

Logit Standardization in Knowledge Distillation

Shangquan Sun^{1,2}, Wenqi Ren^{3†}, Jingzhi Li¹, Rui Wang^{1,2}, Xiaochun Cao³

CVPR 2024
(Highlight)

- Problem/Objective

- Knowledge Distillation
- Classification Task

- Contribution/Key Idea

- Prove the irrelevance between T-S temperature
- Show the drawbacks of the shared Temperature
- Propose Logit Standardization



Shangquan Sun

University of Chinese Academy of Sciences

在 iie.ac.cn 的电子邮件经过验证 - [首页](#)

[Computer Vision](#) [Machine Learning](#)

标题

[Logit Standardization in Knowledge Distillation](#)

S Sun, W Ren, J Li, R Wang, X Cao

CVPR 2024 (Highlight)

[Rethinking image restoration for object detection](#)

S Sun, W Ren, T Wang, X Cao

NeurIPS 2022

[Event-aware video deraining via multi-patch progressive learning](#)

S Sun, W Ren, J Li, K Zhang, M Liang, X Cao

IEEE Transactions on Image Processing

[Restoring Images in Adverse Weather Conditions via Histogram Transformer](#)

S Sun, W Ren, X Gao, R Wang, X Cao

ECCV 2024

[DI-Retinex: Digital-Imaging Retinex Theory for Low-Light Image Enhancement](#)

S Sun, W Ren, J Peng, F Song, X Cao

arXiv preprint arXiv:2404.03327

[EnslR: An Ensemble Algorithm for Image Restoration via Gaussian Mixture Models](#)

S Sun, W Ren, Z Liu, H Park, R Wang, X Cao

NeurIPS 2024

- KD에 대해

$$\mathcal{L}_{KD} = \alpha \times \mathcal{L}_{KL}(y, q(z)) + \beta \times \mathcal{L}_{KL}(q(v; \tau), q(z; \tau))$$



GT → Student: $\mathcal{L}_{KL}(y, q(z))$



Teacher → Student: $\mathcal{L}_{KL}(q(v; \tau), q(z; \tau))$

- Distillation through **Softmax q + Temperature T** to soften pseudo-label
- Temperature은 softmax 시 확률 분포를 조정하는 값(>1이면 고르게, <1이면 차이 극대화)

- KD에 대해

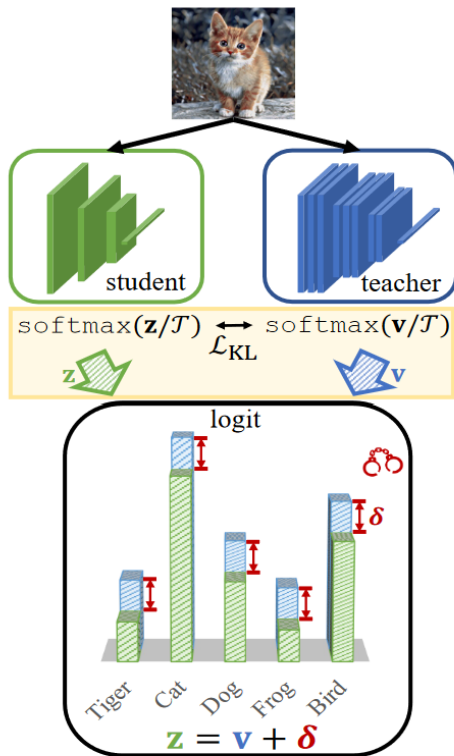
$$\mathcal{L}_{KD} = \alpha \times \mathcal{L}_{KL}(y, q(z)) + \beta \times \mathcal{L}_{KL}(q(v; \tau), q(z; \tau))$$

↓
GT → Student: $\mathcal{L}_{KL}(y, q(z))$

↓
Teacher → Student: $\mathcal{L}_{KL}(q(v; \tau), q(z; \tau))$

- Temperature T 가 공유되는 문제점 (Teacher & Student)

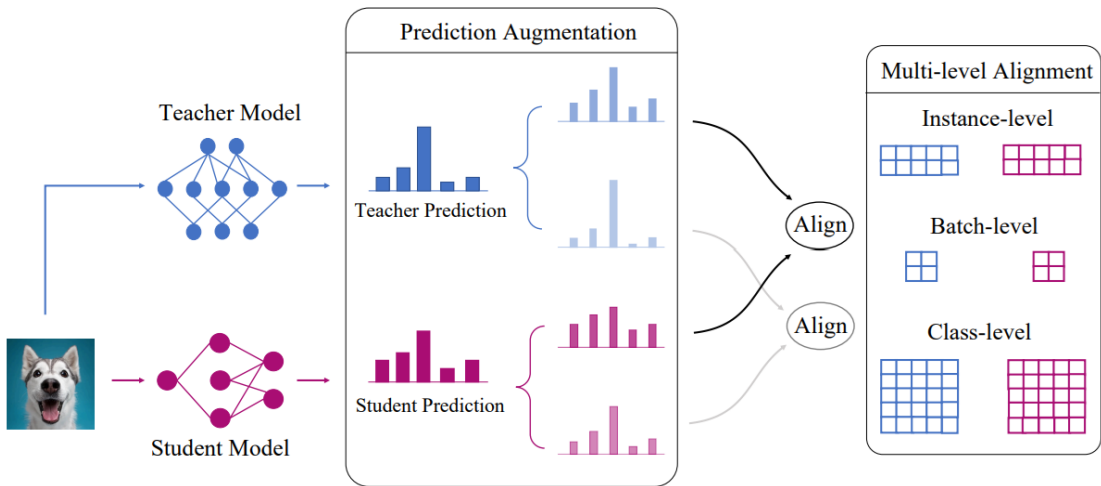
● KD에 대해



(a) Vanilla KD

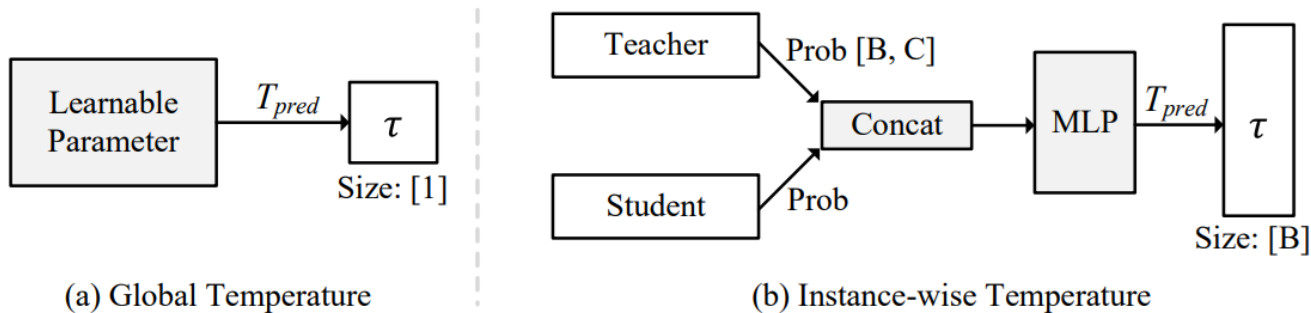
- Student는 독자적인 capability 존재
- Teacher = Student + δ 로 차이가 일정한 관계 (* δ 가 아님!!!)
→ 만약 * δ 였다면 같은 scale의 Temperature로 나눠주는 것이 Ok
- 그래서 Logit 에 Bias가 있는거는 현재 Hinton KD 수식으로 설명 불가

● Related Work



- 여러 Temperature (=Level) 에서 distillation 해보면 더 좋다.

● Related Work



- Sample에 따라서 Temperature를 다르게 학습해서, 각각의 sample에 맞는 Temperature을 주자.
→ 여전히 Teacher / Student가 Temperature 공유하는 문제점

Logit Standardization in Knowledge Distillation

Shangquan Sun^{1,2}, Wenqi Ren^{3†}, Jingzhi Li¹, Rui Wang^{1,2}, Xiaochun Cao³

• Related Work

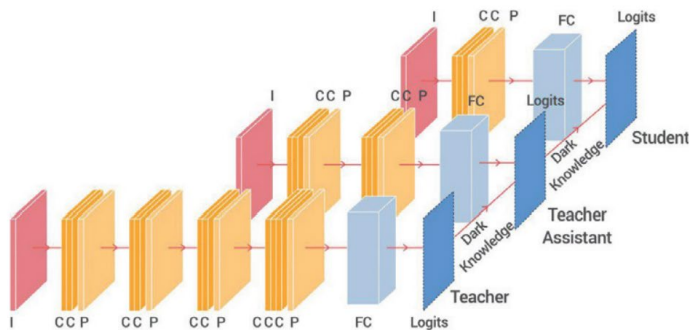


Figure 1: TA fills the gap between student & teacher

- Teacher - TA - Student 를 이용해서 Gap을 줄이려는 노력
→ But, 수식적으로 완벽하진 x

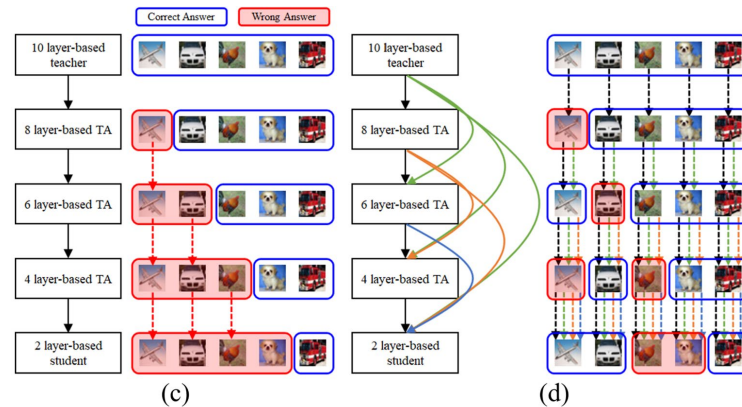


Figure 1. **Problem definition of the large gap between a teacher and a student network.** (a) In general, the difference between

[1] Mirzadeh, Seyed Iman, et al. "Improved knowledge distillation via teacher assistant." *AAAI* 2020

[2] Son, Wonchul, et al. "Densely guided knowledge distillation using multiple teacher assistants." *CVPR* 2021

• Background and Notation

K -class classification task with N samples $\{\mathbf{x}_n, y_n\}_{n=1}^N$.

Each $\mathbf{x}_n \in \mathbb{R}^{H \times W}$ is an input and y_n is its corresponding output.

Teacher f_T and student f_S produce logits $\mathbf{z}_n = f_S(\mathbf{x}_n)$ and $\mathbf{v}_n = f_T(\mathbf{x}_n)$, respectively, for n^{th} sample.

Softmax function $q(\mathbf{z}_n)^{(k)} = \frac{\exp(\mathbf{z}_n^{(k)}/\tau)}{\sum_{m=1}^K \exp(\mathbf{z}_n^{(m)}/\tau)}$, $q(\mathbf{v}_n)^{(k)} = \frac{\exp(\mathbf{v}_n^{(k)}/\tau)}{\sum_{m=1}^K \exp(\mathbf{v}_n^{(m)}/\tau)}$ with temperature τ is general.

Knowledge distillation essentially aims to let $q(\mathbf{z}_n)^{(k)}$ to mimic $q(\mathbf{v}_n)^{(k)}$.

$$\mathcal{L}_{\text{KL}}(q(\mathbf{z}_n) \| q(\mathbf{v}_n)) = \sum_{k=1}^K q(\mathbf{z}_n)^{(k)} \log \left(\frac{q(\mathbf{v}_n)^{(k)}}{q(\mathbf{z}_n)^{(k)}} \right).$$

\mathcal{L}_{KL} is “theoretically” equivalent to $\mathcal{L}_{\text{CE}}(q(\mathbf{z}_n), q(\mathbf{v}_n)) = - \sum_{k=1}^K q(\mathbf{z}_n)^{(k)} \log(q(\mathbf{z}_n)^{(k)})$.

• Methodology - Derivation of Softmax (Classification)

$$\max_q \mathcal{L}_1 = - \sum_{n=1}^N \sum_{k=1}^K q(\mathbf{v}_n)^{(k)} \log q(\mathbf{v}_n)^{(k)} \longrightarrow \text{Classification의 objective function은 entropy를 maximize 하는 문제로 볼 수 있음}$$

• 제약 조건

$$s.t. \begin{cases} \sum_{k=1}^K q(\mathbf{v}_n)^{(k)} = 1, \quad \forall n \longrightarrow q(\mathbf{v}_n) \text{이 probability distribution이 되게끔} \\ \mathbb{E}_q[\mathbf{v}_n] = \sum_{k=1}^K \mathbf{v}_n^{(k)} q(\mathbf{v}_n)^{(k)} = \mathbf{v}_n^{(y_n)}, \quad \forall n. \longrightarrow q(\mathbf{v}_n) \text{을 target one-hot vector (g^n)이 되게끔} \end{cases}$$

the target class. Suppose \hat{q}_n to be the one-hot hard probability distribution whose values are all zero except at the target index $\hat{q}_n^{(y_n)} = 1$. The second constraint is then actu-

• Methodology - Derivation of Softmax (Classification)

$$\mathcal{L}_T = \mathcal{L}_1 + \sum_{n=1}^N \alpha_{1,n} \left(\sum_{k=1}^K q(\mathbf{v}_n)^{(k)} - 1 \right) + \sum_{n=1}^N \alpha_{2,n} \left(\sum_{k=1}^K \mathbf{v}_n^{(k)} q(\mathbf{v}_n)^{(k)} - \mathbf{v}_n^{(y_n)} \right)$$



라그랑지안 multipliers로 미분해서 최대화하는 방향 찾기



편미분

$$\frac{\partial \mathcal{L}_T}{\partial q(\mathbf{v}_n)^{(k)}} = -1 - \log q(\mathbf{v}_n)^{(k)} + \alpha_{1,n} + \alpha_{2,n} \mathbf{v}_n^{(k)}$$

Softmax function $q(\mathbf{z}_n)^{(k)} = \frac{\exp(\mathbf{z}_n^{(k)}/\tau)}{\sum_{m=1}^K \exp(\mathbf{z}_n^{(m)}/\tau)}$

$$q(\mathbf{v}_n)^{(k)} = \exp(\alpha_{2,n} \mathbf{v}_n^{(k)}) / Z_T$$



softmax를 쓰는게 최적의 q(Vn)

$$Z_T = \exp(1 - \alpha_{1,n}) = \sum_{m=1}^K \exp(\alpha_{2,n} \mathbf{v}_n^{(m)})$$

• Methodology - Derivation of Softmax (KD)

$$\begin{aligned} \max_q \mathcal{L}_2 &= - \sum_{n=1}^N \sum_{k=1}^K q(\mathbf{z}_n)^{(k)} \log q(\mathbf{z}_n)^{(k)} \\ \text{s.t. } \begin{cases} \sum_{k=1}^K q(\mathbf{z}_n)^{(k)} = 1, & \forall n \\ \sum_{k=1}^K \mathbf{z}_n^{(k)} q(\mathbf{z}_n)^{(k)} = \mathbf{z}_n^{(y_n)}, & \forall n \\ \sum_{k=1}^K \mathbf{z}_n^{(k)} q(\mathbf{z}_n)^{(k)} = \sum_{k=1}^K \mathbf{z}_n^{(k)} q(\mathbf{v}_n)^{(k)}, & \forall n. \end{cases} \end{aligned}$$

$$\begin{aligned} \mathcal{L}_S &= \mathcal{L}_2 + \sum_{n=1}^N \beta_{1,n} \left(\sum_{k=1}^K q(\mathbf{z}_n)^{(k)} - 1 \right) \\ &\quad + \sum_{n=1}^N \beta_{2,n} \left(\sum_{k=1}^K \mathbf{z}_n^{(k)} q(\mathbf{z}_n)^{(k)} - \mathbf{z}_n^{(y_n)} \right) \\ &\quad + \sum_{n=1}^N \beta_{3,n} \sum_{k=1}^K \mathbf{z}_n^{(k)} \left(q(\mathbf{z}_n)^{(k)} - q(\mathbf{v}_n)^{(k)} \right) \end{aligned}$$

$$\frac{\partial \mathcal{L}_S}{\partial q(\mathbf{z}_n)^{(k)}} = -1 - \log q(\mathbf{z}_n)^{(k)} + \beta_{1,n} + \beta_{2,n} \mathbf{z}_n^{(k)} + \beta_{3,n} \mathbf{z}_n^{(k)}$$

$$q(\mathbf{z}_n)^{(k)} = \exp(\beta_n \mathbf{z}_n^{(k)}) / Z_S \rightarrow \text{softmax를 쓰는게 최적의 } q(\mathbf{Z}_n)$$

$$Z_S = \exp(1 - \beta_{1,n}) = \sum_{k=1}^K \exp(\beta_n \mathbf{z}_n^{(k)})$$

• Methodology - Derivation of Softmax

$$\begin{aligned} \operatorname{argmax}_{q(\mathbf{v}_n)^{(k)}} \mathcal{L}_T &= \frac{\exp(\alpha_{2,n} \mathbf{v}_n^{(k)})}{\sum_{m=1}^K \exp(\alpha_{2,n} \mathbf{v}_n^{(m)})} = \operatorname{softmax}(\mathbf{v}_n^{(k)}; \alpha_{2,n}^{-1}) \\ \operatorname{argmax}_{q(\mathbf{z}_n)^{(k)}} \mathcal{L}_S &= \frac{\exp(\beta_n \mathbf{z}_n^{(k)})}{\sum_{m=1}^K \exp(\beta_n \mathbf{z}_n^{(m)})} = \operatorname{softmax}(\mathbf{z}_n^{(k)}; \beta_n^{-1}) \end{aligned} \quad \left. \vphantom{\begin{aligned} \operatorname{argmax}_{q(\mathbf{v}_n)^{(k)}} \mathcal{L}_T &= \frac{\exp(\alpha_{2,n} \mathbf{v}_n^{(k)})}{\sum_{m=1}^K \exp(\alpha_{2,n} \mathbf{v}_n^{(m)})} = \operatorname{softmax}(\mathbf{v}_n^{(k)}; \alpha_{2,n}^{-1}) \\ \operatorname{argmax}_{q(\mathbf{z}_n)^{(k)}} \mathcal{L}_S &= \frac{\exp(\beta_n \mathbf{z}_n^{(k)})}{\sum_{m=1}^K \exp(\beta_n \mathbf{z}_n^{(m)})} = \operatorname{softmax}(\mathbf{z}_n^{(k)}; \beta_n^{-1}) \end{aligned}} \right\} \text{Optimal } q \text{ equivalents to softmax.}$$

$\beta_n = \alpha_{2,n} = 1/\mathcal{T} \rightarrow$ 일반적인 KD에서 쓰는 Temperature 식

$\beta_n = \alpha_{2,n} = 1 \rightarrow$ 일반적인 Classification에서 쓰는 softmax 식

+ 알파, 베타 달라도 된다 !! Temperature 달라도 된다!!

• Drawback of shared Temperature

$$q(\mathbf{z}_n; a_S, b_S)^{(k)} = \frac{\exp[(\mathbf{z}_n^{(k)} - a_S)/b_S]}{\sum_{m=1}^K \exp[(\mathbf{z}_n^{(m)} - a_S)/b_S]} \rightarrow \text{Softmax를 다음과 같은 식으로 가정 (a,b의 bias가 추가된)}$$

- Knowledge distillation에선 $\mathbf{Z}_n, \mathbf{V}_n$ 에 softmax했을때 같이 지게끔 하면된다

$$\frac{\exp[(\mathbf{z}_n^{(i)} - a_S)/b_S]}{\exp[(\mathbf{z}_n^{(j)} - a_S)/b_S]} = \frac{\exp[(\mathbf{v}_n^{(i)} - a_T)/b_T]}{\exp[(\mathbf{v}_n^{(j)} - a_T)/b_T]} \Rightarrow (\mathbf{z}_n^{(i)} - \mathbf{z}_n^{(j)})/b_S = (\mathbf{v}_n^{(i)} - \mathbf{v}_n^{(j)})/b_T$$

$$(\mathbf{z}_n^{(i)} - \bar{\mathbf{z}}_n)/b_S = (\mathbf{v}_n^{(i)} - \bar{\mathbf{v}}_n)/b_T \rightarrow \frac{\sigma(\mathbf{z}_n)^2}{\sigma(\mathbf{v}_n)^2} = \frac{\frac{1}{K} \sum_{i=1}^K (\mathbf{z}_n^{(i)} - \bar{\mathbf{z}}_n)^2}{\frac{1}{K} \sum_{i=1}^K (\mathbf{v}_n^{(i)} - \bar{\mathbf{v}}_n)^2} = \frac{b_S^2}{b_T^2}$$

1 to K 합

• Drawback of shared Temperature

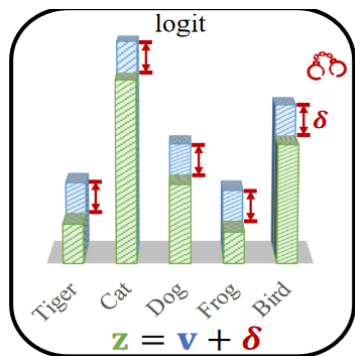
$$\frac{\sigma(\mathbf{z}_n)^2}{\sigma(\mathbf{v}_n)^2} = \frac{\frac{1}{K} \sum_{i=1}^K \left(\mathbf{z}_n^{(i)} - \bar{\mathbf{z}}_n \right)^2}{\frac{1}{K} \sum_{i=1}^K \left(\mathbf{v}_n^{(i)} - \bar{\mathbf{v}}_n \right)^2} = \frac{b_S^2}{b_T^2}$$

$$(b_S = b_T)$$

즉 if) Temperature 같다면 2가지 문제점

$$1) \quad \mathbf{z}_n^{(i)} = \mathbf{v}_n^{(i)} + \Delta_n,$$

$$2) \quad \sigma(\mathbf{z}_n) = \sigma(\mathbf{v}_n)$$



This is another shackle applied to student restricting the standard deviation of its predicted logits. In contrast, since

결론) Teacher - Student 서로 temperature 다르게 해야한다.

• propose method

- Z-score 정규화

Algorithm 1: Weighted \mathcal{Z} -score function.

Input: Input vector \mathbf{x} and Base temperature τ

Output: Standardized vector $\mathcal{Z}(\mathbf{x}; \tau)$

```

1  $\bar{\mathbf{x}} \leftarrow \frac{1}{K} \sum_{k=1}^K \mathbf{x}^{(k)}$ 
2  $\sigma(\mathbf{x}) \leftarrow \sqrt{\frac{1}{K} \sum_{k=1}^K (\mathbf{x}^{(k)} - \bar{\mathbf{x}})^2}$ 
3 return  $(\mathbf{x} - \bar{\mathbf{x}}) / \sigma(\mathbf{x}) / \tau$ 

```

정규분포 Z score
$$Z = \frac{x - \mu}{\sigma}$$

본 논문 Z-score
$$Z(x; \tau) = \frac{x - \bar{x}}{\sigma(x) \cdot \tau}$$

Algorithm 2: \mathcal{Z} -score logit standardization pre-process in knowledge distillation.

Input: Transfer set \mathcal{D} with image-label sample pair $\{\mathbf{x}_n, y_n\}_{n=1}^N$, Base Temperature τ , Teacher f_T , Student f_S , Loss \mathcal{L}_{KD} (e.g., \mathcal{L}_{KL}), loss weight λ , and \mathcal{Z} -score function \mathcal{Z} in Algo. 1

Output: Trained student model f_S

```

1 foreach  $(\mathbf{x}_n, y_n)$  in  $\mathcal{D}$  do
2    $\mathbf{v}_n \leftarrow f_T(\mathbf{x}_n), \mathbf{z}_n \leftarrow f_S(\mathbf{x}_n)$ 
3    $q(\mathbf{v}_n) \leftarrow \text{softmax}[\mathcal{Z}(\mathbf{v}_n; \tau)]$ 
4    $q(\mathbf{z}_n) \leftarrow \text{softmax}[\mathcal{Z}(\mathbf{z}_n; \tau)]$ 
5    $q'(\mathbf{z}_n) \leftarrow \text{softmax}(\mathbf{z}_n)$ 
6   Update  $f_S$  towards minimizing
      $\lambda_{\text{CE}} \mathcal{L}_{\text{CE}}(y_n, q'(\mathbf{z}_n)) + \lambda_{\text{KD}} \tau^2 \mathcal{L}(q(\mathbf{v}_n), q(\mathbf{z}_n))$ 
7 end

```

• Toy case

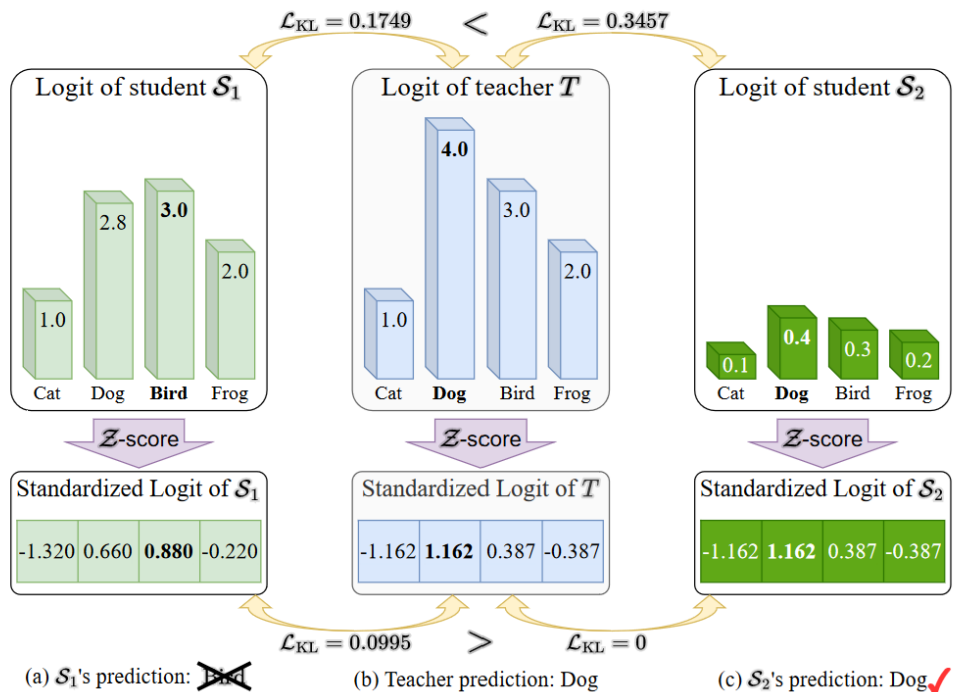


Figure 2. A toy case where two students, \mathcal{S}_1 and \mathcal{S}_2 , learning

Logit Standardization in Knowledge Distillation

Shangquan Sun^{1,2}, Wenqi Ren^{3†}, Jingzhi Li¹, Rui Wang^{1,2}, Xiaochun Cao³

CVPR 2024
(Highlight)

Experiment - CIFAR 100

Table 1. The Top-1 Accuracy (%) of different knowledge distillation methods on the validation set of CIFAR-100 [18]. The teacher and student have distinct architectures. The KD methods are sorted by the types, i.e., feature-based and logit-based. We apply our logit standardization to the existing logit-based methods and use Δ to show its performance gain. The values in **blue** denote slight enhancement and those in **red** non-trivial enhancement no less than 0.15. The best and second best results are emphasized in **bold** and underlined cases.

Type	Teacher	ResNet32×4 79.42	ResNet32×4 79.42	ResNet32×4 79.42	WRN-40-2 75.61	WRN-40-2 75.61	VGG13 74.64	ResNet50 79.34
	Student	SHN-V2 71.82	WRN-16-2 73.26	WRN-40-2 75.61	ResNet8×4 72.50	MN-V2 64.60	MN-V2 64.60	MN-V2 64.60
Feature	FitNet [33]	73.54	74.70	77.69	74.61	68.64	64.16	63.16
	AT [53]	72.73	73.91	77.43	74.11	60.78	59.40	58.58
	RKD [31]	73.21	74.86	77.82	75.26	69.27	64.52	64.43
	CRD [39]	75.65	75.65	78.15	75.24	70.28	69.73	69.11
	OFD [12]	76.82	76.17	79.25	74.36	69.92	69.48	69.04
	ReviewKD [5]	77.78	76.11	78.96	74.34	<u>71.28</u>	70.37	69.89
	SimKD [4]	78.39	77.17	79.29	75.29	70.10	69.44	69.97
	CAT-KD [10]	78.41	76.97	78.59	75.38	70.24	69.13	71.36
	KD [13]	74.45	74.90	77.70	73.97	68.36	67.37	67.35
	KD+Ours	75.56	75.26	77.92	77.11	69.23	68.61	69.02
	Δ	1.11	0.36	0.22	3.14	0.87	1.24	1.67
Logit	CTKD [24]	75.37	74.57	77.66	74.61	68.34	68.50	68.67
	CTKD+Ours	76.18	75.16	77.99	77.03	69.53	68.98	69.36
	Δ	0.81	0.59	0.33	2.42	1.19	0.48	0.69
	DKD [57]	77.07	75.70	78.46	75.56	69.28	69.71	70.35
	DKD+Ours	77.37	76.19	78.95	76.75	70.01	69.98	70.45
	Δ	0.30	0.49	0.49	1.19	0.73	0.27	0.10
	MLKD [17]	78.44	76.52	79.26	<u>77.33</u>	70.78	<u>70.57</u>	71.04
	MLKD+Ours	78.76	77.53	79.66	77.68	71.61	70.94	71.19
	Δ	0.32	1.01	0.40	0.35	0.83	0.37	0.15

Heterogeneous

27 **red** 1 **blue**

Table 2. The Top-1 Accuracy (%) of different knowledge distillation methods on the validation set of CIFAR-100 [18]. The teacher and student have identical architectures but different configurations. The KD methods are sorted by the types. We apply our logit standardization to the existing logit-based methods and use Δ to show its performance gain. The values in **blue** denote slight enhancement and those in **red** non-trivial enhancement no less than 0.15. The best and second best results are emphasized in **bold** and underlined cases.

Type	Teacher	ResNet32×4 79.42	VGG13 74.64	WRN-40-2 75.61	WRN-40-2 75.61	ResNet56 72.34	ResNet110 74.31	ResNet110 74.31
	Student	ResNet8×4 72.50	VGG8 70.36	WRN-40-1 71.98	WRN-16-2 73.26	ResNet20 69.06	ResNet32 71.14	ResNet20 69.06
Feature	FitNet [33]	73.50	71.02	72.24	73.58	69.21	71.06	68.99
	AT [53]	73.44	71.43	72.77	74.08	70.55	72.31	70.65
	RKD [31]	71.90	71.48	72.22	73.35	69.61	71.82	69.25
	CRD [39]	75.51	73.94	74.14	75.48	71.16	73.48	71.46
	OFD [12]	74.95	73.95	74.33	75.24	70.98	73.23	71.29
	ReviewKD [5]	75.63	74.84	75.09	76.12	71.89	73.89	71.34
	SimKD [4]	78.08	74.89	74.53	75.53	71.05	73.92	71.06
	CAT-KD [10]	76.91	74.65	74.82	75.60	71.62	73.62	71.37
	KD [13]	73.33	72.98	73.54	74.92	70.66	73.08	70.67
	KD+Ours	76.62	74.36	74.37	76.11	71.43	74.17	71.48
	Δ	3.29	1.38	0.83	1.19	0.77	1.09	0.81
Logit	CTKD [24]	73.39	73.52	73.93	75.45	71.19	73.52	70.99
	KD+CTKD+Ours	76.67	74.47	74.58	76.08	71.34	74.01	71.39
	Δ	3.28	0.95	0.65	0.63	0.15	0.49	0.40
	DKD [57]	76.32	74.68	74.81	76.24	71.97	74.11	71.06
	DKD+Ours	77.01	74.81	74.89	76.39	72.32	74.29	71.85
	Δ	0.69	0.13	0.08	0.15	0.35	0.18	0.79
	MLKD [57]	77.08	<u>75.18</u>	<u>75.35</u>	<u>76.63</u>	72.19	74.11	71.89
	MLKD+Ours	78.28	75.22	75.56	76.95	72.33	74.32	72.27
	Δ	1.20	0.04	0.21	0.32	0.14	0.21	0.38

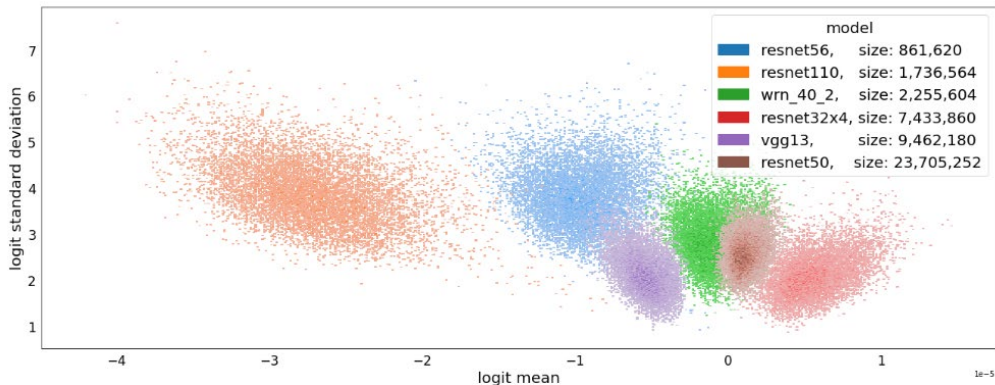
Homogeneous

24 **red** 4 **blue**

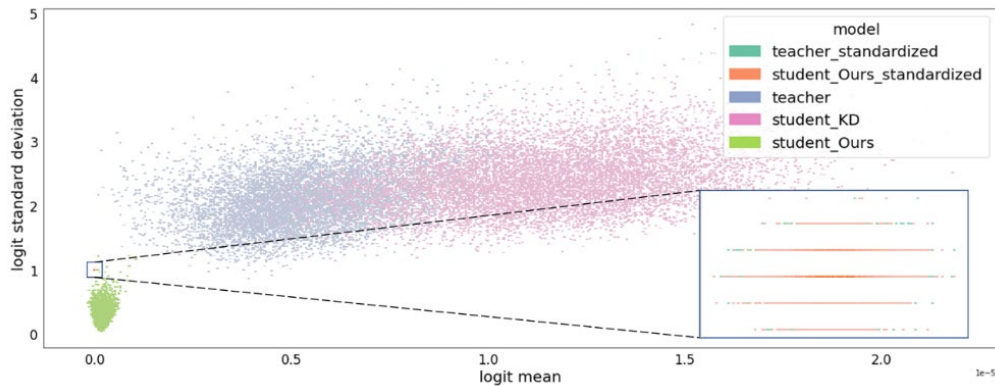
● **Experiment - ImageNet**

Teacher/Student	ResNet34/ResNet18		ResNet50/MN-V1	
Accuracy	top-1	top-5	top-1	top-5
Teacher	73.31	91.42	76.16	92.86
Student	69.75	89.07	68.87	88.76
AT [53]	70.69	90.01	69.56	89.33
OFD [12]	70.81	89.98	71.25	90.34
CRD [39]	71.17	90.13	71.37	90.41
ReviewKD [5]	71.61	90.51	72.56	91.00
SimKD [4]	71.59	90.48	72.25	90.86
CAT-KD [10]	71.26	90.45	72.24	91.13
KD [13]	71.03	90.05	70.50	89.80
KD+Ours	71.42 ^{+0.39}	90.29 ^{+0.24}	72.18 ^{+1.68}	90.80 ^{+1.00}
KD+CTKD [24]	71.38	90.27	71.16	90.11
KD+CTKD+Ours	71.81 ^{+0.43}	90.46 ^{+0.19}	72.92 ^{+1.76}	91.25 ^{+1.14}
DKD [57]	71.70	90.41	72.05	91.05
DKD+Ours	71.88 ^{+0.18}	90.58 ^{+0.17}	72.85 ^{+0.80}	91.23 ^{+0.18}
MLKD [17]	71.90	90.55	73.01	91.42
MLKD+Ours	72.08 ^{+0.18}	90.74 ^{+0.19}	73.22 ^{+0.21}	91.59 ^{+0.17}

● Experiment - logit space 시각화



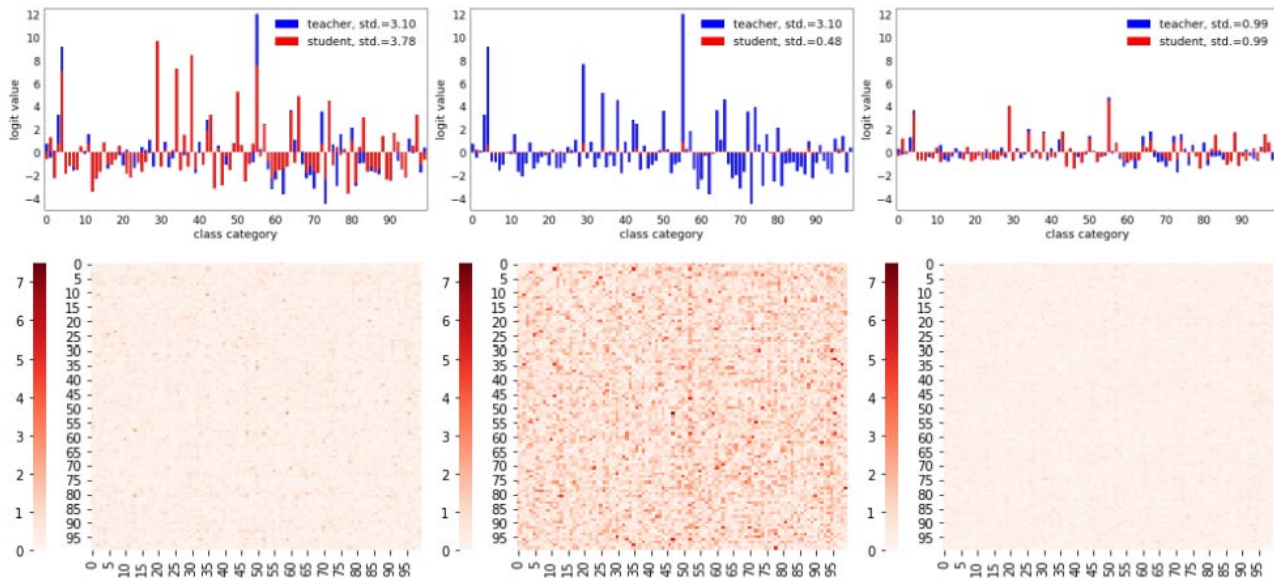
(a) Teacher models of different sizes



Teacher 모델에 따라서 Logit의 scale이 다르다
(parameter마다)

Standardized 했을때 엄청 넓게 퍼져있다가
엄청 좁은 영역에 일정하게 분포함

● Experiment - 시각화



(a) Vanilla KD

(b) Ours w/o \mathcal{Z} -score

(c) Ours w/ \mathcal{Z} -score

Mean: 0.27, Max: 3.03. Mean: 0.94, Max: 7.36. Mean: 0.18, Max:1.18.

평균만 빼면 일정한 격차로 유지

표준편차로 나눠야만 같은 스케일

- Experiment - robust

Table 4. The ablation studies under different settings in our \mathcal{Z} -score. The base temperature τ is set to be 2. By default $\lambda_{\text{CE}} = 0.1$. The logit vector of teacher \mathbf{v}_n and student \mathbf{z}_n are abbreviated as \mathbf{z} for succinctness. The teacher and student are ResNet32 \times 4 and ResNet8 \times 4.

λ_{KD}	\mathbf{z} (KD)	$\mathbf{z} - \bar{\mathbf{z}}$	$\frac{\mathbf{z}}{\sigma(\mathbf{z})}$	$\frac{(\mathbf{z} - \bar{\mathbf{z}})}{\sigma(\mathbf{z})}$ (Ours)
0.9	73.60	73.37	73.79	74.14
3.0	74.38	74.33	75.86	76.11
6.0	74.45	74.82	76.44	76.56
9.0	73.33	73.94	76.30	76.62
12.0	68.29	71.56	76.49	76.56
15.0	65.34	62.01	76.42	76.61
18.0	63.45	61.31	76.18	76.33