# Revisiting Knowledge Distillation via Label Smoothing Regularization — Supplementary Material

## 1. Implementation details for exploratory experiments

In this section, we provide implementation details and experiment settings for the exploratory experiments. We conduct experiments based on the standard implementation of knowledge distillation [1]. The loss function for standard knowledge distillation is:

$$\mathcal{L}_{KD} = (1-\alpha)H(q,p) + \alpha D_{KL}(p_\tau^t, p_\tau) \qquad (S.1)$$

where $q$ is the distribution of ground truth, $p$ is the output distribution of student model, $H(,)$ is cross-entropy loss function and $D_{KL}$ is KL divergence, and $p_\tau^t$ is the output distribution of teacher model soften by temperature $\tau$. The temperature $\tau$ and weight $\alpha$ are hyper-parameters. The temperature $\tau$ and weight $\alpha$ for Reverse KD and Normal KD are given in Tab. 1; $\tau$ and $\alpha$ for De-KD are given in Tab. 2. For fair comparisons, the hyper-parameters of Normal KD, Re-KD and De-KD are grid searched from 70 training epochs (200 epochs in total), and 8 GPUs are used to search these hyper-paramwters (4 NVIDIA V, 4 NVIDIA X). All experiments are conducted with the same experiment settings.

**CIFAR10 and CIFAR100** For exploratory experiments Re-KD and De-KD on CIFAR10 and CIFAR100, we train for 200 epochs, with batch size 128. For the Plain CNN, the The initial learning rate is 0.01 and then be divided by 5 at the 60'th, 120'th, 160'th epoch. For other used models (MobileNetV2, ShuffleNetV2, ResNet, ResNeXt, DenseNet), the initial learning rate is 0.1 and then be divided by 5 at the 60'th, 120'th, 160'th epoch. We use Adam optimizer for the Plain CNN and SGD optimizer with momentum 0.9 for other models, and the weight decay is set to be 5e-4. For hyper-parameters, $\tau$ (temperature) and $\alpha$, we use grid search to find the best value.

The Plain CNN used in exploratory experiments is a 5-layer neural network with 3 convolutional layers and 2 fully-connected layers. On CIFAR10, the architecture of the Plain CNN is: $conv1$(3x3, 32, channels) $\rightarrow batchnorm \rightarrow conv2$(3x3, 64, channels) $\rightarrow$ $batchnorm$ $\rightarrow$ $conv3$(3x3, 128, channels) $\rightarrow batchnorm \rightarrow fc(128) \rightarrow fc(10)$.

**Tiny-ImageNet** For exploratory experiments on Tiny-ImageNet, all models are trained for 200 epochs, with batch size $bn = 128$ for MobileNetV2, ShuffleNetV2, ResNet18 and $bn = 64$ for ResNet50, DenseNet121. The initial learning rate is $\eta = 0.1 * \frac{bn}{128}$ and then be divided by 10 at the 60'th, 120'th, 160'th epoch. We use SGD optimizer with momentum of 0.9, and the weight decay is set to be 5e-4.

**Model complexity** We provide the model complexity (size and FLOPs) of all models we used in this work in Tab. 3, which is the reference to choose teacher and student model. The model size is measured by the total number of learnable parameters within each model. The FLOPs of model is tested with image size of $3 \times 224 \times 224$.

## 2. Comparison between $\text{Tf}_{reg}$ with LSR

We first recall some important equantions for LSR and $\text{Tf}_{reg}$. In LSR, it minimizes the cross-entropy between modified label distribution $q'(k)$ and the network output $p(k)$, where $q'(k)$ is the smoothed label distribution formulated as

$$q'(k) = (1-\alpha)q(k) + \alpha u(k), \qquad (S.2)$$

which is a mixture of $q(k)$ and a fixed distribution $u(k)$, with weight $\alpha$. Usually, the $u(k)$ is uniform distribution as $u(k) = 1/K$. The loss function of LSR can be written as

$$\mathcal{L}_{LS} = (1-\alpha)H(q,p) + \alpha D_{KL}(u,p). \qquad (S.3)$$

For $\text{Tf}_{reg}$, we define a virtual teacher model which will output distribution for classes as the following:

$$p^d(k) = \begin{cases} a & \text{if } k = c, \\ (1-a)/(K-1) & \text{if } k \neq c, \end{cases} \qquad (S.4)$$

where $K$ is the total number of classes, $c$ is the correct label and $a$ is the correct probability for the correct class. And the loss function is

$$L_{reg} = (1-\alpha)H(q,p) + \alpha D_{KL}(p_\tau^d, p_\tau), \qquad (S.5)$$

where $\tau$ is the temperature to soften the manually-designed distribution $p^d$ (as $p_\tau^d$ after softening). We set a high temperature $\tau \geq 20$ to make this virtual teacher output a soft

Table 1. Parameters for Normal KD and Re-KD experiments (Temperature $\tau$ and $\alpha$)

| Dataset | Teacher | Student | Normal KD | Re-KD |
|---|---|---|---|---|
| CIFAR-10 | ResNet18 | Plain CNN | $\tau=20, \alpha=0.90$ | $\tau=20, \alpha=0.01$ |
| | | MobileNetV2 | $\tau=20, \alpha=0.90$ | $\tau=20, \alpha=0.05$ |
| | MobileNetV2 | Plain CNN | $\tau=20, \alpha=0.40$ | $\tau=20, \alpha=0.10$ |
| | ResNeXt29 | ResNet18 | $\tau=6,\quad \alpha=0.95$ | $\tau=20, \alpha=0.10$ |
| CIFAR100 | ResNet18 | MobileNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | ResNet50: | MobileNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | Densenet121 | MobileNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ | $\tau=20, \alpha=0.60$ |
| | ResNeXt29 | MobileNetV2 | $\tau=20, \alpha=0.60$ | $\tau=20, \alpha=0.60$ |
| | | ResNet18 | $\tau=20, \alpha=0.60$ | $\tau=20, \alpha=0.60$ |
| T-ImageNet | ResNet18 | MobileNetV2 | $\tau=20, \alpha=0.10$ | $\tau=20, \alpha=0.60$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.10$ | $\tau=20, \alpha=0.60$ |
| | ResNet50 | MobileNetV2 | $\tau=20, \alpha=0.10$ | $\tau=20, \alpha=0.10$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.10$ | $\tau=20, \alpha=0.50$ |
| | | ResNet18 | $\tau=20, \alpha=0.50$ | $\tau=20, \alpha=0.10$ |

Table 2. Parameters for De-KD (Temperature $\tau$ and $\alpha$)

| Dataset | Poorly-trained Teacher | Student | $\tau$ and $\alpha$ |
|---|---|---|---|
| CIFAR100 | ResNet18: 15.48% | MobileNetV2 | $\tau=20, \alpha=0.95$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ |
| | ResNet50: 45.82% | MobileNetV2 | $\tau=20, \alpha=0.95$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ |
| | | ResNet18 | $\tau=20, \alpha=0.60$ |
| | ResNeXt29: 51.94% | MobileNetV2 | $\tau=20, \alpha=0.95$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.95$ |
| | | ResNet18 | $\tau=20, \alpha=0.60$ |
| Tiny-ImageNet | ResNet18: 9.41% | MobileNetV2 | $\tau=20, \alpha=0.10$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.10$ |
| | ResNet50: 31.01% | MobileNetV2 | $\tau=20, \alpha=0.10$ |
| | | ShuffleNetV2 | $\tau=20, \alpha=0.10$ |

Table 3. Model complexity of all used models in this work.

| Model | MobileNetV2 | ShuffleNetV2 | ResNet18 | ResNet50 |
|---|---|---|---|---|
| # param | $3.50 \times 10^6$ | $2.28 \times 10^6$ | $11.69 \times 10^6$ | $25.56 \times 10^6$ |
| FLOPs | $0.32 \times 10^9$ | $0.15 \times 10^9$ | $1.82 \times 10^9$ | $4.12 \times 10^9$ |

| Model | GoogLeNet | DenseNet121 | ResNeXt29 (8x64d) | ResNeXt101 (32×8d) |
|---|---|---|---|---|
| # param | $13.0 \times 10^6$ | $7.98 \times 10^6$ | $34.52 \times 10^6$ | $88.79 \times 10^6$ |
| FLOPs | $1.51 \times 10^9$ | $2.88 \times 10^9$ | $4.40 \times 10^9$ | $16.51 \times 10^9$ |

probability, in which way it gains the smoothing property as LSR.

We compare $\text{Tf}_{reg}$ with LSR in two aspects.

(1). Some special parameters will make $\text{Tf}_{reg}$ be identical to LSR? For example, let $a = (1 - \alpha) + \alpha/K$, then Eq S.4 will be

$$p^d(k) = \begin{cases} 1 - a + \alpha/K & \text{if } k = c, \\ \alpha/k & \text{if } k \neq c, \end{cases} \quad \text{(S.6)}$$

which seems that this $p^d(k)$ is identical to the soft distribution of original LSR. However, it is still not a uniform distribution as LSR. After we split the loss function of LSR as two part, the distribution of LSR is a uniform distribution, so it is not the same as $p^d(k)$ even we let $a = (1-\alpha)+\alpha/K$.

(2). Check if $\text{Tf}_{reg}$ is an over-parameters version of LSR. If $\text{Tf}_{reg}$ is an over-parameters version of LSR, we can tune the parameters $\alpha$ of LSR to make it reach the same/similar performance as $\text{Tf}_{reg}$. When we tune $\alpha$ in LSR, we find that in most of the cases, the LSR will obtain the optimal performance as $\alpha$ is around 0.1, but which is still lower than our $\text{Tf}_{reg}$. We search $\alpha$ between [0.01, 1] with an interval of 0.02, and the search results of LSR on CIFAR100 are given in Tab 4. It can be observed that the best results of LSR are still worse than $\text{Tf}_{reg}$.

Table 4. Comparison between LSR and Tf-KD$_{reg}$ on CIFAR100.

| Model | Baseline | Tf-KD$_{reg}$ | + LSR (optimal $\alpha$) |
|---|---|---|---|
| MobileNetV2 | 68.38 | 70.88 (**+2.50**) | 69.54 ($\alpha = 0.19$) |
| ShuffleNetV2 | 70.34 | 72.09 (**+1.75**) | 70.98 ($\alpha = 0.23$) |
| ResNet18 | 75.87 | 77.36 (**+1.49**) | 77.49 ($\alpha = 0.15$) |

On the other hand, we cannot let Eq. S.5 be the same as Eq. S.3 because the output of student model ($p_\tau$) also soften by $\tau$. Same as normal KD, the temperature $\tau$ is always very high ($\tau \gg 1$), thus the distribution of $\text{Tf}_{reg}$ obtain the smoothing property. The temperature $\tau \neq 1$, so LSR and Tf-KD$_{reg}$ will not the same when we tune the $\alpha$ in LSR. Actually, in our analyze, even the original KD has a similar format LSR, so the Tf-KD will have a similar format with LSR because our Tf-KD is proposed based on the KD and LSR.

We summary that even $\text{Tf}_{reg}$ has a similar format with LSR, but it is not identical to LSR or an over-parameters version of LSR.

## 3. Visualization of the output distribution of teacher

To better comparing the $p_\tau^t(k)$ (the output distribution of teacher model) and $u(k)$ (the uniform distribution of label smoothing), we visualize the soft targets of ResNet18 (trained on CIFAR10 with 95.12% accuracy) and compare the soft targets in different temperature with $u(k)$. As shown in Fig. S.1, we can observe that with the temperature $\tau$ increasing, the two distributions become closer. In the common experiments of knowledge distillation, we always adopt temperature as 20 [1].

## 4. Experiment settings for $\text{Tf}_{self}$ and $\text{Tf}_{reg}$

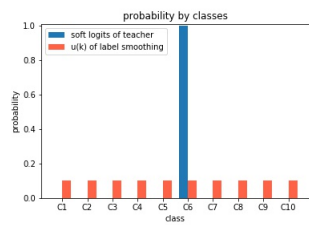For all $\text{TF}_{self}$ experiments on ImageNet, we set temperature $\tau$=20, and weight $\alpha$=0.10. The hyper-parameters on CIFAR100 and Tiny-ImageNet are given in Tab. 5.

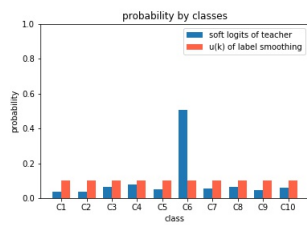Table 5. Hyper-parameters for $\text{TF}_{self}$ experiments (Temperature $\tau$ and $\alpha$)

| Dataset | Model | $\text{TF}_{self}$ |
|---|---|---|
| CIFAR100 | MobileNetV2 | $\tau$=20, $\alpha$=0.95 |
| | ShuffleNetV2 | $\tau$=20, $\alpha$=0.95 |
| | ResNet18 | $\tau$=6, $\alpha$=0.95 |
| | GoogLeNet | $\tau$=20, $\alpha$=0.40 |
| | DenseNet121 | $\tau$=20, $\alpha$=0.95 |
| | ResNeXt29 (8x64d) | $\tau$=20, $\alpha$=0.90 |
| T-ImageNet | MobileNetV2 | $\tau$=20, $\alpha$=0.10 |
| | ShuffleNetV2 | $\tau$=20, $\alpha$=0.10 |
| | ResNet18 | $\tau$=6, $\alpha$=0.10 |
| | ResNet50 | $\tau$=20, $\alpha$=0.10 |
| | DenseNet121 | $\tau$=20, $\alpha$=0.15 |

For all $\text{TF}_{reg}$ experiments on ImageNet, we set temperature $\tau$=20, and weight $\alpha$=0.10. The hyper-parameters on CIFAR100 and Tiny-ImageNet are given in Tab. 6.

Table 6. Hyper-parameters for $\text{TF}_{reg}$ experiments (Temperature $\tau$ and $\alpha$)

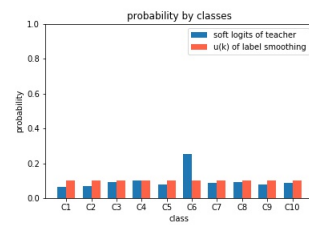| Dataset | Model | $\text{TF}_{reg}$ |
|---|---|---|
| CIFAR100 | MobileNetV2 | $\tau$=40, $\alpha$=0.95 |
| | ShuffleNetV2 | $\tau$=20, $\alpha$=0.95 |
| | ResNet18 | $\tau$=20, $\alpha$=0.10 |
| | GoogLeNet | $\tau$=40, $\alpha$=0.10 |
| T-ImageNet | MobileNetV2 | $\tau$=20, $\alpha$=0.10 |
| | ShuffleNetV2 | $\tau$=20, $\alpha$=0.10 |
| | ResNet50 | $\tau$=20, $\alpha$=0.10 |
| | DenseNet121 | $\tau$=20, $\alpha$=0.10 |

## References

[1] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 3

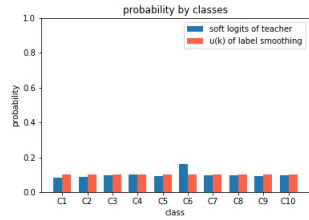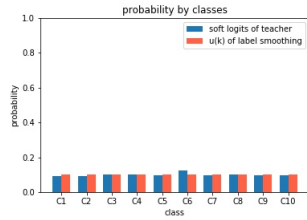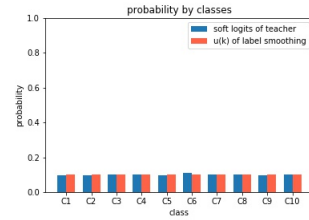Figure S.1. Comparison between label smoothing and soft targets of KD in different temperature $\tau$. C6 is the correct label.