# DYNAMIC WEIGHTING: EXPLOITING THE POTENTIAL OF A SINGLE WEIGHT ACROSS DIFFERENT MODES

Anonymous authors

Paper under double-blind review

## Abstract

Weights play an essential role in determining the performance of deep networks. This paper introduces a new concept termed "Weight Augmentation strategy"(WAS), which emphasizes the exploration of weight spaces rather than traditional network structure design. The core of WAS is the utilization of randomly transformed weight coefficients, referred to as Shadow Weights (SW), for deep networks to calculate the loss function and update the parameters. Differently, stochastic gradient descent is applied to Plain Weights (PW), which is referred to as the original weight of the network before the random transformation. During training, numerous SW collectively form a high-dimensional space, while PW is directly learned from the distribution of SW. To maximize the benefits of WAS, we introduce two operational modes, *i.e.*, the Accuracy-Priented Mode (AOM) and the Desire-Oriented Mode (DOM). To be concrete, AOM relies on PW, which ensures that the network remains highly robust and accurate. Meanwhile, DOM utilizes SW, which is determined by the specific objective of our proposed WAS, such as reduced computational complexity or lower sensitivity to particular data. These dual modes can be switched at any time as needed, thereby providing flexibility and adaptability to different tasks. By extending the concept of augmentation from data to weights, our WAS offers an easy-to-understand and implement technique that can significantly enhance almost all networks. Our experimental results demonstrate that convolutional neural networks, including VGG-16, ResNet-18, ResNet-34, GoogleNet, MobileNetV2, and EfficientNet-Lite, benefit substantially with little to no additional costs. On the CIFAR-100 and CIFAR-10 datasets, model accuracy increases by an average of 7.32% and 9.28%, respectively, with the highest improvements reaching 13.42% and 18.93%. In addition, DOM can reduce floating point operations (FLOPs) by up to 36.33%.

034

037

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

#### 1 INTRODUCTION

038 learning with data augmentation (DA) has achieved great success, which significantly enhances 039 model performance through various preprocessing methods such as rotation, translation, scaling, and 040 random cropping Moreno-Barea et al. (2020); Shorten et al. (2021); Yang et al. (2022). The primary 041 goal of DA is to align the distribution of the original dataset more closely with that of natural scenes, 042 thereby increasing both the diversity and quantity of data. Wang et al. Hao & Zhili (2020) advanced 043 this field by improving the Mosaic data augmentation algorithm, which analyzes synthetic image 044 areas and randomly fills them with a certain number of training images. Additionally, researchers at CMU Trabucco et al. (2023) proposed the DA-Fusion strategy using pre-trained text-to-image diffusion models to generate diverse variants of real images, thereby enhancing data variety and 046 improving model robustness. 047

Although complex DA techniques can enhance model accuracy, they come with notable drawbacks
He et al. (2019); Dosovitskiy et al. (2020); Cubuk et al. (2020); Tian et al. (2020); Zoph et al.
(2020). First, intricate augmentation strategies can make models difficult to implement and customize. Second, it is challenging to avoid generating irrelevant tasks or misleading augmented data.
For example, issues such as data over-enhancement, where the model cannot effectively learn from the augmented data, and deviations of the augmented data from the original distribution can arise Suresh et al. (2021); Zheng et al. (2024); Wang et al. (2022); Jiang et al. (2020); Gou et al. (2021).



Figure 1: This figure contrasts traditional and innovative model training strategies. Panel (a) depicts
 the traditional method where specific parameters are derived directly from the data, represented by
 red circles. In contrast, panel (b) introduces our novel strategy where the distribution of weights,
 rather than the weights themselves, is learned from the data.

Weights are the core part of model deployment and the basis for decision-making in practical ap-073 plications, which do not directly interfere with the distribution of real data Zamfirescu-Pereira et al. 074 (2023); Alberts et al. (2023); Ghosal et al. (2023). Since the advent of deep learning, the dominant 075 approach to model training has been to utilize data to optimize model parameters, aiming to obtain 076 a set of weights that lead to superior performance. However, the outcomes are not always favorable. 077 Izmailov et al. Izmailov et al. (2018) proposed stochastic weight averaging, which integrates the weights generated during the training stage by averaging multiple weights, thereby abandoning the 079 final weight result. In traditional training methods, inference relies on a single weight, which can be directly obtained. However, identifying a set of parameters from the weight space that satisfies the 081 feature space mapping is challenging, particularly when both the weight space and the feature space are high-dimensional. Figure 1(a) illustrates the traditional training method, where the relationship 083 between data and weights is many-to-one. Model training aims to find weights that correspond to the data. The formula can be expressed as follows: 084

085

072

$$: \mathbb{R}^m \to x \tag{1}$$

087 In this context,  $\Phi$  represents the function that maps an *m*-dimensional feature space to the corre-088 sponding weights. Specifically,  $\mathbb{R}^m$  is defined as the *m*-dimensional Euclidean space representing the feature vectors. Likewise,  $\mathbb{R}^n$  denotes the *n*-dimensional space of potential weight vectors, 089 illustrating the multidimensional space where the weights reside. Figure 1(b) illustrates our innova-090 tive strategy, where the relationship between data and weights is conceptualized as many-to-many. 091 Throughout the training process, the weight space becomes increasingly complete. With the ap-092 plication of the proposed WAS, the weight space will be further expanded. Instead of seeking the optimal weights directly, our approach involves learning the distribution of the entire weight space. 094 The relationship is mathematically expressed as follows: 095

Φ

096

$$\Psi: \mathbb{R}^m \to \mathbb{R}^n \tag{2}$$

where  $\Psi$  denotes the transformation that maps the *m*-dimensional feature space to the *n*-dimensional 098 weight space. To address the shortcomings of traditional methods like Dropout Srivastava et al. (2014), which combats overfitting by randomly dropping units during training but does not reduce 100 computational complexity in inference, we introduce Weight Augmentation Strategy (WAS). This 101 technology is influenced by both dropout and data augmentation strategies and takes a novel ap-102 proach by randomly transforming weights during the training phase. These transformed weights, 103 referred to as SW, are employed to compute the loss function and facilitate the update of model 104 parameters. In contrast, during the inference phase, Stochastic Gradient Descent (SGD) is applied 105 to PW. This dual-weight strategy ensures that the model achieves optimal performance only when PW aligns effectively with the majority of SW. As a result, WAS allows for a more robust system 106 that improves upon the limitations of dropout by considering the computational demands during 107 inference.



Figure 2: Sketch of WAS architecture. There exist two modes during inference, namely the accuracy-oriented mode using PW and the desire-oriented mode using SW.

While our system utilizes two different types of weights, only the PW need to be saved. The SW are generated through random transformations applied to PW. This approach allows us to maintain a single set of weights that can operate in two distinct modes, which offers different functionalities at minimal cost. Depending on the task requirements, the operational mode can be switched dynam-ically. As depicted in Figure 2, AOM utilizes PW to ensure high robustness and optimal handling of most data types. In contrast, DOM employs SW, which can be tailored for specific functions, such as reducing computational complexity and enhancing sensitivity to particular data types. Our contributions can be summarized as follows: 

- We propose WAS, which enhances the robustness of model training through the random transformation of weights.
- The concept of dual-mode weights allows for retaining just one set of model weights during inference, enabling the model to adapt to multiple states and meet various task requirements efficiently.
- Our approach redefines traditional training by focusing on the entire distribution of weights rather than optimizing for a specific weight, with the goal of identifying the most effective and robust weight configuration for diverse applications.
- 2 RELATED WORK

With the continuous development of deep learning, researchers have begun to seek more efficient and economical methods to obtain model weights. WiSE-FT Wortsman et al. (2022) enhances the robustness of fine-tuned pre-trained models by integrating the weights of zero-shot and fine-tuned models. This method can maintain high accuracy and adapts well to changes in data distribution. However, the performance of WiSE-FT largely depends on the pre-trained model and fine-tuning data. Guo et al. (2020) introduced a collaborative knowledge distillation approach that trains multiple student models simultaneously. This strategy enhances learning efficacy without the need for a separate teacher model, though the selection and quantity of student models impact the outcome and require substantial computational resources. Zhang et al. (2019) developed a self distillation technique that employs the network as both the student and the teacher, facilitating internal knowledge transfer. It does not require an additional pre-trained teacher model and can improve accuracy without increasing inference time. However, self distillation introduces an additional shallow classifier, which prevents the model convergence and increases the complexity of training. 

To overcome the limitations of these methods, we introduce the Weight Augmentation Strategy (WAS). We randomly transform PW during training to obtain SW. PW is used to be compatible with various SW, in order to enhance the robustness of the network and reduce the sensitivity of the model to noise.



Figure 3: Interaction Triangle of WAS Components. SW are utilized to compute the loss function, which influences the adjustment of PW.

3 Method

171

172 173

174 175 176

188 189 190

199

3.1 WEIGHT AUGMENTATION STRATEGY

177 In traditional training, it is a common occurrence that many weights generated are considered "by-178 products" and typically disregarded as they are not optimal for broader data applications, despite 179 their effectiveness on specific datasets. These weights are usually discarded during the model selection process. However, with the deepening of deep network research, the potential value of 181 these "by-products" has been recognized. For example, ensemble learning methods Dietterich et al. 182 (2002); Krawczyk et al. (2017); Huang et al. (2009) can integrate these seemingly useless weights 183 to form a more powerful and robust model. These conventional approaches of leveraging diverse weights have not fully tapped into the potential of weight augmentation. A robust weight space, where weights are learned from distributions rather than specific data points, can significantly im-185 prove performance. The augmentation involves transforming weights through various methods such as rotation, translation, and scaling, as shown in the equation: 187

$$h(x) = ReLU\left(TW'x\right) = ReLU\left(\sum_{i=1}^{k} T_iW_ix_i + b\right)$$
(3)

where W and W' denote the weight vector and its transposition, respectively, and  $T_i$  represents the transformation applied to each weight component.

However, merely reorganizing weights does not fully capture the weight distribution and can lead to biases and reduced generalization ability. To address the above issues, we introduce randomness into the augmentation process, thereby enhancing the diversity and representativeness of the weight space:

$$h(x) = ReLU\left(\sum_{i=1}^{j} \gamma(T_i) W_i x_i + b\right), \tag{4}$$

where  $\gamma(T_i)$  controls the degree of randomness in the transformation, which is crucial for optimizing the efficacy of WAS. The application of randomness is adjusted based on the requirements of convolutional layers and the specific augmentation strategy chosen, such as the degree of translation or scaling. By integrating these elements, WAS aims to refine the way weights are used and adjusted, promoting a more dynamic and flexible approach to deep learning model training.

WAS is an innovative strategy designed to optimize deep networks by directly applying randomness
 and constraints to weights, thus enhancing the model's capability to explore the weight space. This
 approach not only addresses the limitations faced by traditional weight optimization methods, such
 as a lack of flexibility and autonomy, but also revitalizes the training process by introducing dynamic
 changes in weight configurations.

As training continues, the distribution of SW is refined, enhancing the overall model parameters.
 The updating of model parameters is governed by the following equations:

$$\theta_j = \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta) \tag{5}$$

$$\theta_j^{pw} = \theta_j^{pw} - \alpha \cdot \frac{\partial}{\partial \theta_j^{pw}} J(\theta_j^{sw}) \tag{6}$$

Table 1: Deep networks defined by WAS				
Name	Crop	Translate	Rotate	
Model-C Model-T Model-R Model-CT Model-CTR				
MOUCH-CTK	ν	ν	V	

Here,  $\theta_j^{sw}$  and  $\theta_j^{pw}$  represent the *j*-th parameters of SW and PW, respectively.  $\alpha$  denotes the learning rate and  $J(\theta)$  is the loss function. The derivative  $\frac{\partial}{\partial \theta_j^{sw}}$  reflects the partial derivative of the loss function  $J(\theta_j^{sw})$  with respect to the parameter  $\theta_j^{pw}$ .

Additionally, the interaction between SW and PW within the training process is captured by

$$\frac{\partial J}{\partial \theta_j} = \sum_{x,y} \frac{\partial J}{\partial Z_{xy}} \cdots \frac{\partial Z_{xy}}{\partial \theta_j},\tag{7}$$

234 where  $Z_{xy}$  represents the feature map at position (x, y). The partial derivatives  $\frac{\partial J}{\partial Z_{xy}}$  and  $\frac{\partial Z_{xy}}{\partial \theta_i}$ 235 quantify the sensitivity of the loss function to changes in the feature map and the feature map's 236 responsiveness to changes in the weights, respectively. Here, SW is employed to compute the loss 237 function, which assesses the discrepancy between the model's outputs and the actual results. These 238 computations not only impact PW adjustments but also facilitate learning the data distribution via 239 SW. Thus, PW is primarily updated based on the distribution of SW, which is a significant departure 240 from conventional methods where direct data-driven updates are commonplace. The relationship between PW, SW, and loss function is shown in Figure 3. The WAS incentive model integrates 241 hundreds of weight training iterations, leveraging the advantages of SW while mitigating potential 242 drawbacks. This strategy enhances model performance and allows for mode switching based on task 243 requirements, paving the way for more intelligent and efficient deep network training. 244

In related research, Ding *et al.* Ding et al. (2021) introduced RepVGG, a method that adjusts weights during inference to reduce time complexity while preserving performance. This approach explores the potential relationship between model structures and weights, although it does not directly link the two. Meanwhile, Zheng *et al.* Zheng et al. (2023) introduced the "Learn From Model" concept, emphasizing the need for deeper exploration and modification of foundation models based on model interfaces.

Our approach takes a step further by advocating for weight-based training over traditional datadriven methods. By focusing on weight distributions rather than direct data features, our method allows the model to learn from a broader array of potential configurations, encompassing even flawed weights that can still capture specific data nuances. This two-stage training process—alternating between learning from data distributions and flawed weight distributions—leads to a model that optimizes its structure and weights autonomously without human intervention.

WAS not only encourages exploration of a wider weight space to find high-performance configura tions but also enhances computational efficiency. This reduces resource consumption during both
 training and inference. Moreover, WAS can be tailored to specific model functions, such as pro cessing particular types of data or reducing computational complexity, which aligns with the unique
 needs of various applications.

262 263

264

216

224 225

226 227

228

229

3.2 DUAL WORKING MODE OF NETWORK MODEL

The core of WAS is to encourage the generation of a large number of weights, which enriches the weight distribution and enhances the generalization and robustness of the model. The model operates in two distinct modes during training: AOM and DOM.

AOM operates similarly to traditional deep learning models and is utilized directly in inference. The innovation in AOM lies in its dynamic weight formation mechanism, which emerges through a competitive process during training. This competition drives continual optimization of the model's

CIFAR10		CIFAR100		Average FLOPs (M)
AOM Top1	DOM Top1	AOM Top1	GOM Top1	-
87.42	-	59.01	-	333
89.29	88.95	62.19	60.53	321
91.15	90.60	63.53	61.96	279
85.53	-	59.70	-	608
88.15	87.38	63.42	62.44	542
89.83	88.64	63.40	61.50	315
86.54	-	58.64	-	1214
89.18	88.81	62.19	61.24	1099
91.71	90.73	60.64	58.73	709
90.55	-	72.04	-	1457
91.99	91.43	73.05	72.02	1237
92.68	92.04	73.08	72.01	772
73.2	-	47.02	-	47
82.58	82.43	55.43	54.65	45
83.02	82.88	55.92	53.95	37
73.05	-	43.13	-	8.00
80.96	79.97	49.08	48.47	6.37
83 72	81.23	50.31	47.41	4.66
	CIFA AOM Top1 87.42 89.29 91.15 85.53 88.15 89.83 86.54 89.18 91.71 90.55 91.99 92.68 73.2 82.58 83.02 73.05 80.96 83.72	CIFAR10AOM Top1DOM Top187.42-89.2988.9591.1590.6085.53-88.1587.3889.8388.6486.54-89.1888.8191.7190.7390.55-91.9991.4392.6892.0473.2-82.5882.4383.0282.8873.05-80.9679.9783.7281.23	CIFAR10CIFAAOM Top1DOM Top1AOM Top1 $87.42$ -59.01 $89.29$ $88.95$ $62.19$ $91.15$ 90.60 $63.53$ $85.53$ -59.70 $88.15$ $87.38$ $63.42$ $89.83$ $88.64$ $63.40$ $86.54$ - $58.64$ $89.18$ $88.81$ $62.19$ $91.71$ 90.73 $60.64$ $90.55$ - $72.04$ $91.99$ $91.43$ $73.05$ $92.68$ $92.04$ $73.08$ $73.2$ - $47.02$ $82.58$ $82.43$ $55.43$ $83.02$ $82.88$ $55.92$ $73.05$ - $43.13$ $80.96$ $79.97$ $49.08$ $83.72$ $81.23$ $50.31$	CIFAR10CIFAR100AOM Top1DOM Top1AOM Top1GOM Top1 $87.42$ -59.01- $89.29$ $88.95$ $62.19$ $60.53$ $91.15$ 90.60 $63.53$ $61.96$ $85.53$ - $59.70$ - $88.15$ $87.38$ $63.42$ $62.44$ $89.83$ $88.64$ $63.40$ $61.50$ $86.54$ - $58.64$ - $91.71$ 90.73 $60.64$ $58.73$ $90.55$ - $72.04$ - $91.99$ $91.43$ $73.05$ $72.02$ $92.68$ $92.04$ $73.08$ $72.01$ $73.2$ - $47.02$ - $82.58$ $82.43$ $55.43$ $54.65$ $83.02$ $82.88$ $55.92$ $53.95$ $73.05$ - $43.13$ - $80.96$ $79.97$ $49.08$ $48.47$ $83.72$ $81.23$ $50.31$ $47.41$

Table 2: Comparison of top-1 accuracy (%) and average FLOPs (M) on CIFAR10 and CIFAR100 datasets

270

performance. Within this framework, SW that demonstrate superior performance can significantly
 reduce the loss function, thereby influencing the determination of PW. Weights that perform poorly
 exert less influence on the final weight determination.

301 DOM is specifically tailored to meet unique operational requirements such as enhancing model 302 specificity, sparsity, and computational efficiency. In DOM, we can devise WAS tailored to the particular demands of the task and implement it within training. This adaptive strategy allows for the 303 tuning of PW for use during inference. Although DOM may result in lower prediction accuracy 304 compared to AOM, it offers substantial benefits for specialized tasks. For example, implementing 305 strategies like random weight cropping can increase the sparsity of the weight matrix, thereby dras-306 tically reducing the computational load—sometimes by several orders of magnitude. This reduction 307 is particularly valuable in scenarios involving large-scale data processing or in resource-constrained 308 environments. 309

As a consequence, WAS provides a flexible model management strategy that only requires the storage of one set of weights while enabling two distinct functional modes. During inference, the model can switch between these modes based on the specific requirements of the application scenario. This dual-mode operation not only simplifies the management of model storage and maintenance but also enhances the algorithm's resilience and adaptability to varying conditions.

315 316

## 4 EXPERIMENTS

317 318

To evaluate the effectiveness of the proposed WAS, we conducted a series of ablation studies and comparative analyses on the CIFAR-10 and CIFAR-100 datasets Krizhevsky et al. (2009). The main objective of these experiments was to demonstrate the impact that WAS integration has on the models. Furthermore, we evaluate to delineate the distinctions and performance implications of dual operational modes. The results from these studies consistently confirm the advantages of incorporating WAS, showcasing substantial enhancements in model performance across various metrics.

<sup>296</sup> 297

345

Model	Data Augmentation	Parameters	AOM Top-1 Drop Rate(%)	GOM Top-1 Drop Rate(%)
		$(0^{\circ}, 15^{\circ})$	5.09	4.86
		$(0^{\circ}, 45^{\circ})$	39.52	25.58
VGG16-R	rotate	$(0^{\circ}, 90^{\circ})$	46.29	42.74
$(0^{\circ}, 90^{\circ})$		$(0^{\circ}, 135^{\circ})$	54.61	53.65
		$(0^{\circ}, 180^{\circ})$	58.40	57.76
		(10%,10%)	1.72	2.33
VGG16-T	translate	(20%,20%)	6.62	8.47
(30%,30%)		(30%,30%)	14.42	18.21
		(40%,40%)	25.76	31.76
		(0.8,1.0)	3.34	2.79
VGG16-C	crop	(0.6, 1.0)	6.49	6.71
(0.8, 1.0)	-	(0.4, 1.0)	13.04	12.12
		(0.2, 1.0)	24.19	22.42

Table 3: Comparison of AOM and DOM Top-1 Drop Rates under Different WAS

#### 4.1 WAS FOR CLASSIFICATION

To evaluate the effectiveness of WAS, we have chosen six prominent deep learning architectures as our experimental models: VGG-16 Simonyan & Zisserman (2014), ResNet18 He et al. (2016), ResNet34 He et al. (2016), GoogLeNet Szegedy et al. (2015), EfficientNet-Lite Tan & Le (2019), and MobileNetV2 Sandler et al. (2018). The results of these comparative experiments are detailed in Table 1, showing improvements against baseline models.

We deliberately eschewed the incorporation of additional techniques, primarily to mitigate the impact of extraneous variables. On a solitary GPU, we established a global batch size of 128. We employed the conventional SGD, initializing the learning rate at 0.01. Furthermore, we fine-tuned the SGD optimizer, assigning a momentum coefficient of 0.9, thereby augmenting the model's stability and hastening convergence throughout the training regimen.

356 Table 2 demonstrates the performance improvement achieved by Model-C and Model-CT under 357 different architectures. On the CIFAR-10 dataset, AOG for the models demonstrated discernible en-358 hancements. Specifically, the accuracy of VGG16-C and VGG16-CT saw an increase of 2.13% and 359 4.27% relative to the baseline VGG16, respectively. For ResNet18, ResNet18-C and ResNet18-CT 360 attained an accuracy improvement of 3.06% and 5.03%. Similarly, the ResNet34-C and ResNet34-361 CT models realized respective accuracy improvements of 3.05% and 5.97%. The GoogleNet-C 362 and GoogleNet-CT models recorded accuracy enhancements of 1.59% and 2.35%. Within the Mo-363 bileNetV2 series, the accuracy for MobileNetV2-C and MobileNetV2-CT marked a significant rise of 12.76% and 13.42% over MobileNetV2. Lastly, the EfficientNetLite-C and EfficientNetLite-CT 364 models secured accuracy improvements of 9.52% and 13.23%.

- 366 On the CIFAR-100 dataset, WAS also demonstrated enhanced performance across various mod-367 els. The VGG16-C and VGG16-CT models showed accuracy improvements of 5.39% and 7.66%, 368 respectively. The ResNet18-C and ResNet18-CT models recorded increases of 6.23% and 6.19%, while the ResNet34-C and ResNet34-CT achieved accuracy gains of 6.05% and 3.41%, respectively. 369 GoogleNet-C and GoogleNet-CT noted slight improvements with gains of 1.40% and 1.44%. Re-370 markably, MobileNetV2-C and MobileNetV2-CT marked significant advancements with increases 371 of 17.89% and 18.93%. Finally, the EfficientNetLite-C and EfficientNetLite-CT models displayed 372 notable accuracy enhancements of 13.78% and 18.93%, respectively, illustrating the substantial im-373 pact of WAS on model performance. 374
- The performance of Model-C and Model-CT under DOM is slightly lower than under AOM, both models still outperform their baseline by approximately 1% to 2%. This indicates that both operational modes enhance model performance, with particularly notable gains in lightweight models. Furthermore, while not the primary focus, it is important to acknowledge the reductions in compu-

7

380	enhanc <u>ements</u>				
381	Model	Data	Parameters	AOM Top-1	DOM Top-1
382		Augmentation	I al anicter s	Drop Rate(%)	Drop Rate(%)
383			(0.8, 1.0)	2.33	2.50
384		corp	(0.6, 1.0)	5.93	5.15
385			(0.4, 1.0)	11.35	11.70
386	VGG16-R		(0.2, 1.0)	22.31	23.08
387	(0°,90°)		(10%,10%)	2.56	3.27
388		translate	(20%,20%)	9.96	10.39
389			(30%,30%)	19.95	20.45
390			(40%,40%)	32.45	32.71
391			$(0^{\circ}, 15^{\circ})$	5.09	5.27
392			$(0^{\circ}, 45^{\circ})$	28.39	27.97
393		rotate	$(0^{\circ}, 90^{\circ})$	46.80	47.28
304	VGG16-T		$(0^{\circ}, 135^{\circ})$	55.18	56.28
305	(30%,30%)		$(0^{\circ}, 180^{\circ})$	58.96	59.67
206			(0.8,1.0)	2.69	2.61
390		crop	(0.6, 1.0)	5.16	6.11
397			(0.4, 1.0)	10.12	12.18
398			(0.2, 1.0)	21.53	25.10
399			$(0^{\circ}, 15^{\circ})$	4.80	4.85
400			$(0^{\circ}, 45^{\circ})$	26.68	28.14
401		rotate	$(0^{\circ}, 90^{\circ})$	46.10	45.59
402	VGG16-C		$(0^{\circ}, 135^{\circ})$	54.84	53.04
403	(0.8,1.0)		(0°,180°)	58.98	56.51
404			(10%,10%)	2.56	3.24
405		translate	(20%,20%)	8.50	9.87
406			(30%,30%)	17.70	18.21
407			(40%,40%)	28.81	33.70

Table 4: Comparison of AOM and DOMTop-1 drop rates of different WAS under different data enhancements

408 409

412

378

tational load as measured by FLOPs. For Model-C, FLOPs decreased by about 5% to 20%, whereas
Model-CT saw a more substantial reduction, with FLOPs decreasing by up to 47%.

413 4.2 CHARACTERISTIC OF WAS

WAS endows models with unique capabilities that can be tailored to specific needs during inference. These capabilities encompass a range of enhancements, including reduced computational complexity and decreased sensitivity to particular types of data. This functionality is achieved by switching to SW during inference, allowing for customizable adjustments according to specific requirements. Consequently, WAS enables the model to adapt effectively to unique operational environments.

Table 3 showcases three representative WAS strategies implemented following the official PyTorch example Imambi et al. (2021): rotation (random rotation angles set between 0° and 90°), translation (random translation of weights up to 30% in horizontal and vertical directions), and cropping (random cropping of weights with ratios between 0.8 and 1.0). To assess the effectiveness of these strategies for processing specific data, we conducted random data augmentation on the test set.

After using WAS, the network operates in two modes: AOM, which does not apply WAS during inference, and DOM, which applies WAS during inference. As detailed in Table 3, when employing a random rotation strategy for WAS training (0° to 90°), AOM consistently outperforms DOM in Top-1 accuracy. Notably, with rotations limited to 0° to 45° and 0° to 90°, the accuracy losses in DOM compared to AOM are significantly reduced by 13.84% and 3.55%, respectively. For other rotation angles, the accuracy loss in DOM is approximately 1% lower than in AOM.

431 Additionally, when implementing the cropping strategy, as the weight cropping ratio increases, the accuracy loss in DOM compared to AOM decreases progressively by 0.58%, 0.22%, 0.92%, and

4	3	6
4	3	7
4	3	8
4	3	ŝ
4	4	(

Table 5: Comparison of the impact of Top-1 drop rates of different randomly cropped WAS on rotation data

WAS Strategy	AOM Top-1 acc	DOM Top-1 acc	Parameters	AOM Top-1 Drop Rate(%)	DOM Top-1 Drop Rate(%)
	89.72	88.95	$(0^{\circ}, 15^{\circ})$	4.80	4.85
			$(0^{\circ}, 45^{\circ})$	26.68	28.14
Crop (0.8,1.0)			(0°,90°)	46.10	45.59
			$(0^{\circ}, 135^{\circ})$	54.84	53.04
			(0°,180°)	58.98	56.51
	90.34	89.71	(0°,15°)	5.58	5.04
			$(0^{\circ}, 45^{\circ})$	27.85	26.85
Crop (0.6,0.8)			(0°,90°)	46.54	45.82
			$(0^{\circ}, 135^{\circ})$	56.18	53.37
			(0°,180°)	59.98	56.16
	88.54	87.99	(0°,15°)	5.79	4.51
			$(0^{\circ}, 45^{\circ})$	27.44	24.62
Crop (0.4,0.6)			(0°,90°)	45.67	43.25
			$(0^{\circ}, 135^{\circ})$	54.66	51.57
			(0°,180°)	58.01	53.99
	87.09	86.19	(0°,15°)	5.96	4.43
			$(0^{\circ}, 45^{\circ})$	27.57	25.72
Crop (0.2,0.4)			(0°,90°)	45.85	44.4
			$(0^{\circ}, 135^{\circ})$	54.52	52.77
			(0°,180°)	58.01	54.88

1.77%. This indicates that the disparity in accuracy loss between DOM and AOM first narrows
and then widens, reflecting the impact of increasing random crop ratios during training, which can
degrade model performance if parameters exceed certain thresholds.

In summary, WAS helps to promote the ability to process specific data. However, when WAS involves random translation, the accuracy loss in DOM exceeds that in AOM, contradicting initial expectations. Upon integrating random cropping with translation, it appears that the reduction in model parameters may diminish the fitting capabilities of the model.

As depicted in Table 4, evaluated the effects of random translations within the test dataset through our AD pipeline, which deviates from traditional AD approaches. As the random translation param-eter increased from 10% to 40%, the Top-1 accuracy drop rate in AOM escalated from 2.56% to 32.45%. Concurrently, DOM experienced a rise in the Top-1 drop rate from 3.27% to 33.70%. As the cropping ratio decreased from 1.0 to 0.2, both modes saw an increase in the Top-1 accuracy drop rate. For AOM, the drop rate increased from 2.69% to 21.53%, while for DOM, it rose from 2.61% to 25.10%. Notably, DOM exhibited a more pronounced performance decline, especially at higher cropping ratios. Based on the preliminary analysis of the above data, we can conclude that AOM has stronger generalization when the WAS strategy is inconsistent with the data augmentation strategy. 

Similarly, as the rotation angle of the test set increases from 0° to 180°, using a random translation
WAS strategy, the drop rate for AOM escalates from 5.09% to 58.96%. In DOM, the drop rate rises
from 5.27% to 59.67%. Conversely, when employing a random cropping WAS strategy, the drop rate
for AOM increases from 4.80% to 58.98%, while for DOM, it climbs from 4.85% to 56.51%. Under
extreme augmentation conditions, the drop rate in AOM is unexpectedly higher than in DOM, which
is contrary to initial expectations. This indicates that DOM may handle extreme data manipulations
more robustly than AOM.

To explore the above phenomenon, Table 5 presents the models trained with varying random cropping ratios for WAS, while DA employs random rotation. On the test set, despite changes in training parameters, the drop rate of AOM remains relatively stable. However, as the random cropping ratio in WAS increases, the Top-1 accuracy drop rates for DOM are 2.47%, 3.82%, and 4.02% lower than those for AOM respectively, which suggests that DOM is more adept at processing rotated data.

		0			0
488	Cropping	Translation	Mode Two	FLOPs	Average
489	parameters	parameters	Acc	(M)	sparsity rate(%)
490	(1.0.1.0)	_	87.42	333	-
491	(0.8, 1.0)	-	88.95	320.68	3.70
492	(0.6,0.8)	-	89.71	259.67	22.02
493	(0.4,0.6)	-	87.99	223.08	33.01
494	(0.2, 0.4)	-	86.19	212.02	36.33
495	-	(30%,30%)	90.37	230.63	30.74
496	(0.8, 1.0)	(30%,30%)	90.60	277.42	16.69

Table 6: Training results on CIFAR10 datasets using different WAS

486

487 488 489

As cropping ratios range from 0.2 to 0.4, the disparity in Top-1 accuracy drop rates between AOM 499 and DOM narrows, which is likely due to a decrease in parameters which affects the model's abil-500 ity to fit data. This underscores how data rotation, similar to cropping, can lead to information 501 loss, yet DOM (employing weight-based random cropping) exhibits a lower drop rate than AOM, 502 highlighting its robustness.

Further analysis of the model's structural efficiency is shown in Table 6. When cropping parameters 504 are set at (0.4,0.6), the model achieves a sparsity rate of 33.01% with zero elements and maintains 505 an accuracy of 87.99%, which is comparable to the base model. When the Cropping parameters 506 are reduced to (0.2,0.6), the accuracy decreases by 1.41%, but the proportion of 0 elements in-507 creases by 3.32%, and the FLOPs are reduced by 11.06M. When introducing translation parameters 508 (30%, 30%), even without cropping, the sparsity rate is increased to 30.74%, floating-point opera-509 tions(FLOPs) are reduced to 230.63M, and the accuracy is improved by 3.26% compared to the base 510 model. Combining cropping with translation through WAS not only enhances model accuracy but 511 also significantly lowers computational costs, demonstrating the dual benefits of this approach in 512 improving model performance and efficiency.

513 514

515

522 523

524

534

#### 5 CONCLUSION

516 We introduce WAS as a training method for deep learning models. Central to WAS is the imple-517 mentation of a dual-mode inference system, which allows the weights to be tailored to meet diverse 518 task requirements effectively. This approach facilitates the fine-tuning of weights to address specific 519 needs precisely. In AOM, WAS has demonstrated the capability to enhance model accuracy by as 520 much as 18.93% without additional computational expenses. In DOM, it can reduce FLOPs by up 521 to 36.33%, while maintaining robust accuracy levels.

# REFERENCES

- 525 Ian L Alberts, Lorenzo Mercolli, Thomas Pyka, George Prenosil, Kuangyu Shi, Axel Rominger, and Ali Afshar-Oromieh. Large language models (llm) and chatgpt: what will the impact on nuclear medicine be? European journal of nuclear medicine and molecular imaging, 50(6):1549–1552, 527 2023. 528
- 529 Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated 530 data augmentation with a reduced search space. In Proceedings of the IEEE/CVF conference on 531 computer vision and pattern recognition workshops, pp. 702–703, 2020. 532
  - Thomas G Dietterich et al. Ensemble learning. The handbook of brain theory and neural networks, 2(1):110-125, 2002.
- 535 Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: 536 Making vgg-style convnets great again. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 13733–13742, 2021. 538
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An

548

570

571

572

583

- image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction-tuned llm and latent diffusion model. *arXiv preprint arXiv:2304.13731*, 2023.
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A
   survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- Qiushan Guo, Xinjiang Wang, Yichao Wu, Zhipeng Yu, Ding Liang, Xiaolin Hu, and Ping Luo.
   Online knowledge distillation via collaborative learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11020–11029, 2020.
- Wang Hao and Song Zhili. Improved mosaic: Algorithms for more complex images. In *Journal of Physics: Conference Series*, volume 1684, pp. 012094. IOP Publishing, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 558–567, 2019.
- Faliang Huang, Guoqing Xie, and Ruliang Xiao. Research on ensemble learning. In 2009 International Conference on Artificial Intelligence and Computational Intelligence, volume 3, pp. 249–252. IEEE, 2009.
- Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. Pytorch. *Programming with TensorFlow: Solution for Edge Computing Applications*, pp. 87–104, 2021.
- Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint* arXiv:1803.05407, 2018.
  - Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020.
- Bartosz Krawczyk, Leandro L Minku, João Gama, Jerzy Stefanowski, and Michał Woźniak. Ensemble learning for data stream analysis: A survey. *Information Fusion*, 37:132–156, 2017.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
   2009.
- Francisco J Moreno-Barea, José M Jerez, and Leonardo Franco. Improving classification accuracy
   using data augmentation on small data sets. *Expert Systems with Applications*, 161:113696, 2020.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo bilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Connor Shorten, Taghi M Khoshgoftaar, and Borko Furht. Text data augmentation for deep learning.
   *Journal of big Data*, 8(1):101, 2021.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
   recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
  Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Susheel Suresh, Pan Li, Cong Hao, and Jennifer Neville. Adversarial graph augmentation to improve graph contrastive learning. *Advances in Neural Information Processing Systems*, 34:15920–15933, 2021.

- 594 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du-595 mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In 596 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015. 597 Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural net-598 works. In International conference on machine learning, pp. 6105–6114. PMLR, 2019. 600 Keyu Tian, Chen Lin, Ming Sun, Luping Zhou, Junjie Yan, and Wanli Ouyang. Improving auto-601 augment via augmentation-wise weight sharing. Advances in Neural Information Processing 602 Systems, 33:19088–19098, 2020. 603 Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective data augmen-604 tation with diffusion models. arXiv preprint arXiv:2302.07944, 2023. 605 Huan Wang, Suhas Lohit, Michael N Jones, and Yun Fu. What makes a" good" data augmentation 607 in knowledge distillation-a statistical perspective. Advances in Neural Information Processing 608 Systems, 35:13456-13469, 2022. 609 Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, 610 Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust 611 fine-tuning of zero-shot models. In Proceedings of the IEEE/CVF conference on computer vision 612 and pattern recognition, pp. 7959–7971, 2022. 613 Suorong Yang, Weikang Xiao, Mengcheng Zhang, Suhan Guo, Jian Zhao, and Furao Shen. Image 614 data augmentation for deep learning: A survey. arXiv preprint arXiv:2204.08610, 2022. 615 616 JD Zamfirescu-Pereira, Richmond Y Wong, Bjoern Hartmann, and Qian Yang. Why johnny can't 617 prompt: how non-ai experts try (and fail) to design llm prompts. In Proceedings of the 2023 CHI 618 Conference on Human Factors in Computing Systems, pp. 1–21, 2023. 619 Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your 620 own teacher: Improve the performance of convolutional neural networks via self distillation. In 621 Proceedings of the IEEE/CVF international conference on computer vision, pp. 3713–3722, 2019. 622 623 Chenyu Zheng, Guoqiang Wu, and Chongxuan Li. Toward understanding generative data augmen-624 tation. Advances in Neural Information Processing Systems, 36, 2024. 625 Hongling Zheng, Li Shen, Anke Tang, Yong Luo, Han Hu, Bo Du, and Dacheng Tao. Learn from 626 model beyond fine-tuning: A survey. arXiv preprint arXiv:2310.08184, 2023. 627 628 Barret Zoph, Ekin D Cubuk, Golnaz Ghiasi, Tsung-Yi Lin, Jonathon Shlens, and Quoc V Le. Learn-629 ing data augmentation strategies for object detection. In Computer Vision-ECCV 2020: 16th 630 European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16, pp. 566– 631 583. Springer, 2020. 632 633 **APPENDIX** А 634 635 You may include other additional sections here. 636 637 638
- 640 641

- 642
- 643 644
- 644 645
- 646
- 647