WEAK-TO-STRONG ENHANCED VISION MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advancements in large language and vision models have demonstrated extraordinary capabilities, driving researchers to train increasingly larger models in pursuit of even greater performance. However, smaller, easier-to-train models often exist prior to these larger models. In this paper, we explore how to effectively leverage these smaller, weaker models to assist in training larger, stronger models. Specifically, we investigate the concept of weak-to-strong knowledge distillation within vision models, where a weaker model supervises a stronger one, aiming to enhance the latter's performance beyond the limitations of the former. To this end, we introduce a novel, adaptively adjustable loss function that dynamically calibrates the weaker model's supervision based on the discrepancy between soft labels and hard labels. This dynamic adjustment allows the weaker model to provide more effective guidance during training. Our comprehensive experiments span various scenarios, including few-shot learning, transfer learning, noisy label learning, and common knowledge distillation settings. The results are compelling: our approach not only surpasses benchmarks set by strong-to-strong distillation but also exceeds the performance of fine-tuning strong models on full datasets. These findings highlight the significant potential of weak-to-strong distillation, demonstrating its ability to substantially enhance vision model performance. Code will be released.

031

025

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

"Big things have small beginnings." — Movie "Prometheus".

This adage aptly encapsulates the developmental journey of high-performance models in the fields
 of computer vision and natural language processing. The remarkable models that currently drive
 advancements in these areas did not appear out of nowhere; they evolved incrementally from simpler,
 less powerful architectures.

037 In the realm of NLP, the journey began with models like RNN and LSTM networks. These early models laid the foundation for more advanced architectures, gradually evolving into models like GPT (Radford et al., 2019; Brown et al., 2020). With its 175 billion parameters, GPT-3 showcased 040 the transformative power of scaling up, ultimately paving the way for today's state-of-the-art large 041 language models that excel in various tasks, from translation to creative writing. Similarly, the evo-042 lution in vision began with the pioneering LeNet (LeCun et al., 1998) architecture, designed for digit recognition. ResNet (He et al., 2016) then addressed the vanishing gradient problem, allowing for 043 much deeper networks. Today, large vision models like ViTs (Dosovitskiy et al., 2020) continue to 044 push the boundaries, achieving unprecedented performance surpassing human across various visual tasks. 046

As demonstrated by empirical studies on scaling laws (Kaplan et al., 2020), model performance
typically scales with model size, dataset size, and the amount of compute used for training. This
suggests that training larger models with more data holds the greatest potential for improvement,
while other factors like training recipes or network architectures have relatively minimal impact
across a wide range. However, before training a new large model, there often exists a smaller,
weaker model. It's natural to ask whether these existing weaker models can be leveraged to assist in
training larger ones. In this paper, we focus on addressing the challenge of how to efficiently utilize
these weaker models to optimize and guide the training of more powerful, larger models.

Building on the notion of lever-055 aging existing models, previous 056 work has introduced the con-057 cept of "superalignment" to ad-058 dress the challenge of incorporating human expertise into the supervision of superhuman AI 060 models. This approach seeks 061 to align powerful models with 062 human input to maximize their 063 learning potential. A particu-064 larly relevant study in this con-065 text is Weak-to-Strong Gener-066 alization (Burns et al., 2023), 067 which explores the intriguing 068 possibility of using weaker mod-069 els to supervise stronger ones. The findings are compelling: de-070 spite their inherent limitations, 071 weaker models can provide su-072 pervision that enables stronger 073 models — already equipped with 074 superior generalization and rep-075 resentational power - to surpass 076 their weaker counterparts. Re-077 markably, even when the weaker models offer incomplete or noisy 079 labels, the stronger models are able to transcend these shortcomings, achieving higher per-081

085

087

090 091 092

094

095

096

097



Figure 1: Our proposed AdaptConf achieves the best performance on a broad range of tasks compared with other knowledge distillation based methods. The corresponding values are calculated by averaging results on each task, *i.e.*, classification, transfer learning, few-shot learning, and learning with noisy labels. CLS-CIFAR-S (same model family): Table 2, CLS-CIFAR-D (different model family): Table 4a, CLS-ImageNet-S: Table 3, CLS-ImageNet-D: Table 3, TL-ImageNet: Table 7a, TL-iNat: Table 7b, FSL-miniImageNet: Table 5, LNL-CIFAR: Table 8.

formance. This concept has shown its efficacy in fields such as natural language processing and
 reinforcement learning, affirming the potential of Weak-to-Strong knowledge distillation as a viable
 and effective strategy.

	Modal	CIFAR-100 validation set			Madal	ImageNet validation set				
MO	widdei	#Params	Top-1 (%)	Δ (%)	Win (%)	Widdei	#Params	Top-1 (%)	$\Delta(\%)$	Win (%)
	MobileNet-V2	0.8M	66.9	-	-	ResNet-18	11.7M	69.8	-	-
	ResNet-56	0.9M	72.9	-6.0	8.4	ResNet-34	21.8M	73.5	-3.7	4.8
	ResNet-110	1.7M	74.8	-7.9	7.4	ResNet-50	25.6M	76.2	-6.4	3.9
	VGG-13	9.5M	75.3	-8.4	6.3	DeiT-S	22M	79.9	-10.1	3.4
	ResNet32×4	7.4M	79.9	-13.0	4.9	DeiT-B	86M	81.8	-12.0	3.1
	ViT-B $_{\uparrow 224}$	86M	89.0	-22.1	2.5	DeiT-B _{↑384}	86M	83.0	-13.2	2.8

Table 1: Comparison between models on CIFAR-100 and ImageNet. Δ represents the performance gap between the baseline MobileNet-V2 / ResNet-18 and other stronger models. "Win" indicates the ratio of samples correctly classified by MobileNet-V2 / ResNet-18 but incorrectly classified by the other stronger models.

098 We delve into the benefits brought via "vision superalignment", specifically investigating the appli-099 cability of Weak-to-Strong enhancement (W2S) within the context of vision tasks. Take the image 100 classification task as an example, models typically develop from small to large, and the cost of train-101 ing a small (weak) model is far less than that of training a large (strong) model. As shown in Table 1, 102 larger models generally perform better. However, even when a model is 100 times larger and the 103 top-1 accuracy is 22% higher, there are still many samples correctly identified by the weak model 104 but misclassified by the strong model. This indicates that there are always opportunities to boost 105 performance by leveraging the weak model. Therefore, when attempting to train a large, strong model, a natural question arises: how can we leverage existing weak models to achieve further per-106 formance gains? Our study meticulously designs and examines multiple scenarios in computer vi-107 sion, including few-shot learning, transfer learning, noisy label learning, and traditional knowledge

2

distillation settings. In these scenarios, stronger models are trained to learn from weaker models.
Through detailed validation and comparative experiments, we demonstrate the feasibility of W2S in
the visual domain. Furthermore, we introduce an improved and adaptive confidence scheme to enhance the efficacy of W2S. Our work validates the concept of weak-to-strong boosting in computer
vision, representing a significant advancement in understanding and optimizing the interaction between strong and weak models. This approach has the potential to pave the way for groundbreaking
advancements in achieving human-level expertise and even superhuman artificial intelligence.

115 116

117

2 RELATED WORKS

118 The pursuit of enhancing the performance of deep neural networks in computer vision has led to 119 the development of the teacher-student learning paradigm (Hinton et al., 2015; Roth et al., 2023). 120 This approach typically involves a stronger model (teacher) improving the performance of a weaker 121 model (student), with extensive research focusing on optimizing the capabilities of the weaker model. Various strategies have been proposed to achieve this. For instance, (Romero et al., 2014) 122 suggests that in addition to the output logits, incorporating intermediate layer features for supervi-123 sion can significantly boost the student's learning. (Park et al., 2019) posits that the relationships 124 between samples can serve as valuable supervisory information. 125

In a further refinement of this approach, (Zhao et al., 2022) redefines classical knowledge distillation (KD) loss, segmenting it into target-class and non-target-class distillation to balance the transfer of these two types of information more effectively. (Heo et al., 2019) delves into the details and components of feature distillation, arriving at an improved method for the transfer of feature knowledge. Meanwhile, (Chen et al., 2021a) explores cross-stage feature transfer as an alternative to the conventional same-stage feature transfer. These methods have proven effective for strong-to-weak generalization scenarios.

133 However, with the gradual increase in the size and complexity of vision foundation models, the focus has shifted towards weak-to-strong boosting, *i.e.*, how a weak model can improve a strong model. 134 In this context, (Furlanello et al., 2018) investigates knowledge distillation between teachers and 135 students of equal size, demonstrating the feasibility of distilling models of the same size. Building 136 upon this, (Xie et al., 2020) introduces the use of additional unlabeled data for knowledge distilla-137 tion among models of equal size, further validating the effectiveness of strong-to-strong boosting, 138 especially in scenarios with abundant data availability. This body of research sets the stage for our 139 exploration into weak-to-strong boosting, a relatively uncharted yet promising domain in the field 140 of vision foundation models. 141

Part of our experimental settings are similar to weakly supervised learning (Durand et al., 2017; Joulin et al., 2016), where there are no annotations (ground truth labels) in training machine learning models. However, unlike weak supervision, which focuses on obtaining a large amount of annotated data at a low cost, we are more interested in the weak-to-strong boosting process itself rather than the availability of annotations.

146 147

148 149

3 WEAK-TO-STRONG ENHANCEMENT

To advance towards super-human AGI models, a weak-to-strong approach is critical. This means us-150 ing human-level intelligence as a foundation to guide and refine the development of more advanced, 151 super-human systems. We begin by focusing on foundational tasks and models that can progres-152 sively support the growth of stronger architectures. In the following, we examine the feasibility 153 of *weak-to-strong enhancement*, where simpler, weaker models provide useful supervision to more 154 complex models. This step-by-step approach ensures that as models evolve, they can benefit from 155 earlier stages, steadily improving in both capability and accuracy. One of the key challenges in this 156 weak-to-strong framework is dealing with the noisy or incomplete supervision signals provided by 157 weaker models. To mitigate this, we introduce a novel technique: *adaptive confidence distillation*. 158 This method leverages the insights from weaker models while dynamically adjusting the level of 159 trust placed in them. By modulating the influence of weak supervision based on its confidence, our approach ensures that the stronger model can benefit from imperfect labels without being misled. 160 This adaptive mechanism allows stronger models to distill meaningful knowledge, even from less 161 accurate outputs.

162 3.1 SELECTION OF VISION MODEL

164 In our exploration of weak-to-strong enhancement for vision models, it is essential to first define 165 which models are suitable for this foundational research. Several categories of models can serve as strong candidates for vision foundation models, including text-vision models (Radford et al., 166 2021), image generation models (Rombach et al., 2022; Chen et al., 2020), and general zero-shot 167 models (Bai et al., 2023; Kirillov et al., 2023). Each of these models brings unique strengths and 168 approaches to solving computer vision tasks. To identify the most appropriate models for weak-tostrong boosting, we propose a definition that emphasizes both versatility and effectiveness. Vision 170 foundation models should be capable of addressing a wide range of visual tasks while consistently 171 delivering high-quality performance. 172

Based on these criteria, we propose that backbones pretrained on ImageNet are strong candidates for 173 vision foundation models. The rationale behind this choice is threefold. First, ImageNet-pretrained 174 backbones have consistently proven to be highly adaptable and effective across downstream vision 175 tasks such as classification, detection, segmentation, tracking, colorization, etc.. Fine-tuning these 176 backbones often results in state-of-the-art performance, underscoring their robustness and versatil-177 ity. Second, there is a wealth of pretraining algorithms specifically developed for these models (He 178 et al., 2022a; Xie et al., 2022), positioning them as universal tools for a variety of vision tasks. Fur-179 thermore, these models frequently serve as one branch in vision-language multimodal models (Du et al., 2022), further validating their applicability in cross-modal tasks. Finally, compared to other 181 foundation models, such as CLIP (Radford et al., 2021) or Diffusion-based (Rombach et al., 2022) 182 models trained on massive web-scale datasets, ImageNet-trained backbones are far more accessible. Their moderate computational demands make them a practical choice for a broader range of 183 researchers with limited resources. 184

To validate the feasibility of weak-to-strong enhancement, we focus on these pretrained backbones and use image classification as the fundamental task. By selecting this well-established task, we create a controlled environment to rigorously test the effectiveness of our adaptive confidence distillation method. This will provide a solid baseline for expanding the weak-to-strong paradigm to more complex vision tasks and models in future work.

190 191

197

3.2 Adaptive Confidence Distillation

In this section, we explore the methodology for implementing weak-to-strong boosting in vision foundation models. The central question we address is how a weak vision foundation model can supervise a stronger counterpart effectively. (Burns et al., 2023) proposes an augmented confidence loss approach, which is formulated as:

$$L_{\text{conf}}(f) = (1 - \alpha) \operatorname{CE}(f(x), f_w(x)) + \alpha \operatorname{CE}(f(x), \hat{f}(x)), \tag{1}$$

where f represent the strong model that needs to be optimized, and f_w denote the weak model, $\hat{f}(x)$ refers to the hard label predicted by the strong model for an input image x. The loss function incorporates the cross-entropy loss (CE) and is balanced by a hyperparameter α . In this formulation, the first term of the loss function resembles the traditional knowledge distillation loss, signifying the learning process of the strong model from the weak model. Given that the labels provided by the weak model may not always be accurate, the second term of the loss function encourages the strong model to leverage its superior generalization abilities and prior knowledge to refine its predictions.

The strength of this approach lies in its ability to balance direct learning from the weak model with the strong model's intrinsic capacity for understanding and interpreting the visual data. This method paves the way for the strong model to surpass the limitations of the weak model, utilizing the latter's guidance while simultaneously enhancing its predictions through its advanced capabilities.

Addressing the limitations inherent in the supervision provided by weak models and the inaccuracies of strong models' self-generated hard labels, a more sophisticated approach is required beyond a simple weighted combination of these labels. Given the challenge in directly discerning the accuracy of each label, leveraging confidence as a metric for selecting the most probable correct label emerges as a viable solution.

215 We propose to use the discrepancy between the soft label and the hard label as an indicator of the model's confidence. The underlying rationale is that when a model's soft label closely aligns with

216		RecNet20	PecNet32	PosNot8×1	WPN 16 2	WPN 40-1	VGG8
217	Teacher	68.93	71 72	72 41	72 71	72 30	71 99
218		ResNet56	ResNet110	ResNet 32×4	WRN-40-2	WRN-40-2	VGG13
219	Student	72.94	74.80	79.90	77.20	77.20	75.26
220	KD (Hinton et al., 2015)	73.81	76.45	79.32	78.25	77.97	76.41
201	FitNet (Romero et al., 2014)	70.51	73.15	77.65	76.71	76.12	76.39
221	RKD (Park et al., 2019)	72.98	75.62	80.10	77.27	77.76	76.20
222	ReviewKD (Chen et al., 2021a)	70.15	72.30	77.22	75.86	75.78	74.22
223	DKD (Zhao et al., 2022)	73.90	76.57	79.52	78.18	77.95	76.62
224	AugConf (Burns et al., 2023)	73.86	76.72	80.34	78.34	78.15	76.55
225	AdaptConf (Ours)	74.17	76.86	80.64	78.58	78.40	76.84
226	Δ	+1.23	+2.06	+0.74	+1.38	+1.20	+1.58

Table 2: Results on the CIFAR-100 validation set. Teachers and students are in the same architectures. And Δ represents the performance improvement over the student model trained from scratch. All results are the average over 3 trials.

its hard label, it suggests a higher confidence in its own judgment. To capitalize on this insight, we introduce an adaptive confidence loss that dynamically adjusts based on the model's confidence level. The specific formulation of this loss is as follows:

$$L_{AC}(f) = (1 - \beta(x))CE(f(x), f_w(x)) + \beta(x)CE(f(x), \hat{f}(x)),$$

$$\beta(x) = \frac{\exp(CE(f(x), \hat{f}(x)))}{\exp(CE(f(x), \hat{f}(x))) + \exp(CE(f(x), \hat{f}_w(x)))}.$$
(2)

In this formula, $\beta(x)$ is a function of the input image x that calculates the confidence weight and $\hat{f}_w(x)$ is the hard label of x in the weak model. This weight determines the balance between learning from the weak model and relying on the strong model's own predictions. The cross-entropy loss (CE) is used for both components, with the first term focusing on learning from the weak model and the second term emphasizing the strong model's self-supervision.

This adaptive confidence loss enables a more nuanced approach to weak-to-strong boosting. By adjusting the weight based on confidence levels, it allows the strong model to discern when to prioritize its own predictions over the guidance of the weak model and vice versa. This adaptability is key to overcoming the inaccuracies and limitations of both models, leading to more effective learning and enhanced performance in vision foundation models.

248 249 250

251

252

253 254

255

227

228

229

230

4 EXPERIMENT

In this section, we report our main empirical results on various tasks, including baselines and promising methods. All implementation details are attached in supplementary materials.

4.1 TASKS

Image Classification. Our experiments 256 are primarily focused on two benchmark 257 datasets. CIFAR-100 (Krizhevsky et al., 258 2009) is a widely recognized dataset for 259 image classification, comprising 32×32 260 pixel images across 100 categories, with 261 training and validation sets containing 262 50,000 and 10,000 images, respectively. 263 Conversely, ImageNet (Deng et al., 2009) 264 is a large-scale dataset for classification 265 tasks, encompassing 1.28 million train-266 ing images and 50,000 validation images 267 across 1,000 classes. Additionally, we explore scenarios where only soft labels 268 generated by a weak teacher model are 269 available for training.

Taachar	ResNet18	MobileNet-V1
Teacher	69.75	71.57
Student	ResNet34	ResNet50
Student	73.47	76.22
KD (Hinton et al., 2015)	73.68	76.52
FitNet (Romero et al., 2014)	70.93	73.61
RKD (Park et al., 2019)	73.65	76.45
ReviewKD (Chen et al., 2021a)	72.99	75.28
DKD (Zhao et al., 2022)	73.74	76.72
AugConf (Burns et al., 2023)	<u>73.80</u>	76.64
AdaptConf (Ours)	74.16	76.94
Δ	+0.69	+0.72

Table 3: Top-1 results on the ImageNet validation set. Δ represents the performance improvement over the student model trained from scratch.

270	Tanahar	ShuffleNet-V1	ShuffleNet-V1	MobileNet-V2	MobileNet-V2	ShuffleNet-V2
271	Teacher	72.40	72.40	66.85	66.85	74.44
272	Student	ResNet32×4	WRN-40-2	VGG13	ResNet50	ResNet32 \times 4
212		79.90	77.20	75.26	80.43	79.90
273	KD (Hinton et al., 2015)	80.19	78.02	75.39	78.64	80.31
274	FitNet (Romero et al., 2014)	77.61	75.15	72.36	75.92	78.05
214	RKD (Park et al., 2019)	80.30	77.23	76.21	79.89	80.39
275	ReviewKD (Chen et al., 2021a)	78.43	75.98	73.69	77.05	77.84
276	DKD (Zhao et al., 2022)	80.55	<u>78.10</u>	75.81	79.65	80.67
077	AugConf (Burns et al., 2023)	80.62	77.92	<u>76.43</u>	<u>80.75</u>	80.84
277	AdaptConf (Ours)	80.99	78.55	76.58	80.98	81.06
278	Δ	+1.09	+1.35	+1.32	+0.55	+1.16

(a) Trained with teacher's prediction and GT label. Δ is the improvement over the student trained from scratch.

Teacher Student	ShuffleNet-V1 72.40 ResNet32×4	ShuffleNet-V1 72.40 WRN-40-2	MobileNet-V2 66.85 VGG13	MobileNet-V2 66.85 ResNet50	ShuffleNet-V2 74.44 ResNet32×4
KD (Hinton et al., 2015)	77.92	<u>76.45</u>	72.13	73.32	78.27
FitNet (Romero et al., 2014)	75.74	74.03	70.57	71.45	76.42
RKD (Park et al., 2019)	76.59	75.70	70.28	72.06	77.84
AugConf (Burns et al., 2023)	78.25	76.37	<u>72.51</u>	74.48	<u>78.81</u>
AdaptConf (Ours)	78.48	76.66	72.93	74.67	79.04
$\overline{\Delta}$	+6.08	+4.26	+6.08	+7.82	+4.37

(b) Trained with teacher's prediction only. Δ represents the performance improvement over the teacher model.

Table 4: Results on the CIFAR-100 validation set. Teachers and students are in the different architectures. All results are the average over 3 trials.

291 292

279

280 281

283 284

286 287 288

289

290

293

294 Few-shot learning. We explore few-shot learning across the miniImageNet (Vinyals et al., 2016) 295 dataset which contains 100 classes sampled from ILSVRC-2012 (Russakovsky et al., 2015). We 296 randomly split the dataset into 64, 16, and 20 classes as training, validation, and testing sets, respec-297 tively. And ensure that each class has 600 images of 84×84 image size. We utilize the ResNet36 to 298 explore the weak-to-strong boosting performance in few-shot task. To demonstrate weak-to-strong 299 boosting performance, we follow Meta-Baseline and conduct related experiments on classifier stage and meta stage. 300

301 **Transfer learning.** We explore transfer learning across two benchmark datasets: ImageNet (Deng 302 et al., 2009), and iNaturalist 2018 (Van Horn et al., 2018), the latter comprising 437,513 training 303 images and 24,426 test images distributed across 8,142 species. We utilize the ViT-B (Dosovitskiy 304 et al., 2020) that has been pretrained on the ImageNet training set using the self-supervised MAE (He et al., 2022b) approach, leveraging only image data without labels. Our results are reported for 305 the fine-tuning phase, which is conducted under the guidance of a weak teacher model on each 306 benchmark. Furthermore, we investigate scenarios where only soft labels produced by the weak 307 teacher model are used for training. 308

309 Learning with noisy labels. We eval-310 uate our approach using two datasets with simulated label noise, specifically 311 CIFAR-10 (Krizhevsky et al., 2009) 312 and CIFAR-100 (Krizhevsky et al., 313 2009). Consistent with prior re-314 search (Li et al., 2020; Tanaka et al., 315 2018), we introduce two distinct types 316 of simulated noisy labels: symmetric 317 and asymmetric. Symmetric noise is in-318 troduced by randomly substituting the 319 labels of a certain proportion of the 320 training data with other possible la-321 bels uniformly. In contrast, asymmetric noise involves systematic mislabel-322 ing to mimic real-world errors, such as 323

Tanahar	ResN	let12	ResNet18		
Teacher	59.65	77.80	60.83	78.96	
Student	ResNet36		ResNet36		
Student	60.91	79.01	60.91	79.01	
	1-shot	5-shot	1-shot	5-shot	
KD	60.94	79.14	61.57	<u>79.79</u>	
RKD (Park et al., 2019)	59.74	78.30	60.80	78.82	
AugConf (Burns et al., 2023)	<u>61.38</u>	<u>79.33</u>	<u>61.66</u>	79.46	
AdaptConf (Ours)	61.50	79.52	62.29	79.96	
Δ	+2.59	+2.67	+3.38	+3.11	

Table 5: Average 5-way accuracy (%) with 95% confidence interval on the miniImageNet validation set in Classification **Training stage.** Δ represents the performance improvement over the student model trained from scratch. All results are the average over 3 trials.

flipping the labels to closely related classes. For example, in CIFAR-10, truck is mislabeled as au-

		Class-stage		Meta-s	stage
324	Tarahan	ResNet12	ResNet18	ResNet12	ResNet18
325	Teacher	59.20	60.63	65.26	66.51
326	Student	ResNet36	ResNet36	ResNet36	ResNet36
327	Studelit	65.08	65.08	65.08	65.08
328	KD (Hinton et al., 2015)	63.43	65.04	<u>66.08</u>	<u>65.93</u>
020	RKD (Park et al., 2019)	64.79	65.42	65.96	65.46
329	AugConf (Burns et al., 2023)	65.15	65.59	65.9	65.78
330	AdaptConf (Ours)	65.38	65.74	textbf66.08	65.95
331	$\overline{\Delta}$	+0.30	+0.66	+1.00	+0.87

Table 6: Average 5-way accuracy on miniImageNet validation set at Meta-Learning stage. Δ represents the performance improvement over the student trained from scratch. All results are the average over 3 trials.

tomobile, bird as *airplane*, and *cat* is interchanged with *dog*. For CIFAR-100, similar mislabeling is applied within each of the super-classes in a circular fashion.

Baseline methods. The predominant framework for implementing teacher-student training paradigms is knowledge distillation (Hinton et al., 2015). This approach outlines a method where a larger, more complex teacher network guides the training of a more compact student network. Nonetheless, inspired by the findings of Burns *et al.* (Burns et al., 2023), our work pivots towards a scenario where the student network surpasses the teacher in visual capabilities. Despite this inversion of roles, there remains valuable dark knowledge in the teacher that can be transferred to the student, either through logits or via intermediate representational features. To benchmark our experiments, we employ a range of established (Hinton et al., 2015; Romero et al., 2014; Park et al., 2019; Heo et al., 2019; Chen et al., 2021a; Hao et al., 2023a) and recently proposed (Zhao et al., 2022; Burns et al., 2023) distillation techniques as baseline methods.

Teacher: ResNet50 (80.36)

350

351 352

332

333

334 335

336

337

338

339

340

341

342

343

344

345

4.2 MAIN RESULTS

4.2.1 IMAGE CLASSIFICATION.

CIFAR-100 image classification. We 353 commence our investigation with an 354 exploration of weak-to-strong boosting 355 (W2S) on the CIFAR-100 dataset. The 356 outcomes of this investigation are de-357 lineated in Tables 2 and 4. Specifi-358 cally, Table 2 presents the scenarios in which both teacher and student models 359 share the same network architectures. 360 We examine a range of prevalent vi-361 sion architectures such as ResNet (He 362 et al., 2016), WRN (Zagoruyko & Ko-363 modakis, 2016), and VGG (Simonyan 364 & Zisserman, 2014). We employ various KD methods to assess the poten-366 tial of larger-capacity students guided 367 by limited-capacity teachers. Remark-368 ably, in nearly all cases employing KD-369 based approaches, the student models outperform those trained from scratch. 370

Furthermore, both AugConf (Burns
et al., 2023) and our proposed AdaptConf method surpasses all previous dis-

Student: ViT-B (MAE pretrain) 83.53 83 62 82.32 KD (Hinton et al., 2015) FitNet (Romero et al., 2014) 82.48 81.02 RKD (Park et al., 2019) 82.19 80.98 DKD (Zhao et al., 2022) 83.68 82.38 AugConf (Burns et al., 2023) 83.70 83.86 82.51 AdaptConf (Ours) +0.33+2.15 Δ (a) Top-1 results on the ImageNet validation set. Teacher: ResNet101 (67.42) Teacher + GT Teacher Student: ViT-B (MAE pretrain) 75.28 KD (Hinton et al., 2015) 75.60 71.57 FitNet (Romero et al., 2014) 73.68 70.11 DKD (Zhao et al., 2022) 75.82 71.73 AugConf (Burns et al., 2023) 75.90 AdaptConf (Ours) 76.03 71.99 +0.75Λ +4.57

Teacher + GT

Teacher

(b) Top-1 results on the iNaturalist 2019 test set.



tillation techniques across all teacher-student pairs. This highlights that simply emulating a weak
teacher does not yield the most favorable outcomes. Notably, AdaptConf consistently achieves superior performance compared to AugConf (Burns et al., 2023), underscoring the advantage of our
dynamic adaptive confidence weighting. This approach provides a more refined mechanism for
facilitating weak-to-strong knowledge transfer.

378	dataset	CIFAR-10				CIFAR-100			
379	noise type	asymmetric		symmetric		asymmetric		symmetric	
380	Teacher	PR18		PR18		PR18		PR18	
381	Teacher	92.98	99.56	95.80	99.80	73.20	92.67	76.16	92.90
000	Student	PR34		PR34		PR34		PR34	
382	Student	93.69	99.61	96.13	99.77	74.80	92.94	78.20	93.77
383		Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
384	KD (Hinton et al., 2015)	<u>93.54</u>	<u>99.84</u>	95.90	99.84	75.49	93.67	77.61	93.74
385	RKD (Park et al., 2019)	92.42	99.75	<u>95.99</u>	<u>99.85</u>	74.20	93.54	76.92	93.09
386	AugConf (Burns et al., 2023)	92.60	99.75	95.10	99.83	74.99	<u>93.72</u>	<u>78.34</u>	<u>94.02</u>
387	AdaptConf (Ours)	93.69	99.84	textbf96.13	99.87	75.61	93.78	78.64	94.03
388	Δ	+0.00	+0.23	+0.00	+0.10	+0.81	+0.84	+0.44	+0.26

389 390 391

Table 8: Top-1 and top-5 results on the CIFAR-10/CIFAR-100 noise label validation set. Δ represents the performance improvement over the student model trained from scratch. All results are the average over 3 trials.

392 Table 4 shows the results of teacher-student pairs from different series, such as ShuffleNet (Zhang et al., 2018) and MobileNet (Sandler et al., 2018). Additionally, take the MobileNetV2-ResNet50 393 pair as an example, the experimental results reveal that when the teacher model is significantly 394 weaker, *i.e.*, a substantial performance gap exists between the weak teacher model and the strong 395 student model, none of the KD-based methods were able to effectively enhance the strong student's 396 performance, except for AugConf and AdaptConf. The possible reason is that these methods include 397 the predictions of the strong student in the loss function. This proves that self-training methods, akin 398 to those described in (Lee et al., 2013), can mitigate the bias from a suboptimal teacher model. It 399 is important to note that FitNet (Romero et al., 2014) consistently underperforms when compared 400 to training from scratch. This could be attributed to its sole focus on intermediate features, which 401 may be more misleading for the strong student to learn from than soft predictions, as suggested by 402 (Hao et al., 2023b). Overall, our AdaptConf achieves an improvement of 0.5%-2% on all evaluated teacher-student pairings, whether they are from the same or different series. 403

Furthermore, we investigate a scenario where only the teacher's output is available, as shown in
Table 4b. In this context, it becomes evident that AugConf and AdaptConf yields more significant
improvements compared to other KD-based methods when ground truth is absent. This observation
underscores the suitability of our confidence distillation approach for more extreme W2S scenarios
where ground truth is not available.

ImageNet image classification. Table 3 presents the top-1 accuracy for image classification on the ImageNet dataset. Our AdaptConf method achieves significant improvements across both W2S scenarios, whether employing the same or different architectures.

413 414

4.2.2 Few-shot learning

415 For the few-shot learning task, we conduct distillation experiments separately in the classification 416 (Table 5) and meta-learning (Table 6) stages. We compare and evaluate the performances of student when trained with teachers of different sizes. In the classification experiments, only RKD results in 417 a performance degradation of the student model, while the usage of other methods led to varying 418 degrees of improvement. Notably, our confidence-based method outperforms previous knowledge 419 distillation based ones. In the meta-learning stage, we employ weights from different training stages 420 of the same model as the teacher. Experimental results demonstrate significant advantages of our 421 proposed method. Even when using the Class-stage weight as the teacher, our approach achieves a 422 +0.66% improvement over the baseline set by a weaker ResNet18 (Class-stage) teacher model. Fur-423 thermore, when using the same stage weight as the teacher, our confidence-based method surpasses 424 previous knowledge distillation results to a greater extent.

425 426

427

4.2.3 TRANSFER LEARNING

Table 7 examines the efficacy of transfer learning using the iNaturalist (Van Horn et al., 2018) and ImageNet (Deng et al., 2009) datasets. When our method is trained with ground truth labels on ImageNet, it demonstrates a notable enhancement, achieving an increase of +0.33% in top-1 accuracy on a model with a high precision of 83.5%. Even without ground truth labels, our approach still secures a +2.15% improvement over the baseline set by a weaker ResNet50 teacher model. On



Figure 2: Ablation study examining the impact of hyper-parameter variation on confidence distillation results. The parameter α for AugConf is adjusted across a range from 0.1 to 0.9, while the temperature T for AdaptConf is scaled from 0.1 to 8.



Figure 3: Quantitative analysis about the value of $\beta(x)$ in Eq. 2 on the CIFAR-100 dataset. The evaluation is based on the ShuffleNetV1-ResNet32x4 teacher-student architecture pair.

the iNaturalist dataset, our confidence-based method also surpasses previous knowledge distillation results by a considerable margin.

4.2.4 LEARNING WITH NOISY LABELS

In Table 8, we analyze the effectiveness of weak-to-strong using the CIFAR-10 and CIFAR-100 datasets under two simulated noisy label settings. When training the model on the sample dataset (CIFAR-10), all methods except ours, negatively impact the model given its already high accuracy. This underscores the robustness of our method, irrespective of the performance gap between the teacher and student models. On the CIFAR-100 dataset, our method demonstrates a performance improvement of 0.81% in top-1 accuracy under the asymmetric noise type setting.

4.3 ABLATION STUDY

Robustness of confidence distillation. In this study, we investigate the necessity of devising a method that goes beyond a mere weighted combination of labels. As depicted in Eq. 1, despite its straightforward approach of integrating direct learning from a weaker model with the intrinsic capacity of a stronger model, AugConf (Burns et al., 2023) still requires manual tuning of a hyper-parameter α to balance the ratio of two different objectives. The setting of different α values can have varying impacts across different contexts. Similarly, although our proposed AdaptConf does not require manual adjustment of α to balance the proportions of objectives, we can manipulate the temperature T to control the degree of probability distribution in soft labels during the computa-tion of the cross-entropy $CE(\cdot)$, following a conventional distillation method (Hinton et al., 2015). Therefore, we explore the effects of these two methods under different hyper-parameter settings on

486 the final outcome. Overall, the performance of KD, AugConf, and AdaptConf improves sequentially 487 across various architectural settings. Moreover, it can be observed that AugConf exhibits a larger 488 fluctuation in results compared to AdaptConf, indicating that the influence of α on AugConf is more 489 significant than the effect of T on AdaptConf, which suggests that our AdaptConf has superior ro-490 bustness. Additionally, the average outcomes achieved by AdaptConf are consistently higher than those of AugConf under different hyper-parameter settings. 491

492 **Robustness of confidence distillation.** In this section, we perform a quantitative analysis of the 493 confidence weight determined by our dynamic function $\beta(x)$ as delineated in Eq. 2, with the findings 494 illustrated in Figure 3. We selected checkpoints from four distinct training phases and calculated 495 their specific $\beta(x)$ values on the validation set. It can be observed that as training progresses, the 496 proportion of samples with $\beta = 0.5$ increases, indicating that the student model's performance is improving and being aligned with the weak teacher's correct classifications. A higher temperature 497 setting T reduces the cross-entropy (CE) discrepancy between the teacher and student, promoting a 498 more uniform balance between the weak teacher's guidance and the strong student's own predictions. 499 Consequently, the number of samples where $\beta = 0.5$ also increases with training. These phenomena 500 collectively validate that our proposed AdaptConf can dynamically adjust the learning ratio between 501 the two components. 502

503

5 CONCLUSION

504 505

506 In this paper, we investigate weak-to-strong boosting for vision foundation models and unveil a promising avenue for enhancing the capabilities of artificial intelligence in the visual domain. By 507 leveraging an innovative adaptive confidence loss mechanism, we demonstrate the feasibility and 508 effectiveness of using weaker models to supervise and improve stronger counterparts. Our findings not only validate the potential of weak-to-strong enhancement but also set the stage for future re-510 search endeavors aimed at unlocking further advancements in AI model performance. This work 511 contributes a significant step forward in the pursuit of more sophisticated, efficient, and capable 512 AI systems, emphasizing the importance of nuanced supervision mechanisms in achieving better 513 performance in vision tasks.

514 515

521

529

REFERENCES 516

- 517 Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan Yuille, Trevor Darrell, Jitendra Malik, and 518 Alexei A Efros. Sequential modeling enables scalable learning for large vision models. arXiv preprint arXiv:2312.00785, 2023. 519
- 520 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. 522 Advances in neural information processing systems, 2020. 523
- Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining 524 Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization: Eliciting strong 525 capabilities with weak supervision. arXiv preprint arXiv:2312.09390, 2023. 526
- 527 Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In International conference on machine learning, 2020.
- Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge review. In 530 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021a. 531
- 532 Yinbo Chen, Zhuang Liu, Huijuan Xu, Trevor Darrell, and Xiaolong Wang. Meta-baseline: Exploring simple meta-learning for few-shot learning. In Proceedings of the IEEE/CVF international conference on computer vision, pp. 9062-9071, 2021b. 534
- 535 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical 536 image database. In 2009 IEEE conference on computer vision and pattern recognition, 2009. 537
- 538 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Un-539 terthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
 - 10

540 541	Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. <i>arXiv</i> preprint arXiv:2202.10936, 2022.
542 543 544	Thibaut Durand, Taylor Mordan, Nicolas Thome, and Matthieu Cord. Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation. In <i>Proceedings of the IEEE</i>
545	conference on computer vision and pattern recognition, 2017.
546	Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again
547	neural networks. In International Conference on Machine Learning, pp. 1607–1616. PMLR, 2018.
548 549	Zhiwei Hao, Jianyuan Guo, Kai Han, Han Hu, Chang Xu, and Yunhe Wang. Vanillakd: Revisit the power of vanilla knowledge distillation from small scale to large scale. <i>arXiv preprint arXiv:2305.15781</i> , 2023a.
550	Zhiwei Hao, Jianyuan Guo, Kai Han, Yehui Tang, Han Hu, Yunhe Wang, and Chang Xu. One-for-all: Bridge
551 552	the gap between heterogeneous architectures in knowledge distillation. <i>arXiv preprint arXiv:2310.19444</i> , 2023b.
553	Kaiming Ha Viangun Zhang, Shaaging Dan and Jian Sun. Dean residual learning for image responsition. In
554 555	Proceedings of the IEEE conference on computer vision and pattern recognition, 2016.
556 557 558	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 16000–16009, 2022a.
559	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders
56U	are scalable vision learners. In Proceedings of the IEEE/CVF conference on computer vision and pattern
562	recognition, 2022b.
563	Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive
564 565	Vision, 2019.
566 567	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> , 2015.
568	Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from
569 570	large weakly supervised data. In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14, 2016.
571 572 573	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> , 2020.
575 576	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. <i>arXiv preprint</i> arXiv:2304.02643.2023
577	urxiv.2304.02043, 2023.
578	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
579 580	Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. <i>Proceedings of the IEEE</i> , 1998.
581	Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural
583	networks. In Workshop on challenges in representation learning, ICML, 2013.
584	Junnan Li Richard Socher and Steven CH Hoj Dividemix. Learning with poisy labels as semi-supervised
585	learning. arXiv preprint arXiv:2002.07394, 2020.
586	Shikun Li Visaba Via Shiming Ga and Tangliang Liu. Selective supervised contractive learning with poiss
587 588	labels. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 316– 325, 2022.
589	Wannya Daak Danain Kim Van Lu and Minau Cha. Deletional hereitades distillation. In D
590	<i>EEE/CVF conference on computer vision and pattern recognition.</i> 2019.
591	
592 593	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 2019.

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 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.
- Karsten Roth, Lukas Thede, Almut Sophia Koepke, Oriol Vinyals, Olivier Hénaff, and Zeynep Akata. Fantastic
 gains and where to find them: On the existence and prospect of general knowledge transfer between any
 pretrained model. *arXiv preprint arXiv:2310.17653*, 2023.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition.
 arXiv preprint arXiv:1409.1556, 2014.
- 617 Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for
 618 learning with noisy labels. In *Proceedings of the IEEE conference on computer vision and pattern recogni-* 619 *tion*, 2018.
- Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro
 Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. Advances in neural information processing systems, 29, 2016.
 - Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10687–10698, 2020.
- Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim:
 A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9653–9663, 2022.
- 632 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. arXiv preprint arXiv:1605.07146, 2016.
- Kiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
 - Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2022.
- 638 639 640

641

637

626

627

628

596

A APPENDIX: IMPLEMENTATION DETAILS.

642 A.1 IMAGENET CLASSIFICATION

CIFAR-100. We adopt the vision architectures of the teacher and student models as outlined in
 the traditional distillation papers (Hao et al., 2023b; Zhao et al., 2022). It should be noted that the
 previous codebase (Zhao et al., 2022) conducted experiments on CIFAR-100 using only 1 GPU. To
 expedite our experiments, we leverage the distributed Pytorch framework (Paszke et al., 2019) to
 train and do inference on 8 GPUs. Consequently, some hyperparameter settings and results may not

align exactly with the previous paper. Specifically, we employ the SGD optimizer with a momentum of 0.9. The learning rate starts at 0.2 and decays with a minimum learning rate of 2e-3 using a cosine annealing schedule. We train for 240 epochs with a batch size of 512 spread across 8 GPUs, and apply a weight decay of 0.0005. Standard data augmentation techniques, including random resized crop and horizontal flip are utilized.

ImageNet. we employ the SGD optimizer with a momentum of 0.9. The learning rate starts at 0.1 and decays with a rate of 0.1 every 30 epochs. We train for 100 epochs with a batch size of 512 spread across 8 GPUs, and apply a weight decay of 0.0001. Standard data augmentation techniques, including random resized crop, horizontal flip and label smoothing are utilized.

657 658

659

A.2 TRANSFER LEARNING.

To fine-tune the self-supervised pretrained ViT-B on ImageNet and iNaturalist, we adopt the hyperparameter settings from MAE (He et al., 2022b). The adamw optimizer is employed for this purpose. The learning rate begins at 2e-3 and gradually decays with a minimum learning rate of 1e-6, utilizing a cosine annealing schedule. We conduct training for 100 epochs, utilizing a batch size of 4096 across 8 GPUs. A weight decay of 0.05 is applied to mitigate overfitting. The fine-tuning process incorporates robust data augmentation techniques, including auto-augment, mixup, cutmix, and stochastic drop path.

- 667 668
- A.3 FEW-SHOT LEAERNING

669 We use ResNet12 and follow the setting of (Chen et al., 2021b) on miniImageNet dataset, and 670 created ResNet18 and ResNet36 by increasing the number of layers in original ResNet12. For the 671 classification training stage, we use the SGD optimizer with momentum 0.9. The learning rate 672 starts from 0.1 and the decay factor is set to 0.1. On miniImageNet, we train 100 epochs with the batch size of 128 on 4 GPUs, the learning rate decays at 90 epoch, and the weight decay is 0.0005. 673 Standard data augmentation strategies including random resized crop and horizontal flip are applied. 674 For meta-learning stage, we use the SGD optimizer with momentum 0.9. The learning rate is fixed 675 as 0.001. The batch size is set to 4, *i.e.*, each training batch contains 4 few-shot tasks to compute 676 the average loss. The cosine scaling parameter τ is initialized as 10. For knowledge distillation, the 677 kd loss weight is set to 1, the temperature is set to 10. We use the threshold with 8 and 0.25 for 678 classifier stage and meta stage, respectively. 679

680 681

A.4 LEARNING WITH NOISY LABELS

682 For CIFAR-10/100 datasets, we follow (Li et al., 2022) use a PreAct ResNet18 network, and created 683 PreAct ResNet34 by increasing the number of layers in PreAct ResNet12. We train our models 684 using SGD with a momentum of 0.9, a weight decay of 1e4, and a batch size of 128. The network 685 is trained for 250 epochs and the warm-up epoch is set to 1 dufring training stage. We set the initial learning rate as 0.1, and reduce it by a factor of 10 after 125 and 200 epochs. The fine-tuning stage of 686 Sel-CL+ has 70 epochs, where the learning rate is 0.001. We always set the Mixup hyperparameter 687 to 1, scalar temperature to 0.1, and loss weights to 1. We try two settings of simulated noisy labels: 688 symmetric and asymmetric. And the noise ratio is set to 0.2 and 0.4, respectively. For knowledge 689 distillation, we set the threshold to 0.5 and assign a weight of 1 to the knowledge distillation loss. 690

- 691
- 692 693
- 694
- 695
- 696
- 697
- 698
- 699 700
- 700