

FerKD: Surgical Label Adaptation for Efficient Distillation

Zhiqiang Shen

Mohamed bin Zayed University of AI

Zhiqiang.Shen@mbzuai.ac.ae

Abstract

We present *FerKD*, a novel efficient knowledge distillation framework that incorporates partial soft-hard label adaptation coupled with a region-calibration mechanism. Our approach stems from the observation and intuition that standard data augmentations, such as *RandomResizedCrop*, tend to transform inputs into diverse conditions: easy positives, hard positives, or hard negatives. In traditional distillation frameworks, these transformed samples are utilized equally through their predictive probabilities derived from pretrained teacher models. However, merely relying on prediction values from a pretrained teacher, a common practice in prior studies, neglects the reliability of these soft label predictions. To address this, we propose a new scheme that calibrates the less-confident regions to be the context using softened hard groundtruth labels. Our approach involves the processes of **hard regions mining + calibration**. We demonstrate empirically that this method can dramatically improve the convergence speed and final accuracy. Additionally, we find that a consistent mixing strategy can stabilize the distributions of soft supervision, taking advantage of the soft labels. As a result, we introduce a stabilized *SelfMix* augmentation that weakens the variation of the mixed images and corresponding soft labels through mixing similar regions within the same image. *FerKD* is an intuitive and well-designed learning system that eliminates several heuristics and hyperparameters in former *FKD* solution [37]. More importantly, it achieves remarkable improvement on *ImageNet-1K* and downstream tasks. For instance, *FerKD* achieves 81.2% on *ImageNet-1K* with *ResNet-50*, outperforming *FKD* and *FunMatch* by remarkable margins. Leveraging better pre-trained weights and larger architectures, our finetuned *ViT-G14* even achieves 89.9%. Our code is available at <https://github.com/szq0214/FKD/tree/main/FerKD>.

1. Introduction

Knowledge Distillation (KD) [13] has achieved impressive results in various visual domains, including image clas-

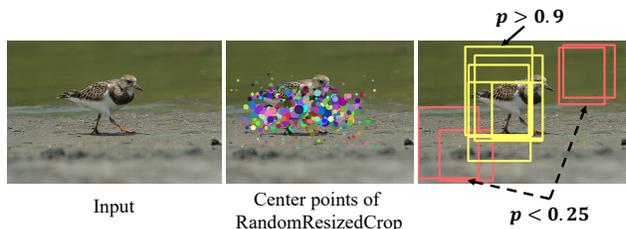


Figure 1: Illustration of motivation for *FerKD*. The left figure depicts the original input, and the middle figure shows the center points of bounding boxes generated using *RandomResizedCrop*. The radius of each circle corresponds to the area of the bounding box. It can be observed that the center points of the bounding boxes are concentrated in the center of the image, and their area increases as they approach the center. The right figure displays several top and bottom confident bounding boxes and their corresponding predictive probabilities from a pre-trained teacher or teachers ensemble. The proposed hard region calibration strategy is established based on these predictions.

sification [52, 36, 5, 37], object detection [6, 46, 11, 8, 53] and semantic segmentation [22, 15, 16]. However, KD methods are often computationally expensive and inefficient due to the additional computational burden imposed by the teacher models. The primary advantage of KD that motivates its usage is its ability to generate precise soft labels that convey more informative details about the input examples. It differs from other label softening techniques, such as label smoothing [43], Mixup [58], and CutMix [56], mainly in two aspects: (1) KD generates soft labels dynamically in each iteration, which is more informative than fixed smoothing patterns used in label smoothing; (2) Mixup and CutMix techniques essentially combine hard labels with coefficients, while KD produces soft labels that are highly correlated with the input sample. This allows KD’s soft labels to become more accurate when different data augmentations, such as *RandomResizedCrop*, *flipping and rotation*, *color jittering*, etc., are applied. In general, mixing-based label softening methods cannot monitor such changes in input content, but KD can address them effortlessly.

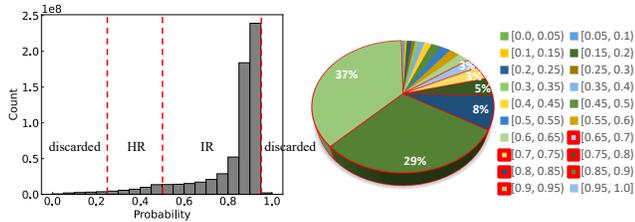


Figure 2: Statistics of soft label max-probability for crops on ImageNet-1K. The soft label is from FKD [37]. In each image, 500 regions are randomly cropped.

To overcome the computational inefficiency of traditional knowledge distillation, FKD [37] was developed to generate region-level soft labels in advance and reuse them across multiple training cycles to eliminate redundant computation. This approach only requires the preparation of soft labels once at the beginning and they can be reused indefinitely. However, this approach overlooks certain critical issues. One is the quality of the soft labels. When using *RandomResizedCrop* to generate regions, some may be cropped from background areas, and the teacher model will still produce a soft label for them based on their similarity to the dataset classes. However, in some cases, these areas may contain irrelevant noise, compensatory information, or context information for the class, and the soft labels may not accurately reflect the context of information they carry. To address this problem, this work proposes to recalibrate these soft labels by incorporating context information from hard ground-truth labels with smoothing.

Furthermore, due to the random nature of the sampling process, a certain proportion of crops that are either excessively easy or difficult do not contribute to the model’s learning capacity. As demonstrated in Fig. 2 and Table 1, these samples can be discarded to expedite the convergence process. The pre-generated soft labels can be utilized as useful indicators to select these specific samples. In our adaptation of surgical soft labeling, we categorize the soft labels into four distinct groups: extreme hard (negative), moderate hard (background or context), hard positive (partial object), and easy positive. Each of these categories is subject to different treatment methodologies.

The Role of Background. The role of the background in images is essential, as it provides critical context and spatial information that aids the model in accurately identifying objects of interest within the scene. Backgrounds can vary in complexity and structure, ranging from simple monochromatic backgrounds to highly cluttered and detailed ones. Soft labels in background areas are typically low, and therefore, it is crucial to handle the background carefully with precise supervision to achieve higher model capability within our surgical label calibration framework.

Hard Regions Mining and Calibration. Hard Regions Mining involves the identification and isolation of chal-

range (P)	ratio	range (P)	agg. ratio
[0.0, 0.1)	0.43%	[0, 0.1)	0.43%
[0.1, 0.2)	0.89%	[0, 0.2)	1.32%
[0.2, 0.3)	1.29%	[0, 0.3)	2.61%
[0.3, 0.4)	2.03%	[0, 0.4)	4.64%
[0.4, 0.5)	3.66%	[0, 0.5)	8.31%
[0.5, 0.6)	4.35%	[0, 0.6)	12.65%
[0.6, 0.7)	5.04%	[0, 0.7)	17.69%
[0.7, 0.8)	7.76%	[0, 0.8)	25.45%
[0.8, 0.85)	8.14%	[0, 0.85)	33.59%
[0.85, 0.9)	28.73%	[0, 0.9)	62.32%
[0.9, 0.95)	37.34%	[0, 0.95)	99.67%
[0.95, 1.0)	0.33%	[0, 1.0)	100%

Table 1: Detailed statistics of soft labels. “range” indicates max-probability of crops, “ratio” indicates the percentage in the whole crops. “agg. ratio” is the aggregated ratio.

lenging or complex regions within an image that the model struggles to identify accurately. These regions can include objects with complex shapes, occlusions, or those with low contrast. By identifying these regions, the model can focus on learning the features and characteristics of these regions, resulting in improved performance. Calibration, on the other hand, involves adjusting the confidence levels of the model’s predictions in challenging regions. The model’s predictions may be less reliable in hard regions, leading to lower confidence scores. Calibrating the predictions can improve the model’s accuracy in these regions by adjusting the confidence levels of the predictions. We found that carefully discarding a portion of negative crops and selecting those hard positive crops by calibrating their labels, can force the training process more efficient and effective.

Stable Training on Soft Labels. Mixture-based augmentations, such as Mixup and CutMix have seen widespread use for training models under hard supervision, where each image is labeled with a single class label. However, in the soft label scenario, we have made a different observation: when employed together with pre-generated soft labels on a typical ResNet, Mixup and CutMix tend to be overly strong, which, conversely, leads to decreased performance. To mitigate label fluctuations and achieve more stable training, we propose a *SelfMix* scheme, which is particularly suitable for cases where data augmentation should not be so strong, such as in finetuning distillation, where mixture-based augmentation is usually disabled. On the other hand, when training ViT models from scratch, stronger data augmentation can yield better results [45, 42], which is consistent with the larger capacity perspective of this type of network.

In summary, our contributions of this work are:

- We present *FerKD*, a sample-calibration framework for *Faster Knowledge Distillation* that achieves state-of-the-art performance. We conduct extensive analysis, ablation, and discussion on the impact of hard and easy samples.

- We make two key observations in the pre-generated soft

label training framework. Firstly, we observe that the few most challenging and simplest crops obtained through the `RandomResizedCrop` operation do not contribute significantly to the model’s learning and can therefore be removed. Secondly, we find that moderately hard crops can provide crucial contextual information that improves the model’s ability to learn robust representations.

- We perform extensive experiments on ImageNet-1K and downstream tasks. On ImageNet-1K, `FerKD` achieves an accuracy of 81.2% using the ResNet-50. When leveraging self-supervised pre-trained weights, our larger model finetuned using ViT-G/14 achieves an accuracy of 89.9%.

2. Related Work

Knowledge Distillation and Fast Knowledge Distillation.

Knowledge Distillation [13] is a learning method in which a “student” model is trained to imitate the predictions of a larger, more complex “teacher” model. A key advantage of this approach is that the teacher model can provide soft supervision that contains more information regarding the input data than traditional one-hot human annotated labels, particularly when the input data is subject to data augmentation. There have been many recent variants and extensions of knowledge distillation [31, 26, 47, 59, 25, 27, 35, 7, 41, 50, 45, 18, 51], including approaches that use internal feature representations [7], adversarial training with discriminators [33], transfer learning techniques [54], fast distillation [37] via instance label preparation, and methods that prioritize patient and consistent learning [5].

Hard Sample Mining. The aim of hard sample mining [23, 38] is to enhance the performance of learning models by selectively focusing on challenging examples that are typically difficult to classify. By prioritizing hard samples during training, models can be trained to better handle a wider range of real-world scenarios and improve their overall performance. One approach [38] to achieve this for object detection is to use Online Hard Example Mining (OHEM) which employs a strategy of selecting challenging examples during training of region-based ConvNet detectors. The motivation behind this approach is that detection datasets typically comprise an overwhelming number of easy examples and a small number of hard examples. Automatic selection of these hard examples can make training more effective and efficient.

Some other techniques that are close to hard sample mining: (1) Curriculum learning [3, 48, 40]: it trains a model on easy examples first and then gradually increasing the difficulty of the examples over time. (2) Active learning [32, 24, 28]: it selects the most informative or uncertain samples for labeling by a human annotator. By focusing on the samples that the model is most uncertain about, the model can learn to better generalize and improve its performance on difficult samples. (3) Loss functions: Attentive

loss functions can be used to emphasize the importance of hard samples during training. For example, focal loss [20] places more weight on the difficult examples during training, helping the model to learn to handle them better.

Data Augmentations. Several studies have incorporated data augmentations into distillation frameworks to improve performance. For instance, FunMatch [5] employed Mixup, and FKD [37] utilized CutMix. Both techniques achieved competitive accuracy on large-scale ImageNet-1K dataset. In this work, we investigate the impact of data augmentation intensity on soft labels. We discover that different network architectures require unique data augmentation levels. Specifically, ResNet necessitates mild data augmentation, while ViTs require stronger data augmentation. However, even for ViTs, finetuning distillation requires a reduction in the intensity of data augmentation, especially for mixture-based methods. Motivated by this observation, we propose a mild `SelfMix` approach for ResNet and finetuning distillation scenarios.

3. Approach

The proposed `FerKD` is a soft label calibration framework for fast and efficient knowledge distillation training. In this section, we aim to provide an elaborated overview of our method, starting with an in-depth analysis of the roles of hard and soft labels in the distillation process. We then present the key components of our approach, which include a region selection strategy and a soft label calibration scheme. Additionally, we explore the data augmentation requirements for soft labels and introduce a simple label ensemble technique to enhance the quality of the soft labels.

3.1. Revisiting Hard and Soft Label in Distillation

The utilization of hard and soft labels is dependent on the comprehension, problem scenarios, and underlying objectives. Various arguments exist regarding their effectiveness. The vanilla knowledge distillation method [13] employs both hard and soft labels to maximize the benefits of both. However, recent researches [36, 35, 57, 37] suggest that the use of hard labels is not necessary in distillation on large-scale datasets as strong teachers can provide more precise soft supervision. Incorporating hard labels may introduce erroneous supervisory signals, ultimately hampering student performance. Other than the above views, this work presents a novel perspective beyond them by acknowledging that both hard and soft labels offer unique advantages, highlights the proper practice of usage that is necessary, and finally proposes an elegant solution to reap benefits derived from both sources of hard and soft labels.

Combination of Soft and Hard Labels. In vanilla KD design [13], for each training example, it will minimize two loss terms for both hard and soft labels. The final objective

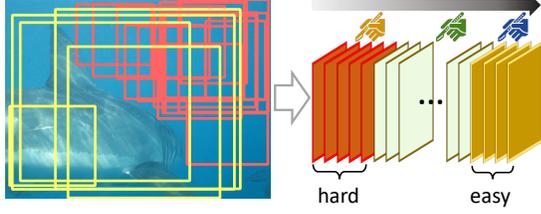


Figure 3: Illustration of region calibration according to their predictive probabilities in FerKD. Left is the input image with RandomResizedCrop. **Bounding box** is with high probability and **bounding box** is with low probability.

\mathcal{L}_{VKD} can be formulated as:

$$\mathcal{L}_{VKD} = \frac{1}{N} \sum_x (\underbrace{\alpha * \mathcal{L}_h(p_\theta(x), y_h(x))}_{\text{CE loss with hard label}} + (1-\alpha) * \underbrace{\mathcal{L}_s(p_\theta(x), y_s(x))}_{\text{KL loss with soft label}}) \quad (1)$$

where α is the coefficient to balance the loss signal intensity from soft and hard labels. $p_\theta(x)$ is the logits prediction from the student model and θ is its parameters. y_h is the hard label and y_s is the soft label from a pre-trained teacher. N is the total number of training samples x .

Full Soft Label Training. Fast KD [37] proposes to use the soft label solely since the prediction from strong teachers is precise enough. Thus, the loss objective \mathcal{L}_{FKD} is:

$$\mathcal{L}_{FKD} = \frac{1}{N} \sum_x \underbrace{\hat{\mathcal{L}}_s(p_\theta(x), y_s(x))}_{\text{SCE loss with soft label}} \quad (2)$$

where ‘‘SCE’’ is the soft version of cross-entropy loss.

3.2. FerKD: Surgical Label Calibration Distillation

Surgical/Partial Soft and Hard Label Adaptive Training. Different from VKD, the proposed FerKD strategy will only involve one loss term but also unlike FKD, FerKD will employ both hard and soft labels in a single objective term and exploit the additional information derived from both sources. The solution is that we only keep their original soft labels for positive regions since they contain fine-grained information regarding the crops, for those background or context regions, we will re-calibrate them by the human-annotated ground-truth labels to avoid misinformation from the soft labels. Hence, the loss function will be:

$$\mathcal{L}_{FerKD} = \frac{1}{N} \sum_x \underbrace{\mathcal{L}_{\text{adap: h or s?}}(p_\theta(x), y_a(x))}_{\text{SCE loss with hard/soft label}} \quad (3)$$

where ‘‘adap:’’ indicates $\{hard\}$ or $\{soft\}$ used for individual training samples. y_a is the calibrated soft labels. As illustrated in Fig. 3, we will calibrate regions’ soft labels y_a using the following rule:

$$y_a = \begin{cases} \mathbf{UR} : \text{discarded} & \text{if } y_s < \mathcal{T}_L \text{ or } y_s > \mathcal{T}_T \\ \mathbf{HR} : 1.0 - \varepsilon & \text{if } \mathcal{T}_L < y_s < \theta_M \\ \mathbf{IR} : y_s & \text{otherwise} \end{cases} \quad (4)$$

where \mathcal{T}_L , \mathcal{T}_M , and \mathcal{T}_T are thresholds at low, middle, and top boundaries. ‘‘UR’’ represents the uninformative regions such as black or white blocks in an image that will be discarded during training. As shown in Fig. 4, ‘‘HR’’ represents the hard regions with a smoothing value ε and ‘‘IR’’ represents the important regions. Thus, the key goal in FerKD becomes to identify the positive or background regions. Thanks to FKD [37], we can have access to all crops’ individual predictions. A quick exploration is performed for verification and the result is shown in Fig. 5. It is clear that by discarding a certain ratio of samples (hardest and easiest), the performance is consistently improved.

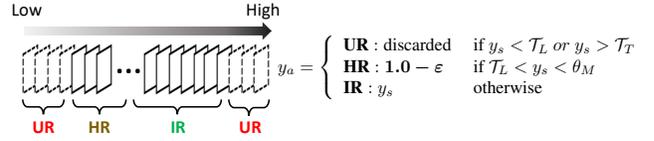


Figure 4: Illustration of region calibration according to their predictive probabilities in FerKD. Left is the input image with RandomResizedCrop. Right is the rule for calibrating the probabilities of regions.

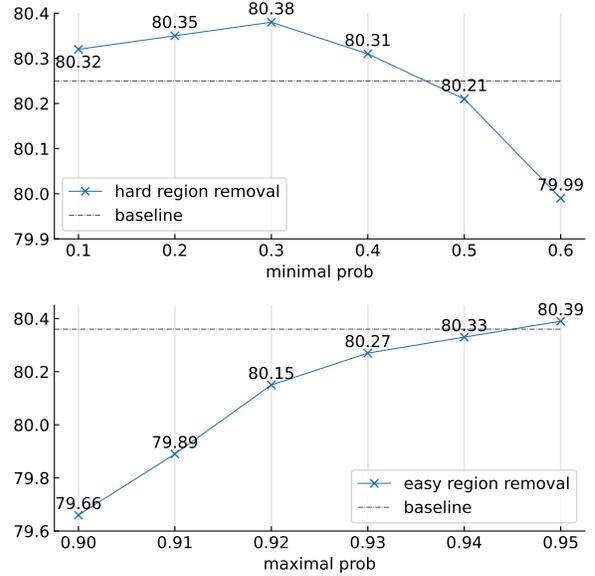


Figure 5: **Minimal and maximal probability.** The upper figure indicates that only regions having the max probability in $[\text{minimal}, 1.0]$ will be trained, and *baseline* indicates that the model is trained with all randomly sampled regions. The bottom figure indicates that only regions having the max probability in $[0.3, \text{maximal}]$ will be trained, and *baseline* indicates that the model is trained with regions in $[0.3, 1.0]$.

3.3. Training Speed of FerKD

When training a model, easy examples allow the model to quickly learn the patterns in the data and update its parameters in a way that minimizes the loss function. This means that the model will converge faster and require fewer

iterations to reach a satisfactory level of accuracy. On the other hand, hard examples can slow down the convergence speed but they can force models to learn more robust classify boundaries. The speed of model convergence can be impacted significantly by the sampling strategy of hard and easy examples. In our FerKD framework, “hard” and “easy” examples refer to the level of difficulty a particular sample presents to a model. This difficulty can be quantified by the probability assigned to the correct label by the teacher. An “easy” sample is one where the probability assigned to the correct label is high, indicating that the teacher is confident in its prediction. On the other hand, a “hard” sample is one where the probability assigned to the correct label is low, indicating that the teacher is uncertain about its prediction. The threshold for what constitutes an “easy” or “hard” sample can vary depending on the specific task and model being used, which is the key for exploring in FerKD.

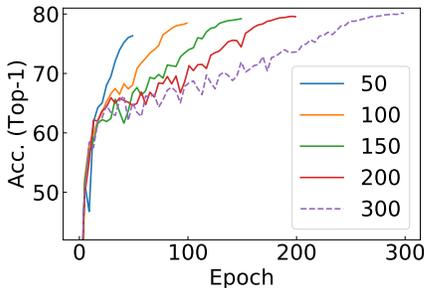


Figure 6: Illustration of the testing accuracy curves. “300” represents the training with the full budget. “50, 100, 150, 200” are the training with reduced budgets.

Our sampling strategy involves the removal of the easiest and hardest training examples to reduce computational costs incurred during uninformative training steps. This strategy offers a clear benefit, as shown in Fig. 6, which presents the test accuracy curves for different training budgets. The results show that, with a training budget reduced by $\frac{2}{3}$ (200 epochs), the achieved accuracy is comparable to that of full-budget training. Although further reducing the budget slightly affects accuracy, our proposed FerKD method is shown to be robust across different training budgets. Moreover, our calibration process will not involve additional training cost since it can be done offline in advance.

3.4. SelfMix: A Mild and Stable Data Augmentation for Soft Labels

Soft label Calibration with SelfMix. The current prevailing data augmentation techniques are designed for hard labels or smoothed hard labels. However, in the dynamic soft label scenario, a different approach is necessary to meet the unique attributes of pre-generated soft labels. Soft labels themselves can mitigate overfitting, making it imperative to develop tailored data augmentation techniques that account for soft labels to achieve improved accuracy. The proposed SelfMix is based on the empirical observations

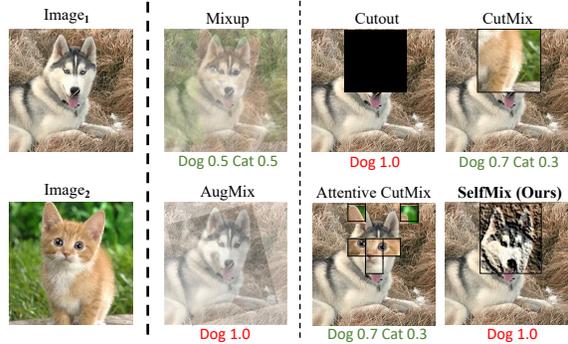


Figure 7: Illustration of the high-level outline for SelfMix to make mixed soft label consistent, reducing the variance. Practically, each label will be a soft distribution for the input instead of the hard label as illustrated, thus it is more moderate on supervision in training. Empirically, it is observed that this strategy is crucial for ConvNet like ResNet but ViT needs intense data augmentations as introduced in [45].

presented in Table 2, which indicate that strong data augmentation does not necessarily improve accuracy if the network is already saturated and may even hurt performance. Consequently, we aim to reduce the intensity of data augmentation while still benefiting from its effects. To achieve this goal, we redesign the data augmentation using a self-mixing approach, as illustrated in Fig. 7.

The specific steps involved in the SelfMix process are illustrated in Fig. 8. Mixing operations are exclusively conducted within each individual image to minimize variations between mixed images and their corresponding mixed soft labels. This straightforward constraint yields significant improvements in the performance of the ResNet backbone and in finetuning distillation.

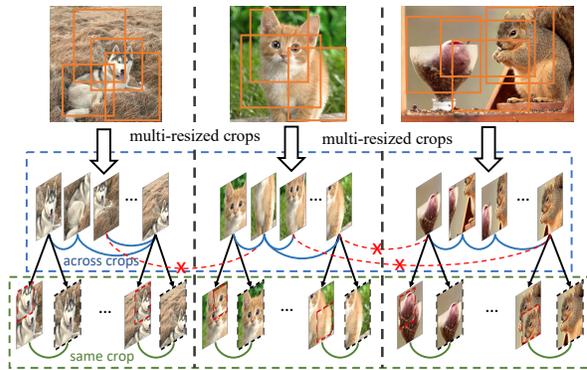


Figure 8: Illustration of the detailed SelfMix augmentation for FerKD. In this strategy, mixture operation only happens within the same image in a mini-batch and cross-image mixing is disabled to preserve stability of soft labels.

Ensemble Supervisions. It has been demonstrated that ensembling more teachers [13, 33] can enhance the performance of distilled students. In this work, we apply this approach to improve the quality of pre-existing soft labels. As

Mixup	CutMix	SelfMix (Ours)	ResNet50	ViT-S/16
✓			78.94	80.32
	✓		79.52	80.95
		✓	80.47	79.29
✓	✓		79.92	81.16
✓	✓	✓	80.26	80.69

Table 2: Top-1 accuracy of data augmentations on different backbones. On ViT, the soft label can benefit from combinations of more mixture operations. The backbones are ResNet-50 and ViT-S/16 and the teacher model is EfficientNet_L2_475. We run three trials and report the means.

Teacher	T _{Top-1} (%)	S _{ResNet50} (%)	S _{ViT-S/16} (%)
Effi_L2_475 [52]	88.14	80.23	81.16
Effi_L2_800 [52]	88.39	80.16	81.30
RegY_128GF_384 [39]	88.24	80.34	81.42
ViT_L16_512 [39]	88.07	80.29	81.43
ViT_H14_518 [39]	88.55	80.18	81.41
BEiT_L_224 [2]	87.52	80.03	81.16
BEiT_L_384 [2]	88.40	80.06	81.11
BEiT_L_512 [2]	88.60	80.09	81.07
ViT_G14_336_30M [10]	89.59	79.03	79.62
ViT_G14_336_CLIP [10]	89.38	79.59	79.48

Table 3: Top-1 accuracy of distillation on ImageNet-1K using a **single** ConvNet or Vision Transformer teacher. Note that Mixup [58] and CutMix [56] are used for training ViT-S/16. We run three trials and report the means.

these soft labels are quantized for efficient storage [37], they must be recovered to their full dimension prior to calibration and averaging for final supervision. Specifically, we define the ensemble soft labels $y_{en} = \frac{1}{M} \sum_{t \in \text{teachers}} \hat{y}_a^t(x)$, where \hat{y}_a^t represents the recovered soft label with calibration from teacher t , x is the input and M is the number of teachers.

4. ImageNet Experiments

We conduct training on the ImageNet-1K (IN1K) [9] training set. We report top-1 validation accuracy of a single 224×224 crop. The default training budget is 300 epochs, and the temperature for both teacher and student is 1.0. The finetuning distillation settings follow their individual designs. All our soft label generation and model training are performed on the A100-GPU High-Flyer cluster with 80GB on each. More details are provided in Appendix.

Baselines: FKD [37] and FunMatch [5]. Regions are randomly sampled in these two approaches. We use ResNet-50 and ViT-S/16 as the backbones for the ablation study. **FKD + Curriculum Sampling:** (i) we sample regions for training from easier ones and gradually increase the level of difficulty. (ii) In contrast, we sample regions from hard ones and gradually decrease the level of difficulty.

4.1. Soft Label from Different Teachers

The soft labels produced by distinct teacher models for the same input image may vary owing to differences in their

calibr. range	[0, 0.2]	[0, 0.3]	[0, 0.4]	[0, 0.5]	[0,0.6]
Top-1	80.31	80.42	80.14	79.84	79.66

Table 4: Ablation of calibration for different probability ranges. The base model is the single teacher FerKD using ResNet-50 without SelfMix data augmentation.

Pre-train	Top-1
vanilla	80.23
+calibration&selfmix	80.68 ^{+0.45}
+multi-teacher ensemble	81.15 ^{+0.47}
+more epochs	81.44 ^{+0.29}

Table 5: Ablation results using ResNet-50 on ImageNet-1K.

distinctive features, architectures, and training strategies. In this section, our objective is to determine which teacher model has the greatest capacity to distill a student. We evaluate two types of student models: ResNet-50 and ViT-S/16. The results are shown in Table 3, where RegY_128GF_384 and ViT_L16_512 achieve the highest student accuracy individually, despite not being the best on their own.

4.2. Ablations

Ablation on Calibration. The results for different calibration ranges are shown in Table 4, it can be observed that [0, 0.3] achieves the best accuracy. In practice, we discard examples in [0, 0.15) and (0.95, 1.0], meanwhile, calibrate examples in [0.15, 0.3] for the final strategy. More visualization of the hardest and easiest regions is shown in Fig. 9. **Ablation on teacher ensemble and training budget.** The results are in Table 5, showing that each design has the consistent improvement. Our final results are shown in Table 8, FerKD performs the best over other SOTA methods.

4.3. Curriculum Distillation

In curriculum distillation, the student model learns from the teachers’ knowledge in a sequential manner, starting from easier examples (high probability regions) and gradually increasing the level of difficulty. This allows the student model to learn from easier to more complex regions. By presenting samples in a curriculum, the student model can learn from the easier samples and build a strong foundation before being exposed to more challenging samples. However, our experimental results, as shown in Table 7, demonstrate that this curriculum learning approach is inferior to our surgical label adaptation strategy, FerKD. Our approach outperforms curriculum distillation by 0.8%. We attribute this performance improvement to the strong ability of soft labels to mitigate overfitting and improve generalization in the initial stages of training. As a result, the advantage of the curriculum learning approach is not as apparent in this scenario. We also notice that “e-to-h” performs slightly better than “h-to-e” demonstrating the effectiveness of curriculum strategy, while both of them are

Effi_L2_475	Effi_L2_800	RegY_128GF_384	ViT_L16_512	ViT_H14_518	BEIT_224	BEIT_384	BEIT_512	Student Acc.
✓	✓							80.23
✓		✓						80.52
✓	✓	✓						80.49
				✓			✓	80.29
			✓	✓			✓	80.53
			✓	✓	✓	✓	✓	80.48
✓			✓					80.35
✓				✓				80.51
✓							✓	80.35
	✓						✓	80.38
	✓			✓				80.51
		✓	✓	✓				80.53
		✓	✓	✓				80.52
✓		✓	✓	✓				80.62
✓		✓	✓	✓				80.74

Table 6: Ablation top-1 accuracy of teacher ensemble on ImageNet-1K with *ConvNet*, *Vision Transformer* or *hybrid teachers*. The left group is ConvNet teachers, the middle is the ViT teachers and the right group is the corresponding student ResNet-50 accuracy. Note that surgical calibration and SelFMix are not used here.

Method	sampling	Top-1
Random (FKD [37])	random	80.2
curriculum distillation	h-to-e	79.7
curriculum distillation	e-to-h	79.9
FerKD (Ours)	surgical	80.7

Table 7: Curriculum and surgical (no ensemble) distillation. “e-to-h” refers to curriculum sampling from easy to hard regions. “h-to-e” refers to sampling from hard to easy.

Model	Epoch	Label	Top-1
Timm [49]	600	Hard _{LS}	80.4
Pytorch (advanced) [44]	600	Hard _{LS}	80.9
MEAL [33]	100 [‡]	Soft	78.2
MEAL V2 [36]	180 [‡]	Soft	80.7
MEAL V2 [36] _{w/ CutMix}	180 [‡]	Soft	81.0
ReLabel [57] _{w/ CutMix}	300	Soft	80.2
FunMatch [5] _{w/ Mixup}	300	Soft	80.5
FKD [37]	300	Soft	80.5
FerKD (Ours)	300	Adap	81.2 ^{+0.7}
FerKD (Ours)	600	Adap	81.4 ^{+0.9}

Table 8: Top-1 accuracy on ImageNet-1K dataset. The backbone network in this table is ResNet-50. [‡] indicates that the model was fine-tuned from hard-label pre-trained weights, resulting in a total training epoch of “300 + ×”.

inferior to the random sampling baseline.

4.4. Finetuning Distillation

Finetuning distillation [36] has been demonstrated as an effective approach to improve the accuracy of knowledge distillation (KD) frameworks. In the case of hard labels, multiple finetuning schemes have been proposed to adapt the model parameters to fit the target dataset, including partial finetuning on selected intermediate layers with varying learning rates [34, 1, 17, 30] or on the last few layers [55]. In this work, we adopt the MEAL V2 [36] protocol by fine-

Pre-train	FT	Top-1
Timm [49]	Hard _{LS}	80.38
Timm [49]	FKD [37]*	80.62
Timm [49]	FerKD* (Ours)	81.06
SWAG [39]	Hard _{LS}	87.22
SWAG [39]	FKD [37] [†]	87.42
SWAG [39]	FerKD [†] (Ours)	87.76
EVA_MIM [10]	Hard _{LS}	89.59
EVA_MIM [10]	FKD [37]	89.67
EVA_MIM [10]	FerKD (Ours)	89.86

Table 9: Fine-tuning distillation results using pre-trained ResNet-50 (Timm), RegNetY-128GF (SWAG) and ViT-G14 (EVA) on ImageNet-1K. [†] We use the same recipe of EVA for SWAG finetuning since SWAG did not provide the complete fine-tuning details. * On ResNet-50, we finetune with 150 epochs for both FKD and FerKD.

tuning the entire network from the pre-trained weights to evaluate the effectiveness of distillation with surgical soft label calibration as the objective. Specifically, we employ three types of pre-trained models:

- (1) Supervised pre-train: ResNet-50 on Timm [49];
- (2) Weakly-supervised pre-train: RegNetY-128GF from SWAG [39];
- (3) Self-supervised pre-train: ViT-G14 from EVA [10];

We verify whether FerKD can continue improving fine-tuning distillation. As shown in Table 9, FerKD achieves consistent improvement across different architectures.

4.5. Single Teacher vs. Teacher Ensemble

The soft labels generated by different teachers for the same input image can be different due to their unique features, architectures, and training strategies. By combining these soft labels from different teachers, the resulting label becomes more informative and can help the student model learn better representations. This approach can be

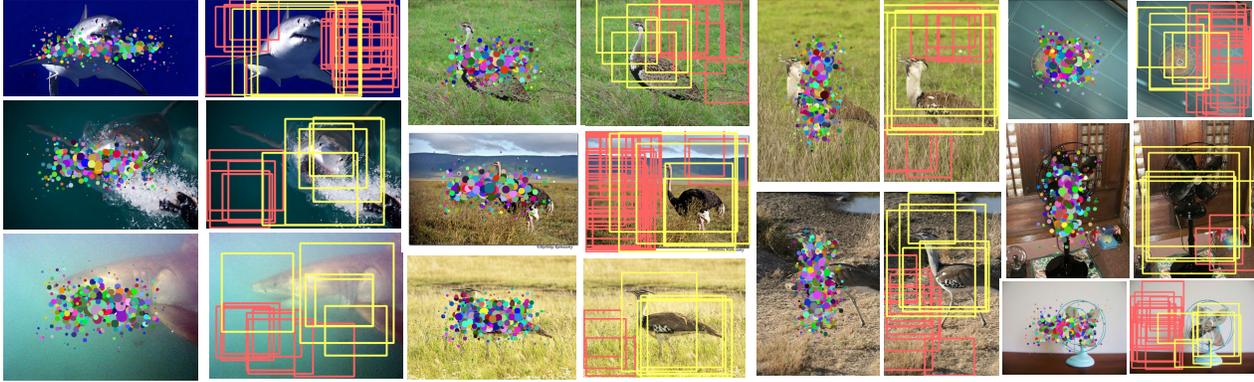


Figure 9: Illustration of the location points from *RandomResizedCrop* and the identified crops by teacher model for **hard** and **easy** samples. Note that we do not involve any localization information, but teacher’s probability can reflect it automatically.

Method	Label	AP ^{box}	AP ^{mask}
ReLabel [57]	LM	39.1	35.2
FKD [37]	Soft	39.7	35.9
FerKD (Ours)	Adap	40.2	36.3

Table 10: COCO object detection and segmentation using a Mask-RCNN with FPN baseline. “LM” represents *label map* during pretraining. The backbone is ResNet-50-300ep.

Method	iNat 2019 ₂₂₄	iNat 2019 ₃₃₆	Places365 ₂₂₄
EVA_MIM [10]	79.9	86.6	61.0
FerKD (Ours)	80.3^{+0.4}	87.1^{+0.5}	61.4^{+0.4}

Table 11: Transfer learning accuracy on various classification datasets. The input sizes are 224×224 and 336×336.

particularly useful in scenarios where the teacher models have complementary strengths or when the input data is challenging and requires multiple perspectives to be accurately labeled. Furthermore, the use of multiple teachers can also help mitigate the effects of overfitting and improve the generalization performance of the student model. The results are shown in Table 6, ensembling four hybrid teachers of *Effi_L2_475*, *RegY_128GF_384*, *ViT_L16_512*, and *ViT_H14_518* achieves the best accuracy.

5. Transfer Learning Experiments

Object detection and segmentation. We conducted evaluations to investigate whether the improvement achieved by FerKD on ImageNet-1K can be transferred to various downstream tasks. Specifically, Table 10 presents the results of object detection and segmentation on COCO dataset [21] using models pre-trained on ImageNet-1K with FerKD. We utilize Mask RCNN [12] with FPN [19] following FKD for the experiment. Our FerKD pre-trained weight consistently outperforms both the baselines ReLabel and FKD on the downstream tasks.

Classification tasks. Table 11 shows the transfer learning result on iNaturalists [14] and Places [60] datasets. On both of these two datasets, our results surpass the baseline EVA pretrained model by significant margins.

Method	IN1K	ReaL	ImageNetV2 Top-images	ImageNetV2 Matched-freq	ImageNetV2 Threshold-0.7
ResNet-50:					
ReLabel [57]	78.9	85.0	80.5	67.3	76.0
FKD [37]	80.1	85.8	81.2	68.2	76.9
FerKD	81.2	86.4^{+0.6}	82.1^{+0.9}	69.5^{+1.3}	77.8^{+0.9}
ViT-G14-336:					
EVA [10]	89.6	90.8	89.0	81.9	86.7
FerKD	89.9	91.3^{+0.5}	89.4^{+0.4}	82.4^{+0.5}	87.1^{+0.4}

Table 12: Results of FerKD on ImageNet ReaL [4] and ImageNetV2 [29] with ResNet-50 and ViT-G14 backbones.

6. Robustness

We provide comparisons on ImageNet ReaL [4] and ImageNetV2 [29] datasets to examine the robustness of FerKD trained models. On ImageNetV2 [29], we verify our FerKD models on three metrics “Top-Images”, “Matched Frequency”, and “Threshold 0.7” following FKD [37]. We perform experiments on two network structures: ResNet-50 and ViT-G14. The results are shown in Table 12, we achieve consistent improvement over ReLabel and FKD on ResNet-50 (224×224) and better accuracy than EVA on ViT-G14 (336×336).

7. Conclusion

In this work, we have presented a new paradigm of *faster knowledge distillation* (FerKD), which employs label adaptation on randomly cropped regions. The proposed method outperforms existing state-of-the-art distillation approaches in terms of both training speed and convergence. Additionally, we make two key observations that could be leveraged in future studies. Firstly, we notice that the most challenging and easiest few crops obtained through the *RandomResizedCrop* operation do not contribute to the model’s learning and can thus be discarded. Secondly, we find that moderate hardness crops can provide crucial context information that helps calibrate the model to learn more robust representations, which in turn benefit downstream tasks.

References

- [1] Armen Aghajanyan, Akshat Shrivastava, Anchit Gupta, Naman Goyal, Luke Zettlemoyer, and Sonal Gupta. Better fine-tuning by reducing representational collapse. *arXiv preprint arXiv:2008.03156*, 2020. 7
- [2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2022. 6
- [3] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, 2009. 3
- [4] Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaoohua Zhai, and Aaron van den Oord. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020. 8
- [5] Lucas Beyer, Xiaoohua Zhai, Amélie Royer, Larisa Markeeva, Rohan Anil, and Alexander Kolesnikov. Knowledge distillation: A good teacher is patient and consistent. In *CVPR*, 2022. 1, 3, 6, 7
- [6] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. Learning efficient object detection models with knowledge distillation. *NeurIPS*, 30, 2017. 1
- [7] Inseop Chung, SeongUk Park, Jangho Kim, and Nojun Kwak. Feature-map-level online adversarial knowledge distillation. In *ICML*, 2020. 3
- [8] Xing Dai, Zeren Jiang, Zhao Wu, Yiping Bao, Zhicheng Wang, Si Liu, and Erjin Zhou. General instance distillation for object detection. In *CVPR*, 2021. 1
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 6
- [10] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. *arXiv preprint arXiv:2211.07636*, 2022. 6, 7, 8, 11
- [11] Jianyuan Guo, Kai Han, Yunhe Wang, Han Wu, Xinghao Chen, Chunjing Xu, and Chang Xu. Distilling object detectors via decoupled features. In *CVPR*, 2021. 1
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 8
- [13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 3, 5
- [14] Grant Van Horn, Oisín Mac Aodha, Yang Song, Alexander Shepard, Hartwig Adam, Pietro Perona, and Serge J. Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018. 8
- [15] Yuenan Hou, Xinge Zhu, Yuexin Ma, Chen Change Loy, and Yikang Li. Point-to-voxel knowledge distillation for lidar semantic segmentation. In *CVPR*, 2022. 1
- [16] Deyi Ji, Haoran Wang, Mingyuan Tao, Jianqiang Huang, Xian-Sheng Hua, and Hongtao Lu. Structural and statistical texture knowledge distillation for semantic segmentation. In *CVPR*, 2022. 1
- [17] Jaejun Lee, Raphael Tang, and Jimmy Lin. What would else do? freezing layers during transformer fine-tuning. *arXiv preprint arXiv:1911.03090*, 2019. 7
- [18] Lujun Li. Self-regulated feature learning via teacher-free feature distillation. In *European Conference on Computer Vision*, pages 347–363. Springer, 2022. 3
- [19] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 8
- [20] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 3
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 8
- [22] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured knowledge distillation for semantic segmentation. In *CVPR*, 2019. 1
- [23] Ilya Loshchilov and Frank Hutter. Online batch selection for faster training of neural networks. *arXiv preprint arXiv:1511.06343*, 2015. 3
- [24] Joel Michael. Where’s the evidence that active learning works? *Advances in physiology education*, 2006. 3
- [25] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? In *NeurIPS*, 2019. 3
- [26] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE symposium on security and privacy (SP)*, 2016. 3
- [27] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *CVPR*, 2019. 3
- [28] Michael Prince. Does active learning work? a review of the research. *Journal of engineering education*, 93(3):223–231, 2004. 3
- [29] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers generalize to imagenet? In *ICML*, 2019. 8
- [30] Youngmin Ro and Jin Young Choi. Autolr: Layer-wise pruning and auto-tuning of learning rates in fine-tuning of deep networks. In *AAAI*, 2021. 7
- [31] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. 3
- [32] Burr Settles. Active learning literature survey. 2009. 3
- [33] Zhiqiang Shen, Zhankui He, and Xiangyang Xue. Meal: Multi-model ensemble via adversarial learning. In *AAAI*, 2019. 3, 5, 7
- [34] Zhiqiang Shen, Zechun Liu, Jie Qin, Marios Savvides, and Kwang-Ting Cheng. Partial is better than all: Revisiting fine-tuning strategy for few-shot learning. In *AAAI*, 2021. 7
- [35] Zhiqiang Shen, Zechun Liu, Dejia Xu, Zitian Chen, Kwang-Ting Cheng, and Marios Savvides. Is label smoothing truly incompatible with knowledge distillation: An empirical study. In *ICLR*, 2021. 3
- [36] Zhiqiang Shen and Marios Savvides. Meal v2: Boosting vanilla resnet-50 to 80%+ top-1 accuracy on imagenet without tricks. *arXiv preprint arXiv:2009.08453*, 2020. 1, 3, 7

- [37] Zhiqiang Shen and Eric Xing. A fast knowledge distillation framework for visual recognition. In *ECCV*, 2022. 1, 2, 3, 4, 6, 7, 8, 11
- [38] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, 2016. 3
- [39] Mannat Singh, Laura Gustafson, Aaron Adcock, Vinicius de Freitas Reis, Bugra Gedik, Raj Prateek Kosaraju, Dhruv Mahajan, Ross Girshick, Piotr Dollár, and Laurens van der Maaten. Revisiting weakly supervised pre-training of visual perception models. In *CVPR*, 2022. 6, 7, 11
- [40] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. *International Journal of Computer Vision*, 130(6):1526–1565, 2022. 3
- [41] Samuel Stanton, Pavel Izmailov, Polina Kirichenko, Alexander A Alemi, and Andrew Gordon Wilson. Does knowledge distillation really work? *arXiv preprint arXiv:2106.05945*, 2021. 3
- [42] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv preprint arXiv:2106.10270*, 2021. 2
- [43] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 1
- [44] PyTorch TorchVision IMAGENET1K_V2. https://pytorch.org/vision/stable/models/generated/torchvision.models.resnet50.html#torchvision.models.ResNet50_Weights. 2022. 7
- [45] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2, 3, 5
- [46] Tao Wang, Li Yuan, Xiaopeng Zhang, and Jiashi Feng. Distilling object detectors with fine-grained feature imitation. In *CVPR*, 2019. 1
- [47] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018. 3
- [48] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4555–4576, 2021. 3
- [49] Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm. *arXiv preprint arXiv:2110.00476*, 2021. 7
- [50] Kan Wu, Jinnian Zhang, Houwen Peng, Mengchen Liu, Bin Xiao, Jianlong Fu, and Lu Yuan. Tinyvit: Fast pretraining distillation for small vision transformers. In *ECCV*, 2022. 3
- [51] Liu Xiaolong, Li Lujun, Li Chao, and Anbang Yao. Norm: Knowledge distillation via n-to-one representation matching. In *ICLR*, 2023. 3
- [52] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *CVPR*, 2020. 1, 6
- [53] Zhendong Yang, Zhe Li, Xiaohu Jiang, Yuan Gong, Zehuan Yuan, Danpei Zhao, and Chun Yuan. Focal and global knowledge distillation for detectors. In *CVPR*, 2022. 1
- [54] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *CVPR*, 2017. 3
- [55] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *NeurIPS*, 27, 2014. 7
- [56] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019. 1, 6
- [57] Sangdoon Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, Junsuk Choe, and Sanghyuk Chun. Re-labeling imagenet: from single to multi-labels, from global to localized labels. In *CVPR*, 2021. 3, 7, 8
- [58] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 1, 6
- [59] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *ICCV*, 2019. 3
- [60] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. *NeurIPS*, 2014. 8

Appendix

In the appendix, we provide more details omitted in the main paper, including:

- Section A: Implementation details.
- Section B: More visualization of identified crops.

Backbone	ResNet-50	ViT-S/16
Epoch	300	300
Batch size	1,024	1,024
Optimizer	AdamW	AdamW
Init. lr	0.002	0.002
lr scheduler	cosine	cosine
Weight decay	0.05	0.05
Warmup epochs	5	5
Num crops	4	4
Label smoothing	\times	\times
Dropout	\times	\times
Stoch. Depth	\times	0.1
Repeated Aug	\times	\times
Gradient Clip.	\times	\times
Rand Augment	\times	\times
Mixup prob.	\times	0.8
Cutmix prob.	\times	1.0
SelfMix prob.	1.0	\times
Random erasing	\times	\times

Table 13: Pre-training setting for ImageNet-1K.

A. Implementation Details

Training details for ResNet-50 and ViT-S/16 in the main text. We elaborate the detailed training settings and hyper-parameters of FerKD for pre-training from scratch on ImageNet-1K with ResNet-50 and ViT-S/16 backbones, as provided in Table 13. Generally, the training protocol follows FKD [37]’s training strategy on ViT, DeiT and SReT. We employ SelfMix for ResNet-50, Mixup and CutMix for ViT-S/16 separately. We also use 4 as the number of crops in each image, batch size = 1,024 during training.

Training details for finetuning ViT-G/14 and RegY-128GF in the main text. The finetuning settings and hyper-parameters of FerKD with ViT-G/14 [10] and RegY-128GF [39] backbones are provided in Table 14, which are similar to the training protocol in EVA [10]. We employ SelfMix for both of the two pretrained backbones.

Data augmentation details for Mixup, Cutmix and SelfMix. The data augmentation configurations adopted in training are: for Mixup, we use probability 0.8 to generate the *Beta distribution*, and 1.0 for CutMix and SelfMix.

B. More Visualization

Fig. 10 illustrates the identified crops by teacher for **hard** and **easy** samples. We do not involve any localization infor-

Backbone	ViT-G/14 [10] RegY-128GF [39]
Peak learning rate	3e-5
Optimizer	AdamW
Optimizer hyper-parameters	$\beta_1, \beta_2, \epsilon = 0.9, 0.999, 1e-8$
Layer-wise lr decay	0.95
Learning rate schedule	cosine decay
Weight decay	0.05
Input resolution	336
Batch size	512
Warmup epochs	2
Training epochs	15
Num crops	2
Drop path	0.4 0.0
Augmentation	RandAug (9, 0.5)
Label smoothing	\times
Cutmix	\times
Mixup	\times
Random erasing	\times
SelfMix prob.	1.0
Random resized crop	(0.08, 1)
Ema	0.9999
Test crop ratio	1.0

Table 14: Fine-tuning setting for ImageNet-1K.

mation, but the teacher’s probability can reflect object and background areas automatically based on their magnitudes.

