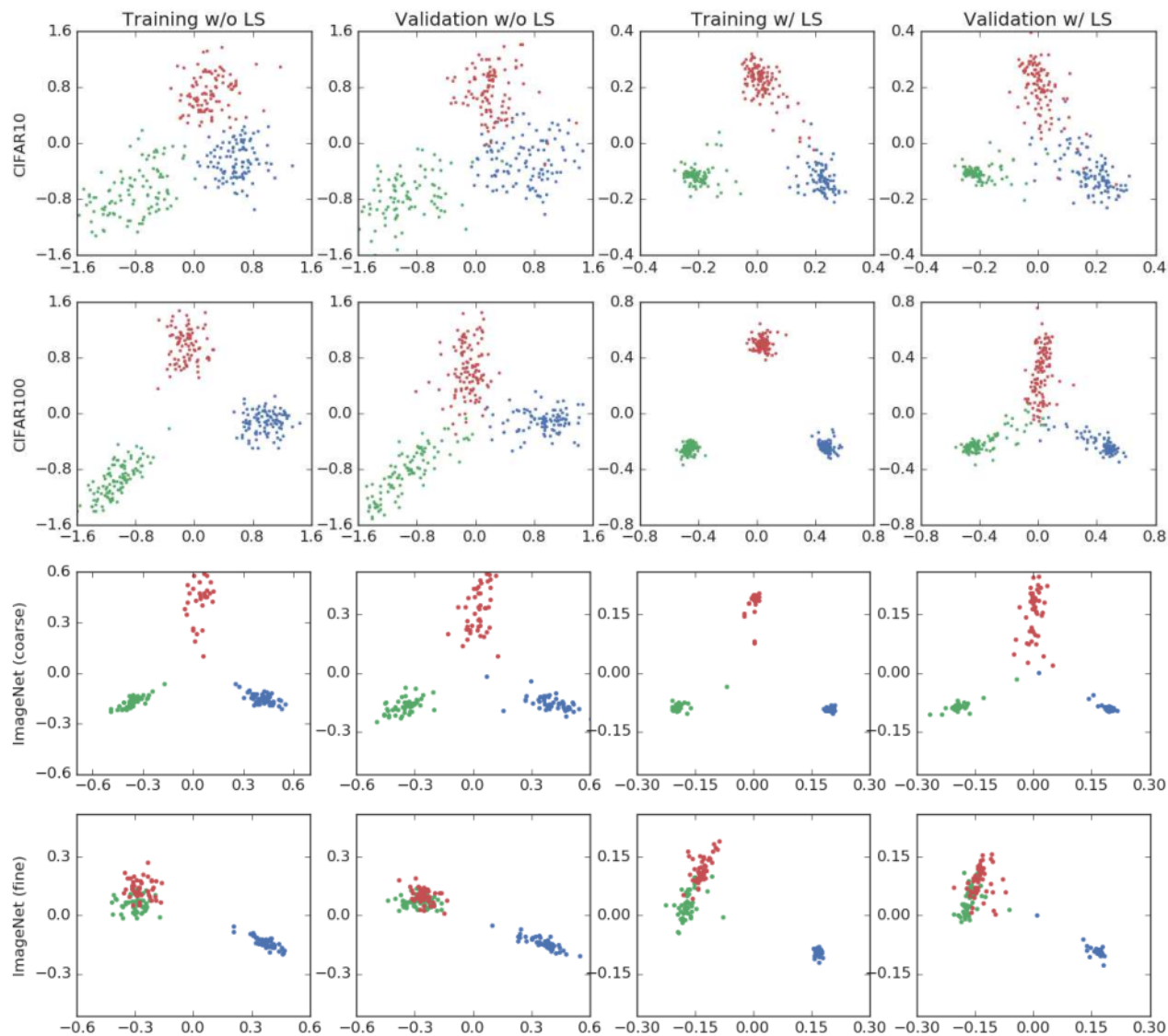
$$H(\mathbf{y}, \mathbf{p}) = \sum_{k=1}^K -y_k \log(p_k)$$

$$\|\mathbf{x} - \boldsymbol{\omega}_k\|^2 = \boxed{\mathbf{x}^T \mathbf{x}} - \mathbf{2x}^T \boldsymbol{\omega}_k + \boxed{\boldsymbol{\omega}_k^T \boldsymbol{\omega}_k}$$

Factored out

Usually constant across classes

# Penultimate Layer Representations

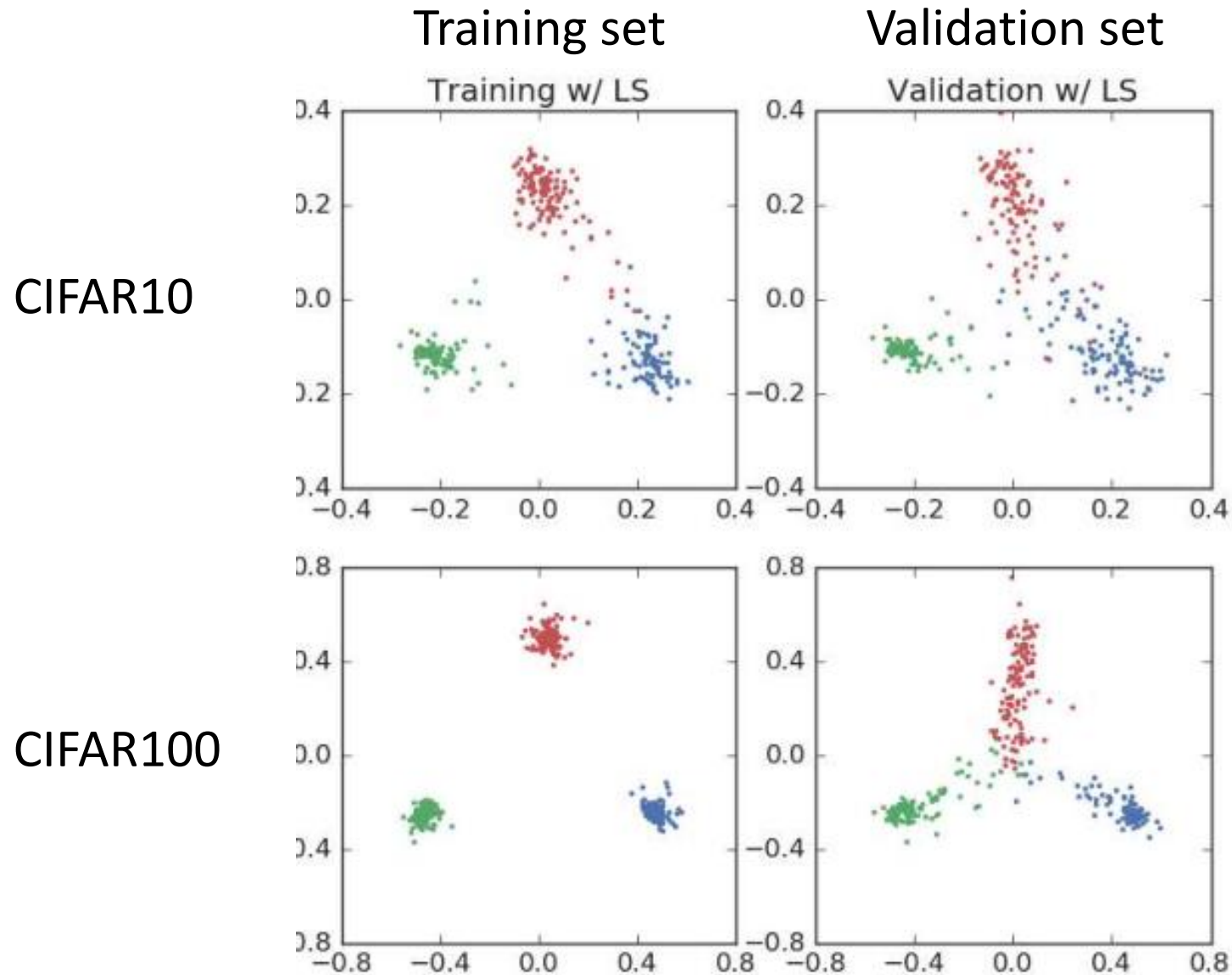


3 semantically different classes

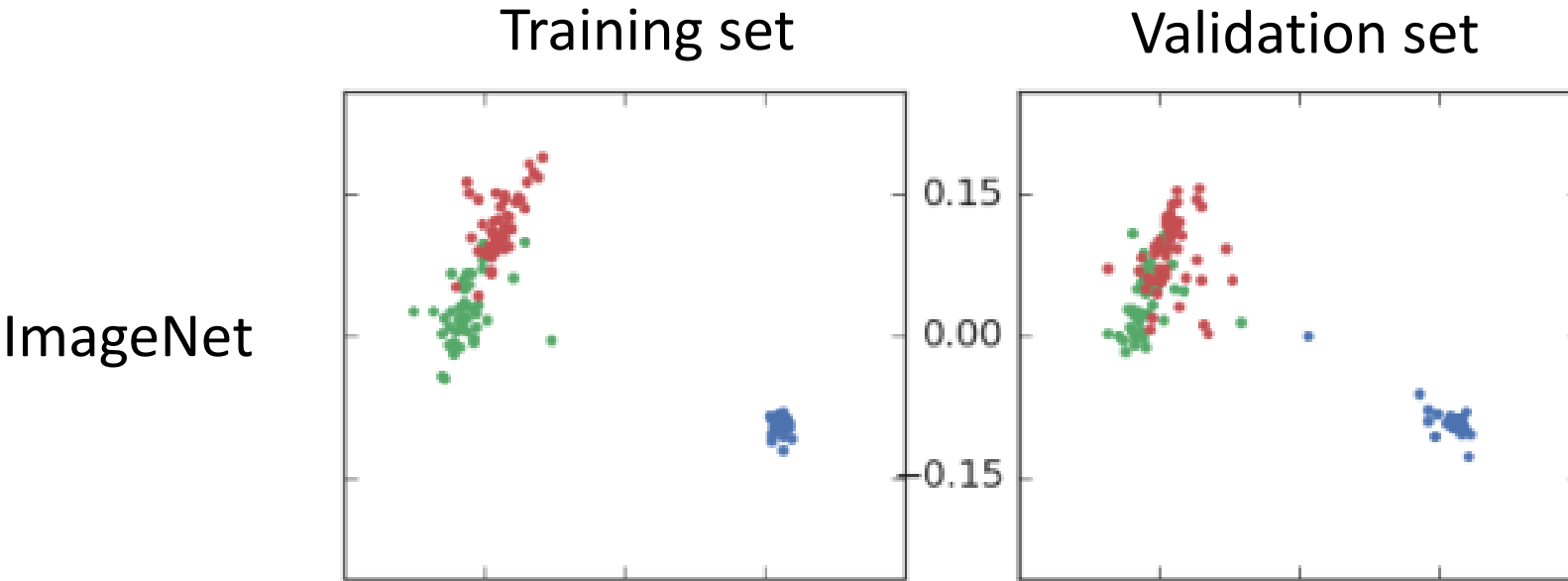
2 similar classes + 1 class

Figure 1: Visualization of penultimate layer's activations of: AlexNet/CIFAR-10 (first row), CIFAR-100/ResNet-56 (second row) and ImageNet/Inception-v4 with three semantically different classes (third row) and two semantically similar classes plus a third one (fourth row).

# Penultimate Layer Representations



# Penultimate Layer Representations



2 similar classes + 1 class

# Penultimate Layer Representations

DATA SET	ARCHITECTURE	ACCURACY ( $\alpha = 0.0$ )	ACCURACY ( $\alpha = 0.1$ )
CIFAR-10	ALEXNET	$86.8 \pm 0.2$	$86.7 \pm 0.3$
CIFAR-100	RESNET-56	$72.1 \pm 0.3$	$72.7 \pm 0.3$
IMAGENET	INCEPTION-V4	80.9	80.9

# Implicit Model Calibration

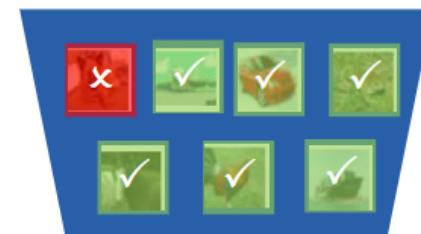
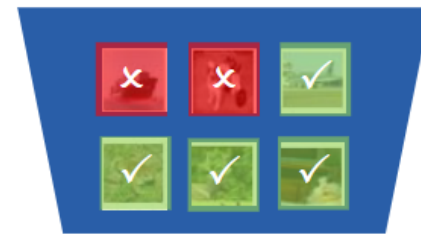
$0 \leq \hat{p}(\mathbf{x}) < 0.2$

$0.2 \leq \hat{p}(\mathbf{x}) < 0.4$

$0.4 \leq \hat{p}(\mathbf{x}) < 0.6$

$0.6 \leq \hat{p}(\mathbf{x}) < 0.8$

$0.8 \leq \hat{p}(\mathbf{x}) < 1$



Avg conf:

0.55

0.74

0.86

— Accuracy:

0.50

0.67

0.71

Gap:

2 x 0.05

6 x 0.17

7 x 0.15

15

Expected Calibrated Error (ECE) = 0.11

# Implicit Model Calibration

DATA SET	ARCHITECTURE	BASELINE	TEMP. SCALING	LABEL SMOOTHING
		ECE ( $T=1.0, \alpha = 0.0$ )	ECE / T ( $\alpha = 0.0$ )	ECE / $\alpha$ ( $T=1.0$ )
CIFAR-100	RESNET-56	0.150	0.021 / 1.9	0.024 / 0.05
IMAGENET	INCEPTION-V4	0.071	0.022 / 1.4	0.035 / 0.1
EN-DE	TRANSFORMER	0.056	0.018 / 1.13	0.019 / 0.1

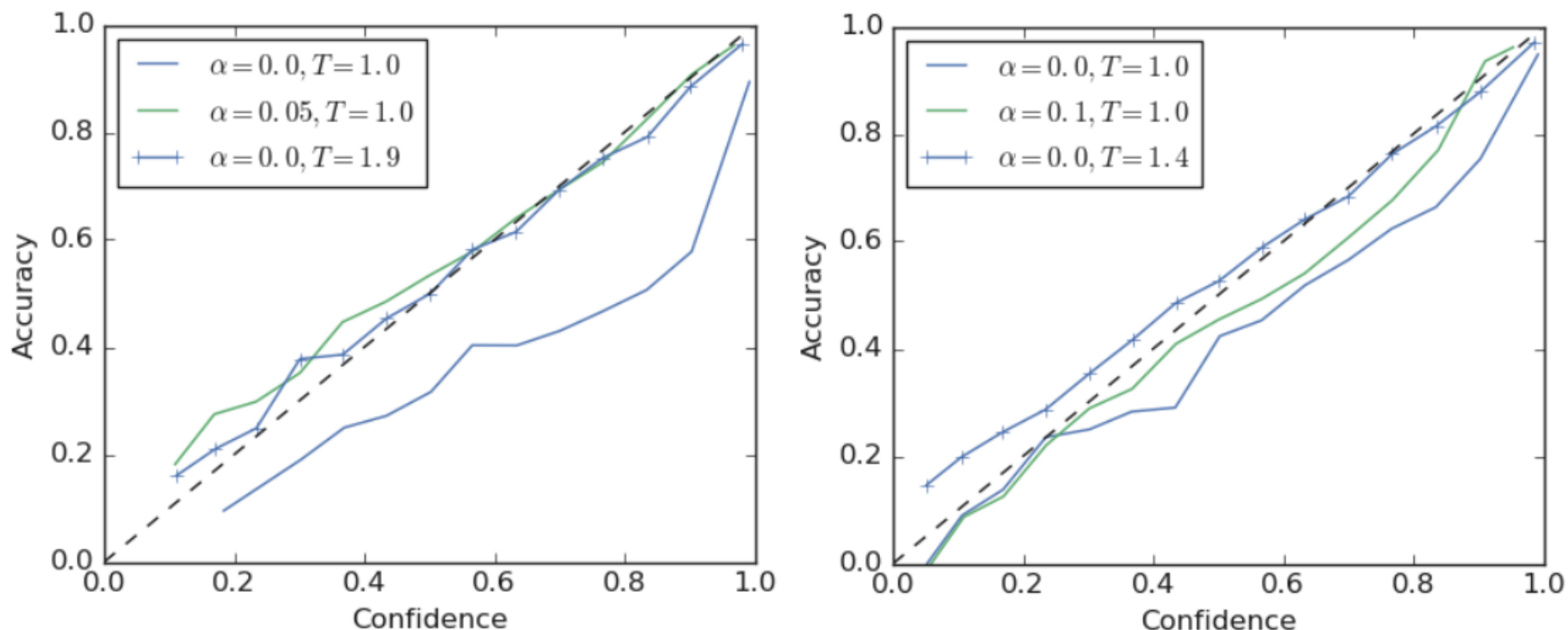


Figure 2: Reliability diagram of ResNet-56/CIFAR-100 (left) and Inception-v4/ImageNet (right).

The setting of machine translation experiment is interesting for two reasons:

1. Label Smoothing  $\rightarrow$  BLEU $\uparrow$  but Perplexity $\uparrow$

2. CV: only care about accuracy

NLP: the network's soft outputs are inputs to a second algorithm (beam-search) which is affected by calibration

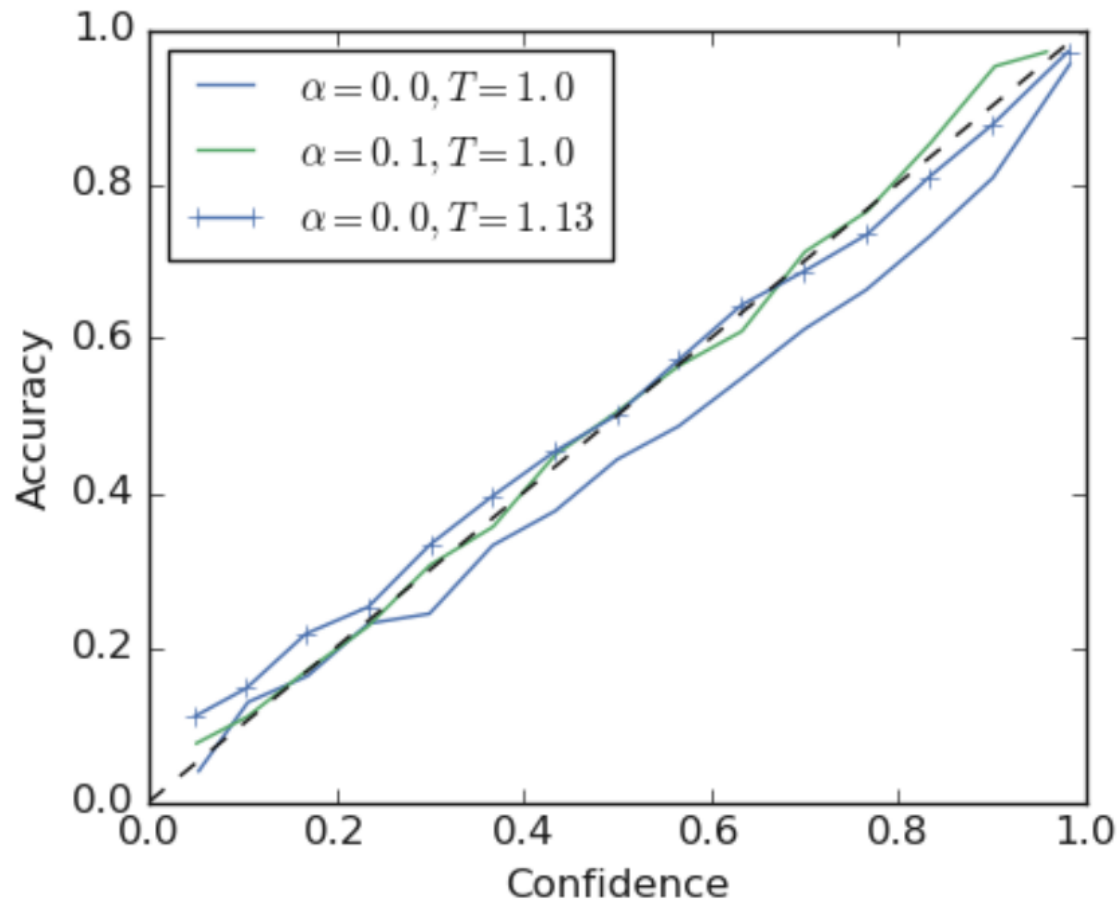


Figure 3: Reliability diagram of Transformer trained on EN-DE dataset.

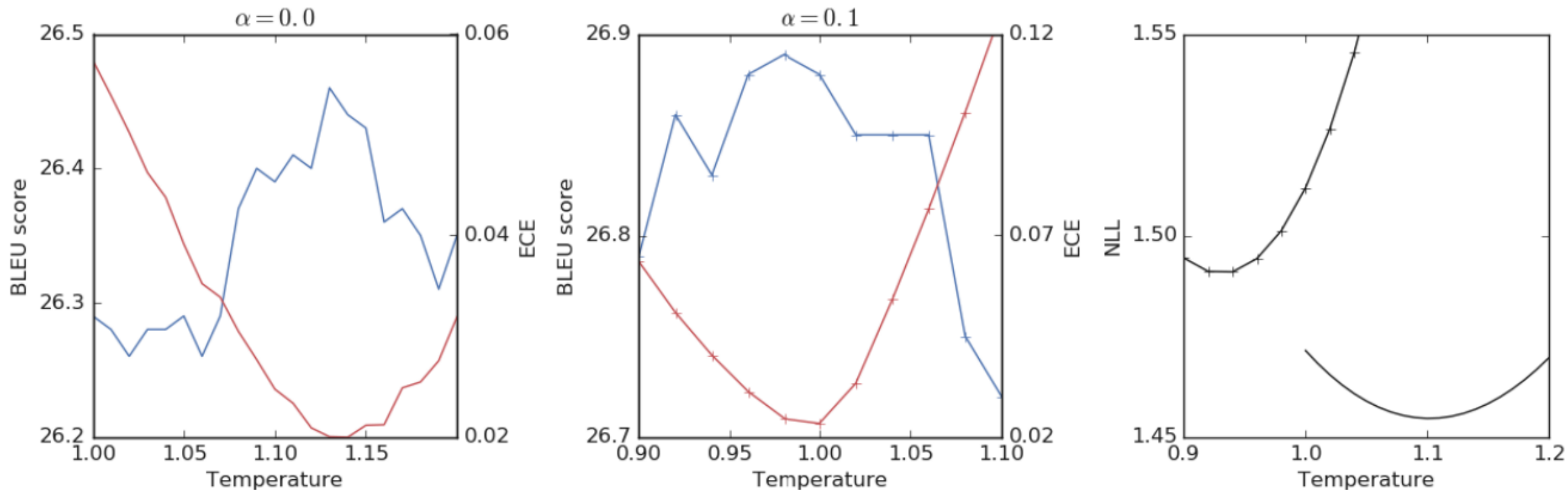


Figure 4: Effect of calibration of Transformer upon BLEU score (blue lines) and NLL (red lines). Curves without markers reflect networks trained without label smoothing while curves with markers represent networks with label smoothing.

## Experiments on MNIST

T/S	Settings	Test error
Teacher	non-convolutional + hard targets + <b>dropout</b>	0.67%
Student	/	<b>0.74%</b>
Teacher	non-convolutional + hard targets + <b>label smoothing</b>	<b>0.59%</b>
Student	/	0.91%

$$(1 - \beta)H(\mathbf{y}, \mathbf{p}) + \beta H(\mathbf{p}^t(T), \mathbf{p}(T))$$

Focus on four results:

1. the teacher's accuracy as a function of the label smoothing factor
2. the student's baseline accuracy as a function of the label smoothing factor without distillation
3. the student's accuracy after distillation with temperature scaling to control the smoothness of the teacher's provided targets (teacher trained with hard targets)
4. the student's accuracy after distillation with fixed temperature ( $T = 1.0$  and teacher trained with label smoothing to control the smoothness of the teacher's provided targets)

$$\gamma = \mathbb{E} \left[ \sum_{k=1}^K (1 - y_k) p_k^t(T) \underline{K/(K-1)} \right]$$

the mass allocated by the teacher to incorrect examples over training set

$$y_k^{LS} = y_k(1 - \alpha) + \alpha/K$$

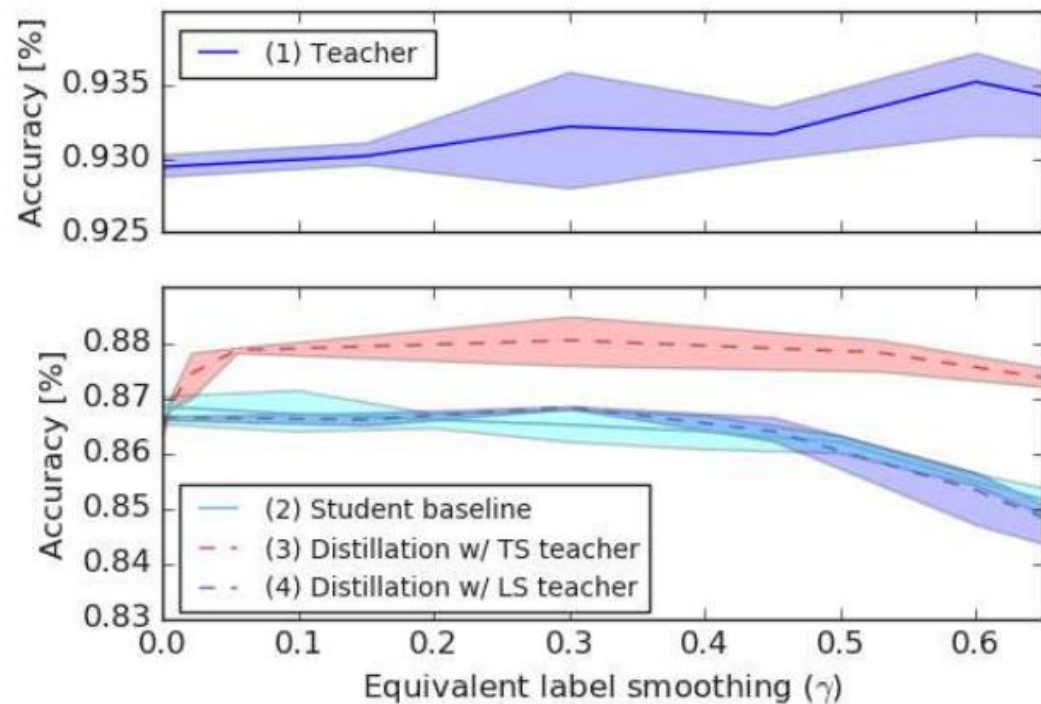


Figure 5: Performance of distillation from ResNet-56 to AlexNet on CIFAR-10.

$$y = f(d(\mathbf{z}_x))$$

$$I(X; Y) = E_{X, Y} [\log(p(y|x)) - \log(\sum_x p(y|x))]$$

$$\hat{I}(X; Y) = \sum_{x=1}^N \left[ - (f(d(\mathbf{z}_x)) - \mu_x)^2 / (2\sigma^2) - \log\left(\sum_{x=1}^N e^{- (f(d(\mathbf{z}_x)) - \mu_x)^2 / (2\sigma^2)}\right) \right]$$

where

$$\mu_x = \sum_{l=1}^L f(d(\mathbf{z}_x)) / L, \sigma^2 = \sum_{x=1}^N (f(d(\mathbf{z}_x)) - \mu_x)^2 / N$$

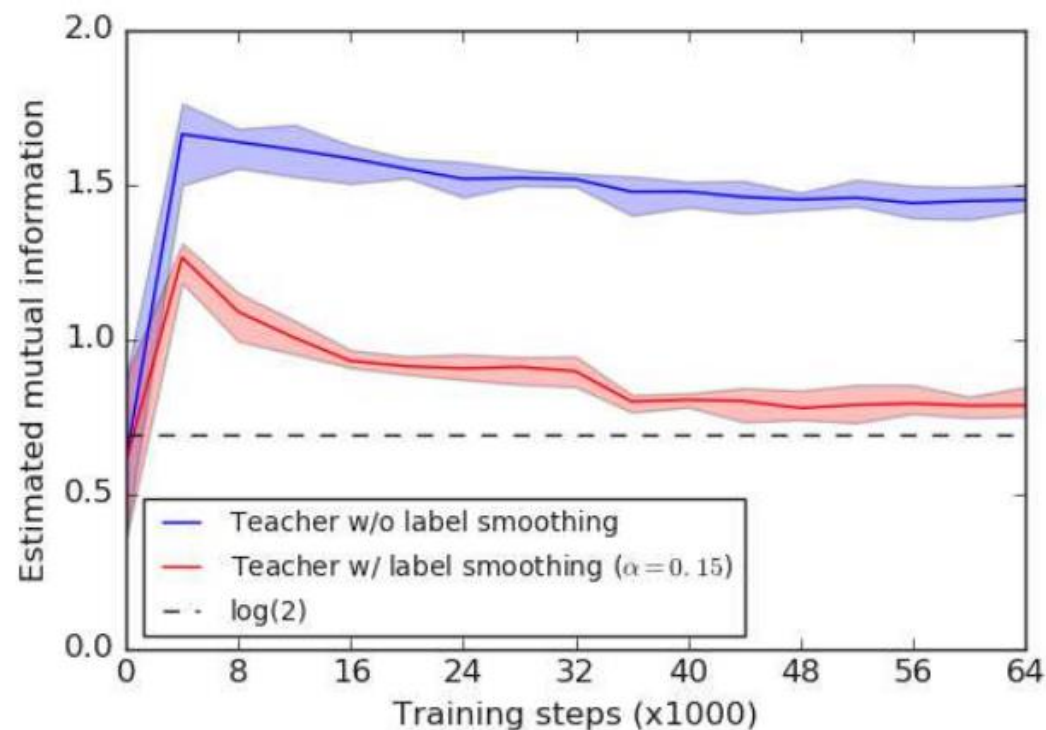


Figure 6: Estimated mutual information evolution during teacher training.



# IS LABEL SMOOTHING TRULY INCOMPATIBLE WITH KNOWLEDGE DISTILLATION: AN EMPIRICAL STUDY

**Zhiqiang Shen**  
CMU

**Zechun Liu**  
CMU & HKUST

**Dejia Xu**  
Peking University

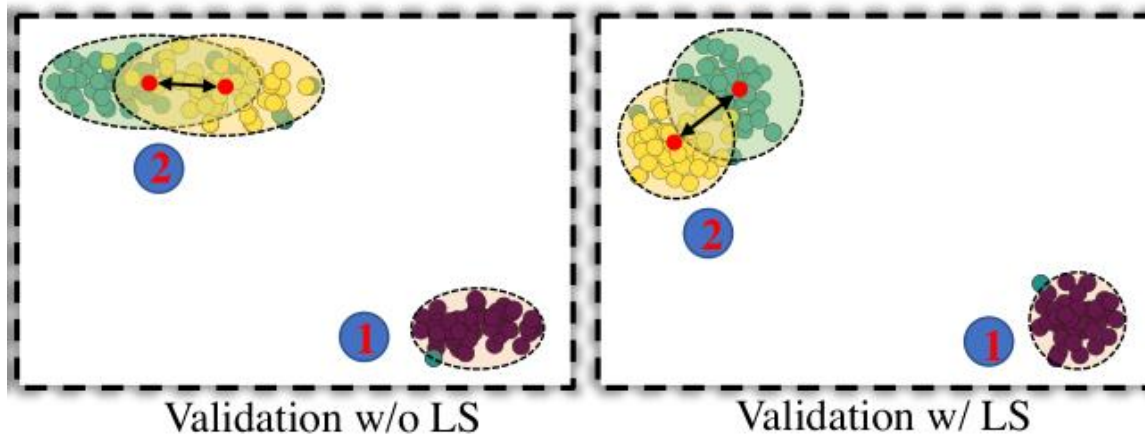
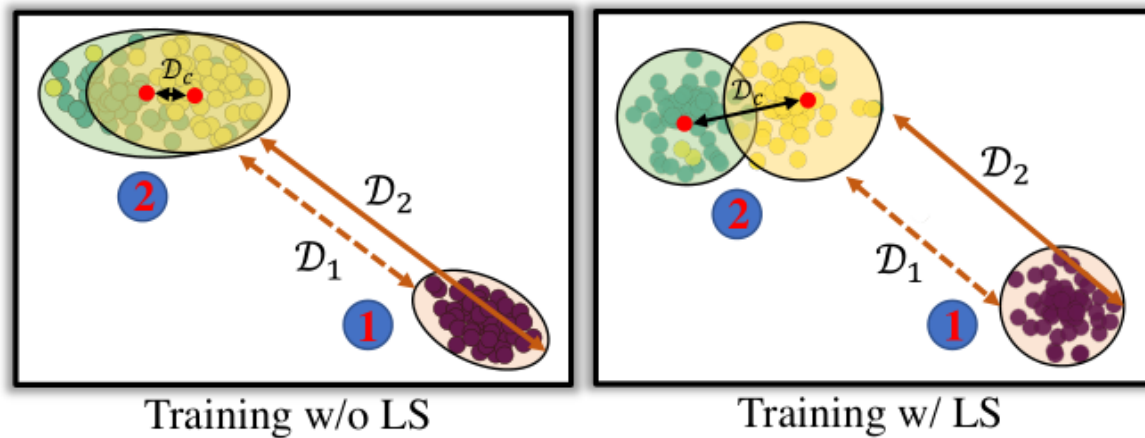
**Zitian Chen**  
UMass Amherst

**Kwang-Ting Cheng**  
HKUST

**Marios Savvides**  
CMU

ICLR 2021

# The Most Important Picture in this Article



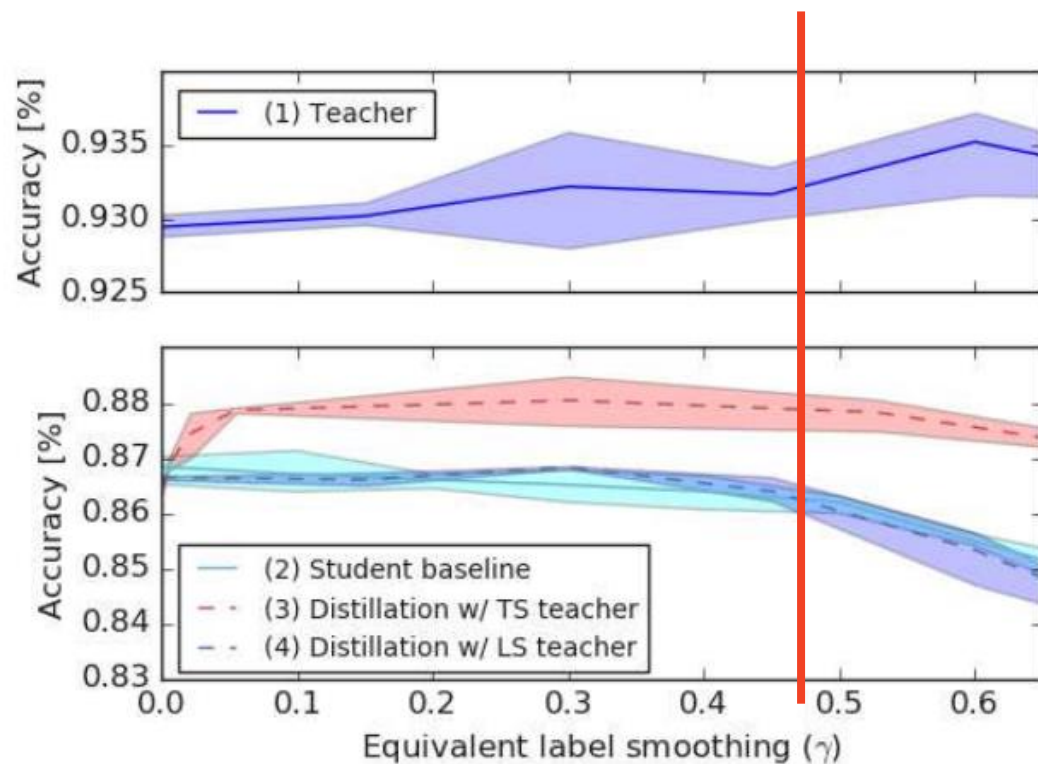
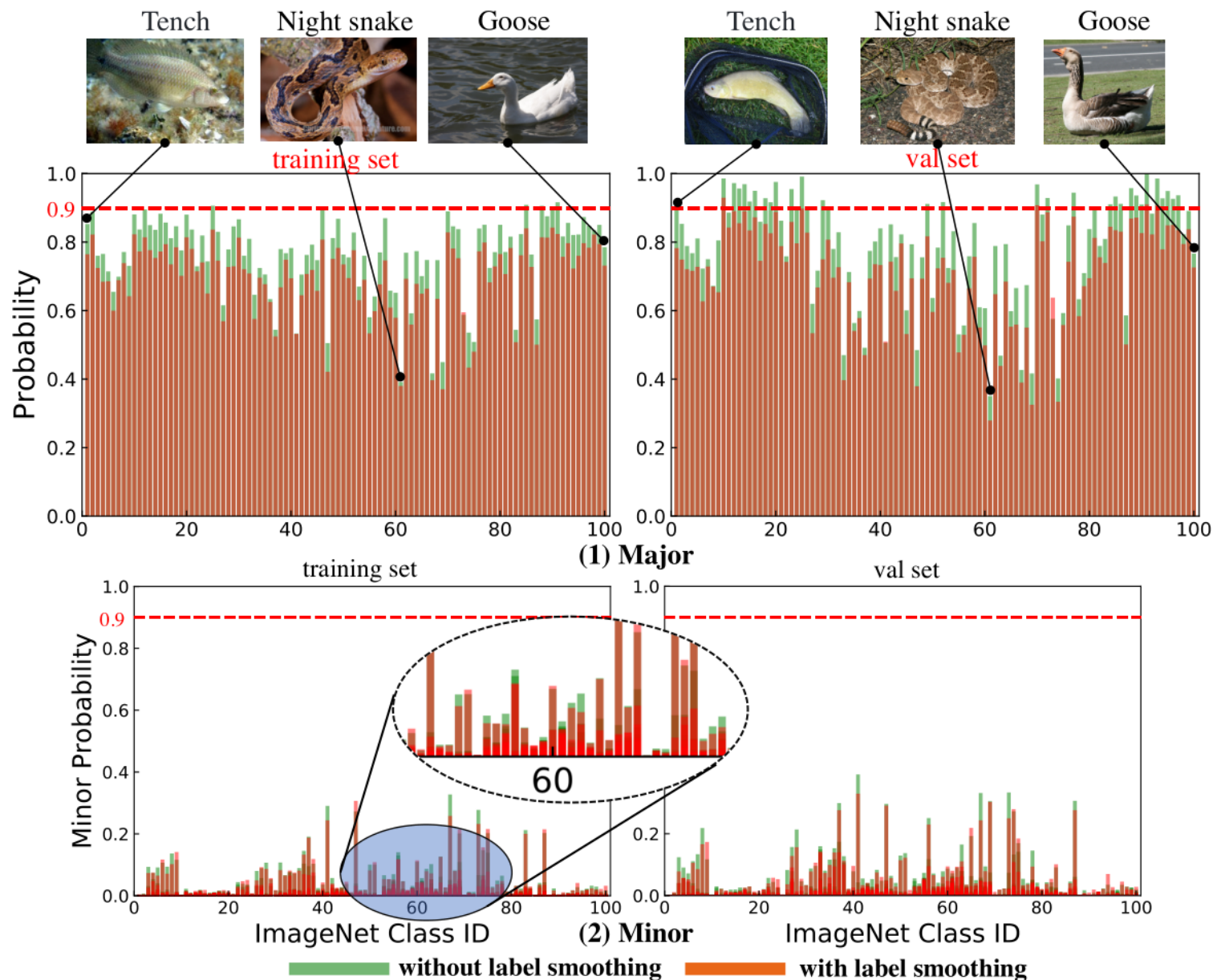


Figure 5: Performance of distillation from ResNet-56 to AlexNet on CIFAR-10.

$$\gamma = \mathbb{E} \left[ \sum_{k=1}^K (1 - y_k) p_k^t(T) K / (K - 1) \right]$$

# The “Erase Information” Effect by Label Smoothing



# The “Erase Information” Effect by Label Smoothing

$$\mathcal{S}_{\text{Stability}} = 1 - \frac{1}{K} \sum_{c=1}^K \left( \frac{1}{n_c} \sum_{i=1}^{n_c} \|\mathbf{p}_{\{i,c\}}^{\mathcal{T}_w} - \bar{\mathbf{p}}_{\{i,c\}}^{\mathcal{T}_w}\|^2 \right)$$

Some interesting facts(the authors do not give specific experimental data):

- 1) It can measure the degree of erased information quantitatively and further help discover more interesting phenomena, e.g., they observe that data augmentation method like CutMix (Yun et al., 2019) together with longer training erases the relative information on logits dramatically and can further be reinforced by label smoothing.
- 2) We found that the proposed metric is highly aligned with model accuracy, thus such metric can be used as a complement for accuracy to evaluate the quality of teacher’s supervision for knowledge distillation.

<del>Networks</del> Networks	Acc. (%) w/o LS	( $1-\mathcal{S}_{\text{Stability}}$ ) w/o LS	Acc. (%) w/ LS	( $1-\mathcal{S}_{\text{Stability}}$ ) w/ LS
ResNet-18 (He et al., 2016)	69.758/89.078	0.3359	<b>69.774/89.122</b>	<b>0.3358 (-0.0001)</b>
ResNet-50 (He et al., 2016)	75.888/92.642	0.3217	<b>76.130/92.936</b>	<b>0.3106 (-0.0111)</b>
ResNet-101 (He et al., 2016)	77.374/93.546	0.3185	<b>77.726/93.830</b>	<b>0.3070 (-0.0115)</b>
MobileNet v2 (Sandler et al., 2018)	71.878/90.286	0.3341	–	–
DenseNet121 (Huang et al., 2017)	74.434/91.972	0.3243	–	–
ResNeXt50 32×4d (Xie et al., 2017)	77.618/93.698	0.3229	<b>77.774/93.642</b>	<b>0.3182 (-0.0047)</b>
Wide ResNet50 (Zagoruyko & Komodakis, 2016)	<b>78.468/94.086</b>	0.3201	77.808/93.682	<b>0.3155 (-0.0046)</b>
ResNeXt101 32×8d (Xie et al., 2017)	79.312/94.526	0.3177	<b>79.698/94.768</b>	<b>0.3116 (-0.0061)</b>
ResNet50+Long	76.526/93.070	0.3222	<b>77.106/93.340</b>	<b>0.3090 (-0.0132)</b>
ResNet50+Long+CutMix (Yun et al., 2019)	76.874/93.500	0.2999	<b>77.274/93.304</b>	<b>0.2890 (-0.0109)</b>

<i>ImageNet-1K (Standard):</i>				
Teacher	w/ LS	Acc. (Top1/Top5)	Student	Acc. (Top1/Top5)
ResNet-50	✗	76.056 ± 0.119/92.791 ± 0.106	ResNet-18	71.425 ± 0.038/90.185 ± 0.075
			ResNet-50	76.325 ± 0.068/92.984 ± 0.043
	✓	<b>76.128 ± 0.069/92.977 ± 0.030</b>	ResNet-18	<b>71.816 ± 0.017/90.466 ± 0.074</b>
			ResNet-50	<b>77.052 ± 0.030/93.376 ± 0.015</b>
ResNet-101	✗	77.374/93.546	ResNet-50	77.428/93.712
			ResNet-101	78.270/94.152
	✓	<b>77.836/93.662</b>	ResNet-50	<b>77.624/93.862</b>
			ResNet-101	<b>78.476/94.008</b>
<i>CUB200-2011 (Fine-grained):</i>				
Teacher	w/ LS	Acc. (Top1/Top5)	Student	Acc. (Top1/Top5)
ResNet-50	✗	79.931 ± 0.037/94.370 ± 0.064	ResNet-18	77.116 ± 0.086/93.241 ± 0.108
			ResNet-50	80.910 ± 0.033/94.738 ± 0.114
	✓	<b>81.497 ± 0.035/95.043 ± 0.112</b>	ResNet-18	<b>78.382 ± 0.099/93.621 ± 0.120</b>
			ResNet-50	<b>82.355 ± 0.050/95.440 ± 0.075</b>
ResNet-101	✗	80.380/94.491	ResNet-50	81.261/94.905
			ResNet-101	81.572/95.371
	✓	<b>82.332/94.970</b>	ResNet-50	<b>82.263/95.320</b>
			ResNet-101	<b>82.522/95.199</b>
<i>iMaterialist-2019_P (Noisy):</i>				
Teacher	w/ LS	Acc. (Top1/Top3)	Student	Acc. (Top1/Top3)
ResNet-50	✗	66.241/91.015	ResNet-18	65.250/90.243
			ResNet-50	67.420/92.155
	✓	<b>66.825/91.669</b>	ResNet-18	<b>65.359/90.530</b>
			ResNet-50	<b>67.528/92.551</b>
ResNet-101	✗	66.726/91.263	ResNet-50	67.905/92.481
			ResNet-101	68.281/92.580
	✓	<b>67.370/91.877</b>	ResNet-50	<b>67.925/92.789</b>
			ResNet-101	<b>68.618/92.907</b>

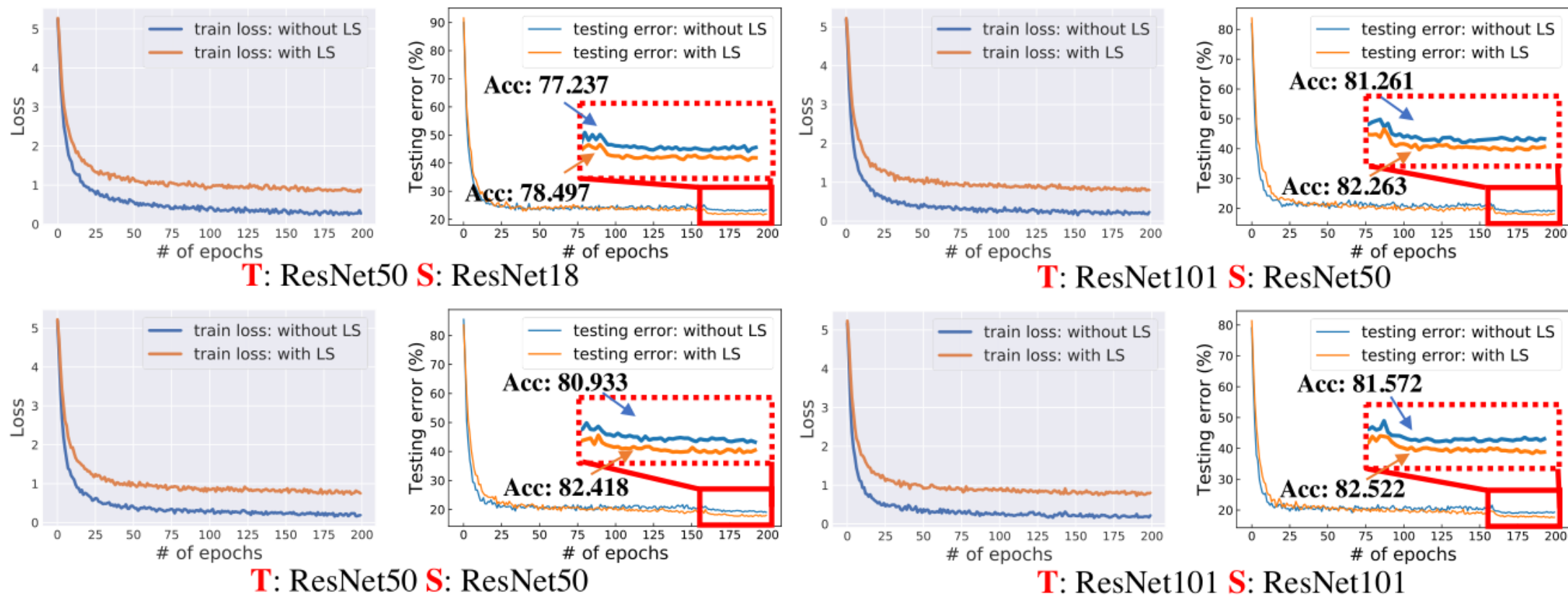
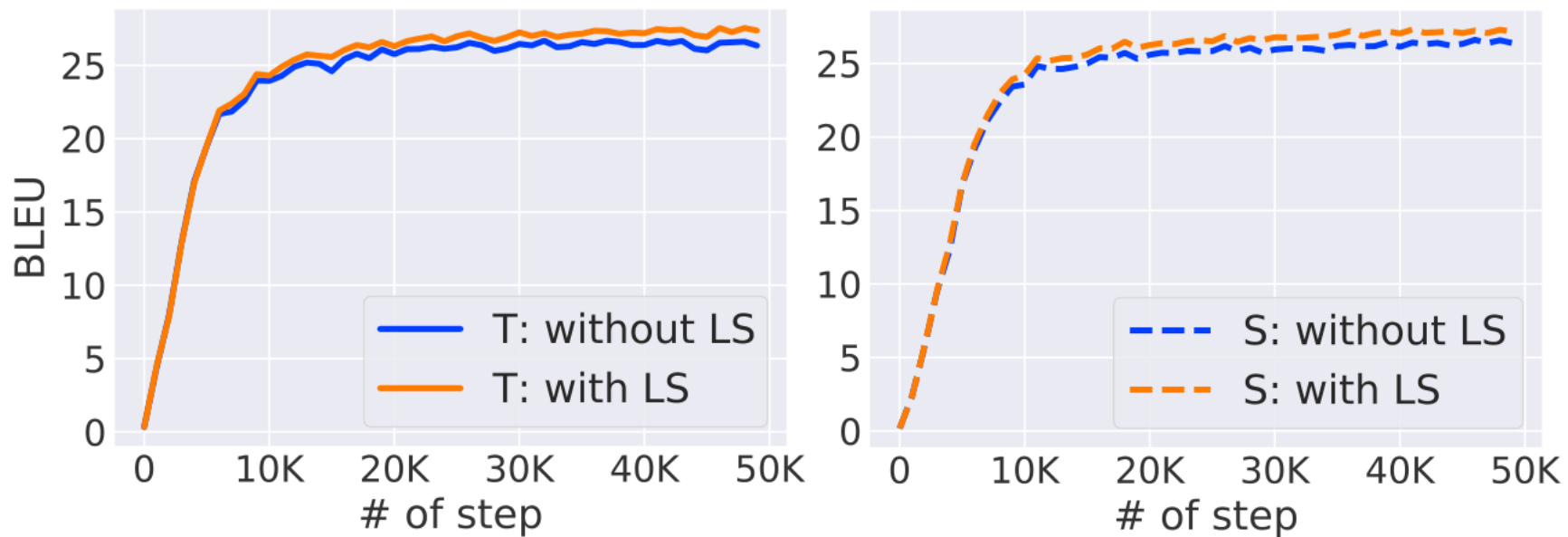


Figure 5: The training and testing curves of knowledge distillation on CUB200-2011 when teachers are trained w/ and w/o label smoothing. The specific teacher and student architectures are given below each subfigure, therein, **T** indicates the teacher architecture and **S** indicates the student.



# What is a Better Teacher in Knowledge Distillation?

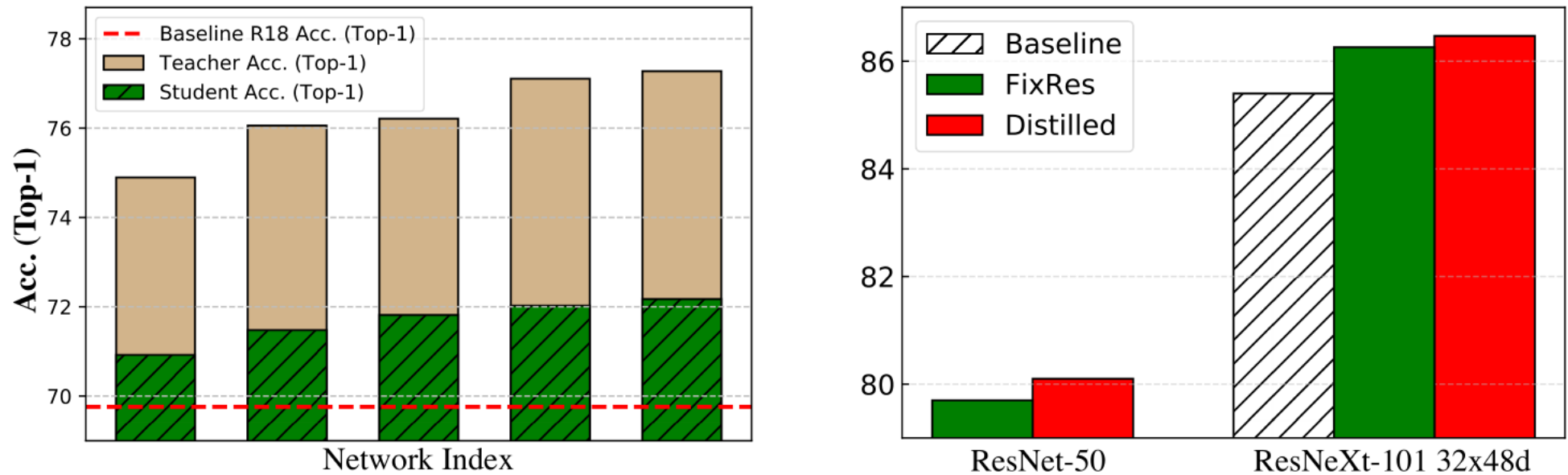


Figure 8: Left is the accuracy relationship between teachers and students, wherein, all teachers are trained with label smoothing. Right is the accuracy of knowledge distillation by using strong teacher to fine-tune the student, FixRes (Touvron et al., 2019) is adopted in both teacher and student networks.

# What is a Better Teacher in Knowledge Distillation?

<b>Teacher (same arch)</b>	<b>Acc. (Top-1)</b>	<b>Student</b>	<b>Acc. (Top-1)</b>
R50	76.056	R18	71.478
R50+LS	76.212	R18	71.816
R50+LS+Long	77.106	R18	72.024
R50+LS+Long+CutMix	77.274	R18	72.172

<b>Teacher (different archs)</b>	<b>Acc. (Top-1)</b>	<b>Student</b>	<b>Acc. (Top-1)</b>
MobileNet V2	71.878	R18	70.054
DenseNet-121	74.894	R18	70.922
Wide ResNet-50-2	77.808	R18	72.232
ResNeXt-101-32x8d	79.698	R18	72.412

# What is a Better Teacher in Knowledge Distillation?

ImageNet		
Ratio (hard label – soft label)	w/ LS	Acc. (Top1/Top5)
0.3 – <del>0.5</del> 0.7	<b>x</b>	71.592/90.386
	<b>✓</b>	<b>71.752/90.412</b>
0.5 – 0.5	<b>x</b>	71.484/90.218
	<b>✓</b>	<b>71.748/90.454</b>
0.7 – 0.3	<b>x</b>	71.164/90.196
	<b>✓</b>	<b>71.314/90.200</b>

# What Circumstances will make LS less Effective?

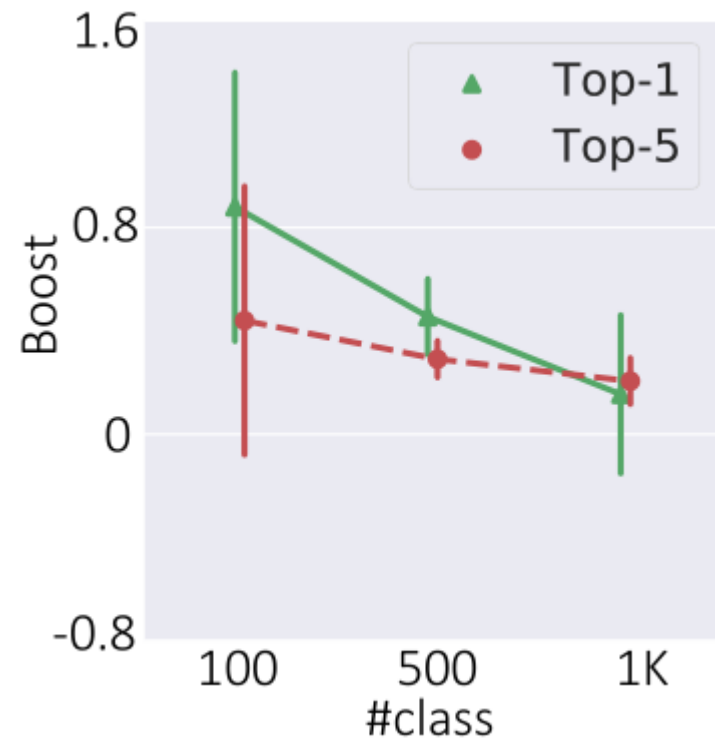
## 1. Long-Tailed Distribution

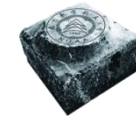
Teacher	w/ LS	ImageNet-LT	Place-LT	iNaturalist 2019
		Acc. (Top1/Top5)	Acc. (Top1/Top3)	Acc. (Top1/Top3)
ResNet-18	<b>✗</b>	<b>39.975/64.645</b>	<b>26.479/47.233</b>	<b>67.195/83.465</b>
	<b>✓</b>	39.115/63.655	25.877/46.260	66.700/83.432
ResNet-34	<b>✗</b>	<b>41.150/66.205</b>	<b>27.329/48.753</b>	<b>70.165/86.304</b>
	<b>✓</b>	40.965/65.850	26.863/48.110	69.406/86.073
ResNet-50	<b>✗</b>	<b>40.985/66.030</b>	27.384/ <b>48.740</b>	<b>73.729/88.845</b>
	<b>✓</b>	39.965/65.195	<b>27.562</b> /47.945	72.904/87.954
ResNet-101	<b>✗</b>	—/—	<b>28.096/50.164</b>	<b>74.389/88.416</b>
	<b>✓</b>	—/—	27.466/48.781	<b>73.597/88.779</b>

# What Circumstances will make LS less Effective?

## 2. More Classes

		ImageNet-100	ImageNet-500	ImageNet-1K
Teacher	w/ LS	Acc. (Top1/Top5)	Acc. (Top1/Top5)	Acc. (Top1/Top5)
ResNet-18	✗	82.380/95.520	73.521/91.642	<b>69.758/89.076</b>
	✓	<b>82.740/95.440</b>	<b>74.123/92.004</b>	<b>69.606/89.372</b>
ResNet-101	✗	82.000/94.340	81.712/95.080	77.374/93.546
	✓	<b>83.400/95.300</b>	<b>82.020/95.300</b>	<b>77.836/93.662</b>
Average (↑)		↑ <b>0.880/0.440</b>	↑0.455/0.291	↑0.155/0.206





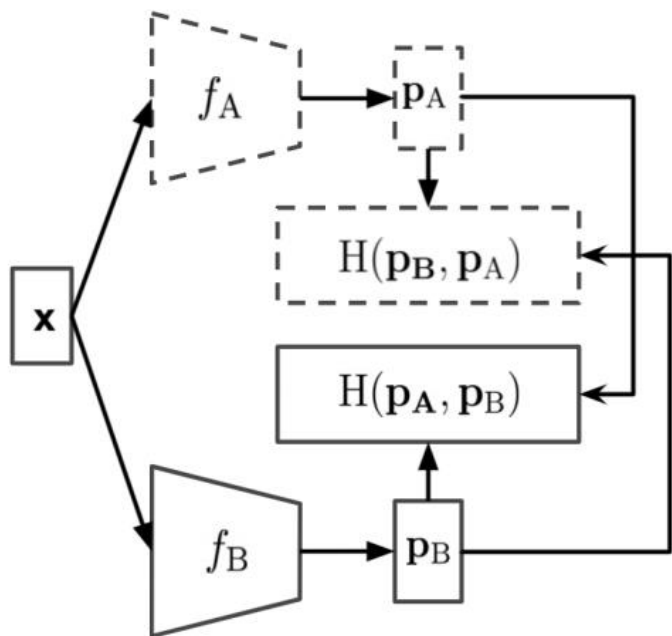
Another paper I want to talk about

---

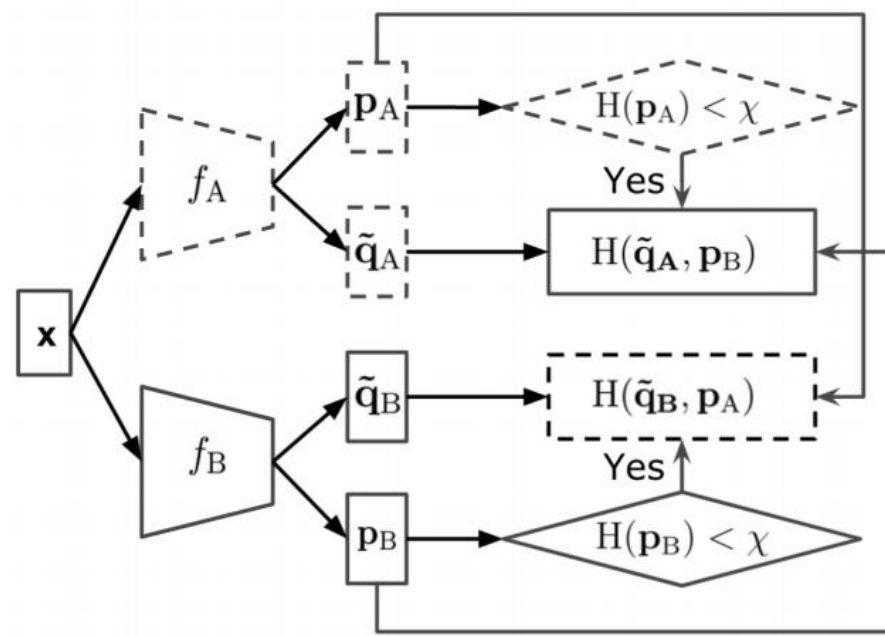
# Not All Knowledge Is Created Equal

---

**Ziyun Li<sup>1</sup>, Xinshao Wang<sup>2,6,\*</sup>, Haojin Yang<sup>1</sup>, Di Hu<sup>3</sup>,  
Neil M. Robertson<sup>4,6</sup>, David A. Clifton<sup>2,5,†</sup>, Christoph Meinel<sup>1</sup>**



(a) Conventional MKD.



(b) CMD.

$$L_A = L_{A_{\text{SelfKD}}} + L_{B2A} = \begin{cases} H(\tilde{\mathbf{q}}_A, \mathbf{p}_A) + H(\tilde{\mathbf{q}}_B, \mathbf{p}_A) = 2E_{\frac{\tilde{\mathbf{q}}_A + \tilde{\mathbf{q}}_B}{2}}(-\log \mathbf{p}_A), & H(\mathbf{p}_B) < \chi, \\ H(\tilde{\mathbf{q}}_A, \mathbf{p}_A) = E_{\tilde{\mathbf{q}}_A}(-\log \mathbf{p}_A), & H(\mathbf{p}_B) \geq \chi. \end{cases}$$

$$L_B = L_{B_{\text{SelfKD}}} + L_{A2B} = \begin{cases} H(\tilde{\mathbf{q}}_B, \mathbf{p}_B) + H(\tilde{\mathbf{q}}_A, \mathbf{p}_B) = 2E_{\frac{\tilde{\mathbf{q}}_A + \tilde{\mathbf{q}}_B}{2}}(-\log \mathbf{p}_B), & H(\mathbf{p}_A) < \chi, \\ H(\tilde{\mathbf{q}}_B, \mathbf{p}_B) = E_{\tilde{\mathbf{q}}_B}(-\log \mathbf{p}_B), & H(\mathbf{p}_A) \geq \chi. \end{cases}$$



南京航空航天大学  
Nanjing University of Aeronautics and Astronautics

ParNeC

模式识别与神经计算研究组  
PAttern Recognition and NEural Computing

Thanks for Listening

