

---

# Explore Positive Noise in Deep Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In computer vision, noise is conventionally viewed as a harmful perturbation  
2 in various deep learning architectures, such as convolutional neural networks  
3 (CNNs) and vision transformers (ViTs), as well as different tasks like image  
4 classification and transfer learning. However, this paper aims to rethink whether  
5 the conventional proposition always holds. We demonstrate that specific noise can  
6 boost the performance of various deep architectures under certain conditions. We  
7 theoretically prove the enhancement gained from positive noise by reducing the task  
8 complexity defined by information entropy and experimentally show the significant  
9 performance gain in large image datasets, such as the ImageNet. Herein, we use  
10 the information entropy to define the complexity of the task. We categorize the  
11 noise into two types, positive noise (PN) and harmful noise (HN), based on whether  
12 the noise can help reduce the complexity of the task. Extensive experiments of  
13 CNNs and ViTs have shown performance improvements by proactively injecting  
14 positive noise, where we achieve an unprecedented top 1 accuracy over 95% on  
15 ImageNet. Both theoretical analysis and empirical evidence have confirmed that the  
16 presence of positive noise, can benefit the learning process, while the traditionally  
17 perceived harmful noise indeed impairs deep learning models. The different roles  
18 of noise offer new explanations for deep models on specific tasks and provide a  
19 new paradigm for improving model performance. Moreover, it reminds us to utilize  
20 noise rather than suppress noise.

## 21 1 Introduction

22 Noise, conventionally regarded as a hurdle in machine learning and deep learning tasks, is universal  
23 and unavoidable due to various reasons, e.g., environmental factors, instrumental calibration, and  
24 human activities [23] [37]. In computer vision, noise can be generated from different phases: (1)  
25 Image Acquisition: Noise can arise from a camera sensor or other imaging device [33]. For example,  
26 electronic or thermal noise in the camera sensor can result in random pixel values or color variations  
27 that can be visible in the captured image. (2) Image Preprocessing: Noise can be introduced during  
28 preprocessing steps such as image resizing, filtering, or color space conversion [1]. For example,  
29 resizing an image can introduce aliasing artifacts, while filtering an image can result in the loss of  
30 detail and texture. (3) Feature Extraction: Feature extraction algorithms can be sensitive to noise  
31 in the input image, which can result in inaccurate or inconsistent feature representations [2]. For  
32 example, edge detection algorithms can be affected by noise in the image, resulting in false positives  
33 or negatives. (4) Algorithms: algorithms used for computer vision tasks, such as object detection or  
34 image segmentation, can also be sensitive to noise in the input data [6]. Noise can cause the algorithm  
35 to learn incorrect patterns or features, leading to poor performance.

36 Since noise is an unavoidable reality in engineering tasks, existing works usually make the assumption  
37 that noise has a consistently negative impact on the current task [30] [24]. Nevertheless, is the above  
38 assumption always valid? As such, it is crucial to address the question of whether noise can ever

39 have a positive influence on deep learning models. This work aims to provide a comprehensive  
40 answer to this question, which is a pressing concern in the deep learning community. We recognize  
41 that the imprecise definition of noise is a critical factor leading to the uncertainties surrounding the  
42 identification and characterization of noise. To address these uncertainties, an in-depth analysis  
43 of the task’s complexity is imperative for arriving at a rigorous answer. By using the definition of  
44 task entropy, it is possible to categorize noise into two distinct categories: positive noise (PN) and  
45 harmful noise (HN). PN decreases the complexity of the task, while HN increases it, aligning with  
46 the conventional understanding of noise.

## 47 1.1 Scope and Contribution

48 Our work aims to investigate how various types of noise affect deep learning models. Specifically,  
49 the study focuses on three common types of noise, i.e., Gaussian noise, linear transform noise, and  
50 salt-and-pepper noise. Gaussian noise refers to random fluctuations that follow a Gaussian distribution  
51 in pixel values at the image level or latent representations in latent space [29]. Linear transforms, on  
52 the other hand, refer to affine elementary matrix transformations to the dataset of original images  
53 or latent representations, where the elementary matrix is row equivalent to an identity matrix [36].  
54 Salt-and-pepper noise is a kind of image distortion that adds random black or white values at the  
55 image level or to the latent representations [7].

56 This paper analyzes the impact of these types of noise on the performance of deep learning models for  
57 image classification and domain adaptation tasks. Two popular model families, Vision Transformers  
58 (ViTs) and Convolutional Neural Networks (CNNs), are considered in the study. Image classification  
59 is one of the most fundamental tasks in computer vision, where the goal is to predict the class label of  
60 an input image. Domain adaptation is a practically meaningful task where the training and test data  
61 come from different distributions, also known as different domains. By investigating the effects of  
62 different types of noise on ViTs and CNNs for typical deep learning tasks, the paper provides insights  
63 into the influences of noises on deep models. The findings presented in this paper hold practical  
64 significance for enhancing the performance of various types of deep learning models in real-world  
65 scenarios.

66 The contributions of this paper are summarized as follows:

- 67 • We re-examined the conventional view that noise, by default, has a negative impact on deep  
68 learning models. Our theoretical analysis and experimental results show that noise can be a  
69 positive support for deep learning models and tasks.
- 70 • We implemented extensive experiments with different deep models, such as CNNs and  
71 ViTs, and on different deep learning tasks. Empowered by positive noise, we achieved  
72 state-of-the-art (SOTA) results in all the experiments presented in this paper.
- 73 • Instead of operating on the image level, our injecting noise operations are performed in the  
74 latent space. We theoretically analyze the difference between injecting noise on the image  
75 level and in the latent space.
- 76 • The theory and framework of reducing task complexity via positive noise in this work can  
77 be applied to any deep learning architecture. There is great potential for exploring the  
78 application of positive noise in other deep-learning tasks beyond the image classification  
79 and domain adaptation tasks examined in this study.

## 80 1.2 Related Work

81 **Positive Noise** In fact, within the signal-processing society, it has been demonstrated that random  
82 noise helps stochastic resonance improve the detection of weak signals [4]. Noises can have positive  
83 support and contribute to less mean square error compared to the best linear unbiased estimator when  
84 the mixing probability distribution is not in the extreme region [28]. Also, it has been reported that  
85 noise could increase the model generalization in natural language processing (NLP) [27]. Recently,  
86 the perturbation, a special case of positive noise, has been effectively utilized to implement self-  
87 refinement in domain adaptation and achieved state-of-the-art performance [36]. The latest research  
88 shows that by proactively adding specific noise to partial datasets, various tasks can benefit from the  
89 positive noise [19]. Besides, noises are found to be able to boost brain power and be useful in many  
90 neuroscience studies [21] [22].

91 **Deep Model** Convolutional Neural Networks have been widely used for image classification, object  
 92 detection, and segmentation tasks, and have achieved impressive results [18][15]. However, these  
 93 networks have limitations in terms of their ability to capture long-range dependencies and extract  
 94 global features from images. Recently, Vision Transformers has been proposed as an alternative to  
 95 CNNs [13]. ViT relies on self-attention mechanisms and a transformer-based architecture to enable  
 96 global feature extraction and modeling of long-range dependencies in images [40]. The attention  
 97 mechanism allows the model to focus on the most informative features of the input image, while  
 98 the transformer architecture facilitates information exchange between different parts of the image.  
 99 ViT has demonstrated impressive performance on a range of image classification tasks and has the  
 100 potential to outperform traditional CNN-based approaches. However, ViT currently requires a large  
 101 number of parameters and training data to achieve state-of-the-art results, making it challenging to  
 102 implement in certain settings [45].

## 103 2 Preliminary

104 In information theory, the entropy [32] of a random variable  $x$  is defined as:

$$H(x) = \begin{cases} -\int p(x) \log p(x) dx & \text{if } x \text{ is continuous} \\ -\sum_x p(x) \log p(x) & \text{if } x \text{ is discrete} \end{cases} \quad (1)$$

105 where  $p(x)$  is the distribution of the given variable  $x$ . And the mutual information (MI) of two  
 106 random discrete variables  $(x, y)$  is denoted as [8]:

$$\begin{aligned} MI(x, y) &= D_{KL}(p(x, y) \parallel p(x) \otimes p(y)) \\ &= H(x) - H(x|y) \end{aligned} \quad (2)$$

107 where  $D_{KL}$  is the Kullback–Leibler divergence [16], and  $p(x, y)$  is the joint distribution. The  
 108 conditional entropy is defined as:

$$H(x|y) = -\sum p(x, y) \log p(x|y) \quad (3)$$

109 The above definitions can be readily expanded to encompass continuous variables through the  
 110 substitution of the sum operator with the integral symbol. In this work, the noise is denoted by  $\epsilon$  if  
 111 without any specific statement.

112 Before delving into the correlation between task and noise, it is imperative to address the initial  
 113 crucial query of the mathematical measurement of a task  $\mathcal{T}$ . With the assistance of information  
 114 theory, the complexity associated with a given task  $\mathcal{T}$  can be measured in terms of the entropy of  $\mathcal{T}$ .  
 115 Therefore, we can borrow the concepts of information entropy to explain the difficulty of the task.  
 116 For example, a smaller  $H(\mathcal{T})$  means an easier task and vice versa.

117 Since the entropy of task  $\mathcal{T}$  is formulated, it is not difficult to define the mutual information of task  $\mathcal{T}$   
 118 and noise  $\epsilon$ ,

$$MI(\mathcal{T}, \epsilon) = H(\mathcal{T}) - H(\mathcal{T}|\epsilon) \quad (4)$$

119 Formally, if the noise can help reduce the complexity of the task, i.e.,  $H(\mathcal{T}|\epsilon) < H(\mathcal{T})$  then the noise  
 120 has positive support. Therefore, a noise  $\epsilon$  is defined as **positive noise** (PN) when the noise satisfies  
 121  $MI(\mathcal{T}, \epsilon) > 0$ . On the contrary, when  $MI(\mathcal{T}, \epsilon) \leq 0$ , the noise is considered as the conventional  
 122 noise and named **harmful noise** (HN). The positive noise can be perceived as an augmentation of  
 123 information gain brought by  $\epsilon$ .

$$\begin{cases} MI(\mathcal{T}, \epsilon) > 0 & \epsilon \text{ is positive noise} \\ MI(\mathcal{T}, \epsilon) \leq 0 & \epsilon \text{ is harmful noise} \end{cases} \quad (5)$$

124 **Moderate Model Assumption:** The positive noise may not work for deep models with severe  
 125 problems. For example, the model is severely overfitting where models begin to memorize the  
 126 random fluctuations in the data instead of learning the underlying patterns. In that case, the presence  
 127 of positive noise will not have significant positive support in improving the models' performance.  
 128 Besides, when the models are corrupted under brute force attack, the positive noise also can not work.

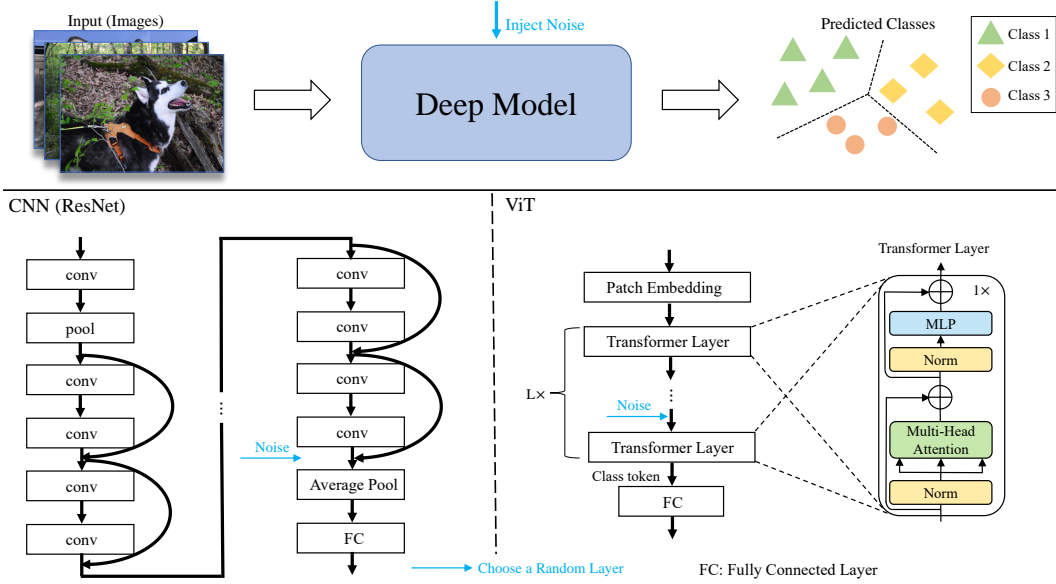


Figure 1: An overview of the proposed method. Above the black line is the standard pipeline for image classification. The deep model can be CNNs or ViTs. The noise is injected into a randomly chosen layer of the model represented by the blue arrow.

### 129 3 Methods

130 The idea of exploring the influence of noise on the deep models is straightforward. The framework is  
 131 depicted in Fig. 1. This is a universal framework where there are different options for deep models,  
 132 such as CNNs and ViTs. Through the simple operation of injecting noise into a randomly selected  
 133 layer, a model has the potential to gain additional information to reduce task complexity, thereby  
 134 improving its performance. It is sufficient to inject noise into a single layer instead of multiple layers  
 135 since it imposes a regularization on multiple layers simultaneously.

136 For a classification problem, the dataset  $(\mathbf{X}, \mathbf{Y})$  can be regarded as samplings derived from  $D_{\mathcal{X}, \mathcal{Y}}$ ,  
 137 where  $D_{\mathcal{X}, \mathcal{Y}}$  is some unknown joint distribution of data points and labels from feasible space  $\mathcal{X}$  and  
 138  $\mathcal{Y}$ , i.e.,  $(\mathbf{X}, \mathbf{Y}) \sim D_{\mathcal{X}, \mathcal{Y}}$  [31]. Hence, given a set of  $k$  data points  $\mathbf{X} = \{X_1, X_2, \dots, X_k\}$ , the label  
 139 set  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_k\}$  is regarded as sampling from  $\mathbf{Y} \sim D_{\mathcal{Y}|\mathcal{X}}$ . The complexity of  $\mathcal{T}$  on dataset  
 140  $\mathbf{X}$  is formulated as [19]:

$$H(\mathcal{T}; \mathbf{X}) = - \sum_{\mathbf{Y} \in \mathcal{Y}} p(\mathbf{Y}|\mathbf{X}) \log p(\mathbf{Y}|\mathbf{X}) \quad (6)$$

141 The operation of adding noise at the image level can be formulated as:

$$\begin{cases} H(\mathcal{T}; \mathbf{X} + \epsilon) = - \sum_{\mathbf{Y} \in \mathcal{Y}} p(\mathbf{Y}|\mathbf{X} + \epsilon) \log p(\mathbf{Y}|\mathbf{X} + \epsilon) & \epsilon \text{ is additive noise} \\ H(\mathcal{T}; \mathbf{X}\epsilon) = - \sum_{\mathbf{Y} \in \mathcal{Y}} p(\mathbf{Y}|\mathbf{X}\epsilon) \log p(\mathbf{Y}|\mathbf{X}\epsilon) & \epsilon \text{ is multiplicative noise} \end{cases} \quad (7)$$

142 While the operation of proactively injecting noise in the latent space can be formulated as:

$$\begin{cases} H(\mathcal{T}; \mathbf{X} + \epsilon) \stackrel{*}{=} H(\mathbf{Y}; \mathbf{X} + \epsilon) - H(\mathbf{X}) & \epsilon \text{ is additive noise} \\ H(\mathcal{T}; \mathbf{X}\epsilon) \stackrel{*}{=} H(\mathbf{Y}; \mathbf{X}\epsilon) - H(\mathbf{X}) & \epsilon \text{ is multiplicative noise} \end{cases} \quad (8)$$

143 Step  $\star$  differs from the conventional definition of conditional entropy, as our method injects the noise  
 144 into the latent representations instead of the original images. The Gaussian noise is additive, the  
 145 linear transform noise is also additive, and the salt-and-pepper is a multiplicative noise.

146 **Gaussian Noise** The Gaussian noise is one of the most common additive noises that appeared in  
 147 computer vision tasks. The Gaussian noise is independent and stochastic, obeying the Gaussian  
 148 distribution. Without loss of generality, defined as  $\mathcal{N}(\mu, \sigma^2)$ . Since our injection happens in the  
 149 latent space, therefore, the complexity of the task is:

$$H(\mathcal{T}; \mathbf{X} + \epsilon) \stackrel{*}{=} H(\mathbf{Y}; \mathbf{X} + \epsilon) - H(\mathbf{X}). \quad (9)$$

150 According to the definition in Equation 4, and making the distribution of  $\mathbf{X}$  and  $\mathbf{Y}$  multivariate  
 151 normal distribution [5] [14], the mutual information with Gaussian noise is:

$$\begin{aligned}
 MI(\mathcal{T}, \epsilon) &= H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{X}) - (H(\mathbf{Y}; \mathbf{X} + \epsilon) - H(\mathbf{X})) \\
 &= H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{Y}; \mathbf{X} + \epsilon) \\
 &= \frac{1}{2} \log \frac{|\Sigma_{\mathbf{X}}| |\Sigma_{\mathbf{Y}} - \Sigma_{\mathbf{Y}\mathbf{X}} \Sigma_{\mathbf{X}}^{-1} \Sigma_{\mathbf{X}\mathbf{Y}}|}{|\Sigma_{\mathbf{X}+\epsilon}| |\Sigma_{\mathbf{Y}} - \Sigma_{\mathbf{Y}\mathbf{X}} \Sigma_{\mathbf{X}+\epsilon}^{-1} \Sigma_{\mathbf{X}\mathbf{Y}}|} \\
 &= \frac{1}{2} \log \frac{1}{(1 + \sigma_{\epsilon}^2 \sum_{i=1}^k \frac{1}{\sigma_{X_i}^2}) (1 + \lambda \sum_{i=1}^k \frac{\text{cov}^2(X_i, Y_i)}{\sigma_{X_i}^2 (\sigma_{X_i}^2 \sigma_{Y_i}^2 - \text{cov}^2(X_i, Y_i))}}
 \end{aligned} \tag{10}$$

152 where  $\lambda = \frac{\sigma_{\epsilon}^2}{1 + \sum_{i=1}^k \frac{1}{\sigma_{X_i}^2}}$ ,  $\sigma_{\epsilon}^2$  is the variance of the Gaussian noise,  $\text{cov}(X_i, Y_i)$  is the covariance of  
 153 sample pair  $X_i, Y_i$ ,  $\sigma_{X_i}^2$  and  $\sigma_{Y_i}^2$  are the variance of data sample  $X_i$  and data label  $Y_i$ , respectively.  
 154 The detailed derivations can be found in section 1.1.2 of the supplementary. Given a dataset, the  
 155 variance of the Gaussian noise, and statistical properties of data samples and labels control the mutual  
 156 information, we define the function:

$$\begin{aligned}
 f(\sigma_{\epsilon}^2) &= 1 - (1 + \sigma_{\epsilon}^2 \sum_{i=1}^k \frac{1}{\sigma_{X_i}^2}) (1 + \lambda \sum_{i=1}^k \frac{\text{cov}^2(X_i, Y_i)}{\sigma_{X_i}^2 (\sigma_{X_i}^2 \sigma_{Y_i}^2 - \text{cov}^2(X_i, Y_i))}) \\
 &= -\sigma_{\epsilon}^2 \sum_{i=1}^k \frac{1}{\sigma_{X_i}^2} - \sigma_{\epsilon}^2 \sum_{i=1}^k \frac{1}{\sigma_{X_i}^2} \cdot \lambda \sum_{i=1}^k \frac{\text{cov}^2(X_i, Y_i)}{\sigma_{X_i}^2 (\sigma_{X_i}^2 \sigma_{Y_i}^2 - \text{cov}^2(X_i, Y_i))} - \lambda \sum_{i=1}^k \frac{\text{cov}^2(X_i, Y_i)}{\sigma_{X_i}^2 (\sigma_{X_i}^2 \sigma_{Y_i}^2 - \text{cov}^2(X_i, Y_i))}
 \end{aligned} \tag{11}$$

157 Since  $\epsilon^2 \geq 0$  and  $\lambda \geq 0$ ,  $\sigma_{X_i}^2 \sigma_{Y_i}^2 - \text{cov}^2(X_i, Y_i) = \sigma_{X_i}^2 \sigma_{Y_i}^2 (1 - \rho_{X_i Y_i}^2) \geq 0$ , where  $\rho_{X_i Y_i}$  is the  
 158 correlation coefficient, the sign of  $f(\sigma_{\epsilon}^2)$  is negative. We can conclude that Gaussian noise injected  
 159 into the latent space is harmful to the task. More details and the Gaussian noise added to the image  
 160 level are provided in the supplementary.

161 **Linear Transform Noise** This type of noise is obtained by elementary transformation of the features  
 162 matrix, i.e.,  $\epsilon = Q\mathbf{X}$ , where  $Q$  is an elementary matrix. We name the  $Q$  the quality matrix since it  
 163 controls the property of linear transform noise and determines whether positive or harmful. In the  
 164 linear transform noise injection in the latent space case, the complexity of the task is:

$$H(\mathcal{T}; \mathbf{X} + Q\mathbf{X}) \stackrel{*}{=} H(\mathbf{Y}; \mathbf{X} + Q\mathbf{X}) - H(\mathbf{X}) \tag{12}$$

165 The mutual information is then formulated as:

$$\begin{aligned}
 MI(\mathcal{T}, Q\mathbf{X}) &\stackrel{*}{=} H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{X}) - (H(\mathbf{Y}; \mathbf{X} + Q\mathbf{X}) - H(\mathbf{X})) \\
 &= H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{Y}; \mathbf{X} + Q\mathbf{X}) \\
 &= \frac{1}{2} \log \frac{|\Sigma_{\mathbf{X}}| |\Sigma_{\mathbf{Y}} - \Sigma_{\mathbf{Y}\mathbf{X}} \Sigma_{\mathbf{X}}^{-1} \Sigma_{\mathbf{X}\mathbf{Y}}|}{|\Sigma_{(I+Q)\mathbf{X}}| |\Sigma_{\mathbf{Y}} - \Sigma_{\mathbf{Y}\mathbf{X}} \Sigma_{\mathbf{X}}^{-1} \Sigma_{\mathbf{X}\mathbf{Y}}|} \\
 &= \frac{1}{2} \log \frac{1}{|I + Q|^2} \\
 &= -\log |I + Q|
 \end{aligned} \tag{13}$$

166 Since we want the mutual information to be greater than 0, we can formulate Equation 13 as an  
 167 optimization problem:

$$\begin{aligned}
 &\max_Q MI(\mathcal{T}, Q\mathbf{X}) \\
 &s.t. \text{rank}(I + Q) = k \\
 &\quad Q \sim I \\
 &\quad [I + Q]_{ii} \geq [I + Q]_{ij}, i \neq j \\
 &\quad \|[I + Q]_i\|_1 = 1
 \end{aligned} \tag{14}$$

168 where  $\sim$  means the row equivalence. The key to determining whether the linear transform is positive  
 169 noise or not lies in the matrix of  $Q$ . The most important step is to ensure that  $I + Q$  is reversible,

170 which is  $|(I + Q)| \neq 0$ . The third constraint is to make the trained classifier get enough information  
 171 about a specific image and correctly predict the corresponding label. For example, for an image  $X_1$   
 172 perturbed by another image  $X_2$ , the classifier obtained dominant information from  $X_1$  so that it can  
 173 predict the label  $Y_1$ . However, if the perturbed image  $X_2$  is dominant, the classifier can hardly predict  
 174 the correct label  $Y_1$  and is more likely to predict as  $Y_2$ . The fourth constraint is to maintain the norm  
 175 of latent representations. More in-depth discussion and linear transform noise added to the image  
 176 level are provided in the supplementary.

177 **Salt-and-pepper Noise** The salt-and-pepper noise is a common multiplicative noise for images. The  
 178 image can exhibit unnatural changes, such as black pixels in bright areas or white pixels in dark  
 179 areas, specifically as a result of the signal disruption caused by sudden strong interference or bit  
 180 transmission errors. In the Salt-and-pepper noise case, the mutual information is:

$$\begin{aligned}
 MI(\mathcal{T}, \epsilon) &\stackrel{*}{=} H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{X}) - (H(\mathbf{Y}; \mathbf{X}\epsilon) - H(\mathbf{X})) \\
 &= H(\mathbf{Y}; \mathbf{X}) - H(\mathbf{Y}; \mathbf{X}\epsilon) \\
 &= - \sum_{\mathbf{X} \in \mathcal{X}} \sum_{\mathbf{Y} \in \mathcal{Y}} p(\mathbf{X}, \mathbf{Y}) \log p(\mathbf{X}, \mathbf{Y}) - \sum_{\mathbf{X} \in \mathcal{X}} \sum_{\mathbf{Y} \in \mathcal{Y}} \sum_{\epsilon \in \mathcal{E}} p(\mathbf{X}\epsilon, \mathbf{Y}) \log p(\mathbf{X}\epsilon, \mathbf{Y}) \\
 &= \mathbb{E} \left[ \log \frac{1}{p(\mathbf{X}, \mathbf{Y})} \right] - \mathbb{E} \left[ \log \frac{1}{p(\mathbf{X}\epsilon, \mathbf{Y})} \right] \\
 &= \mathbb{E} \left[ \log \frac{1}{p(\mathbf{X}, \mathbf{Y})} \right] - \mathbb{E} \left[ \log \frac{1}{p(\mathbf{X}, \mathbf{Y})} \right] - \mathbb{E} \left[ \log \frac{1}{p(\epsilon)} \right] \\
 &= -H(\epsilon)
 \end{aligned} \tag{15}$$

181 Obviously, the mutual information is smaller than 0, which indicates the complexity is increasing  
 182 when injecting salt-and-pepper noise into the deep model. As can be foreseen, the salt-and-pepper  
 183 noise is pure detrimental noise. More details and Salt-and-pepper added to the image level are in the  
 184 supplementary.

## 185 4 Experiments

186 In this section, we conduct extensive experiments to explore the influence of various types of noises  
 187 on deep learning models. We employ popular deep learning architectures, including both CNNs  
 188 and ViTs, and show that the two kinds of deep models can benefit from the positive noise. We  
 189 employ deep learning models of various scales, including ViT-Tiny (ViT-T), ViT-Small (ViT-S),  
 190 ViT-Base (ViT-B), and ViT-Large (ViT-L) for Vision Transformers (ViTs), and ResNet-18, ResNet-34,  
 191 ResNet-50, and ResNet-101 for ResNet architecture. The details of deep models are presented in the  
 192 supplementary. Without specific instructions, the noise is injected at the last layer of the deep models.  
 193 Note that for ResNet models, the number of macro layers is 4, and for each macro layer, different  
 194 scale ResNet models have different micro sublayers. For example, for ResNet-18, the number of  
 195 macro layers is 4, and for each macro layer, the number of micro sublayers is 2. The noise is injected  
 196 at the last micro sublayer of the last macro layer for ResNet models. More experimental settings for  
 197 ResNet and ViT are detailed in the supplementary.

### 198 4.1 Noise Setting

199 We utilize the standard normal distribution to generate Gaussian noise in our experiments, ensuring  
 200 that the noise has zero mean and unit variance. Gaussian noise can be expressed as:

$$200 \epsilon \sim \mathcal{N}(0, 1) \tag{16}$$

201 For linear transform noise, we use a quality matrix of as:

$$201 Q = -\alpha I + \alpha f(I) \tag{17}$$

202 where  $I$  is the identity matrix,  $\alpha$  represents the linear transform strength and  $f$  is a row cyclic shift  
 203 operation switching row to the next row, for example, in a  $3 \times 3$  matrix,  $f$  will move Row 1 to Row  
 204 2, Row 2 to Row 3, and Row 3 to Row 1. For salt-and-pepper noise, we also use the parameter  $\alpha$  to  
 205 control the probability of the emergence of salt-and-pepper noise, which can be formulated as:

$$\begin{cases} \max(X) & \text{if } p < \alpha/2 \\ \min(X) & \text{if } p > 1 - \alpha/2 \end{cases} \tag{18}$$

Table 1: ResNet with different kinds of noise on ImageNet. Vanilla means the vanilla model without noise. Accuracy is shown in percentage. Gaussian noise used here is subjected to standard normal distribution. Linear transform noise used in this table is designed to be positive noise. The difference is shown in the bracket.

Model	ResNet-18	ResNet-34	ResNet-50	ResNet-101
Vanilla	63.90 (+0.00)	66.80 (+0.00)	70.00 (+0.00)	70.66 (+0.00)
+ Gaussian Noise	62.35 (-1.55)	65.40 (-1.40)	69.62 (-0.33)	70.10 (-0.56)
+ Linear Transform Noise	<b>79.62 (+15.72)</b>	<b>80.05 (+13.25)</b>	<b>81.32 (+11.32)</b>	<b>81.91 (+11.25)</b>
+ Salt-and-pepper Noise	55.45 (-8.45)	63.36 (-3.44)	45.89 (-24.11)	52.96 (-17.70)

Table 2: ViT with different kinds of noise on ImageNet. Vanilla means the vanilla model without injecting noise. Accuracy is shown in percentage. Gaussian noise used here is subjected to standard normal distribution. Linear transform noise used in this table is designed to be positive noise. The difference is shown in the bracket. Note **ViT-L is overfitting on ImageNet** [13] [34].

Model	ViT-T	ViT-S	ViT-B	ViT-L
Vanilla	79.34 (+0.00)	81.88 (+0.00)	84.33 (+0.00)	88.64 (+0.00)
+ Gaussian Noise	79.10 (-0.24)	81.80 (-0.08)	83.41 (-0.92)	85.92 (-2.72)
+ Linear Transform Noise	<b>80.69 (+1.35)</b>	<b>87.27 (+5.39)</b>	<b>89.99 (+5.66)</b>	<b>88.72 (+0.08)</b>
+ Salt-and-pepper Noise	78.64 (-0.70)	81.75 (-0.13)	82.40 (-1.93)	85.15 (-3.49)

206 where  $p$  is a probability generated by a random seed,  $\alpha \in [0, 1)$ , and  $X$  is the representation of an  
 207 image.

## 208 4.2 Image Classification Results

209 We implement extensive experiments on large-scale datasets such as ImageNet [11] and small-scale  
 210 datasets such as TinyImageNet [17] using ResNets and ViTs.

### 211 4.2.1 CNN Family

212 The results of ResNets with different noises on ImageNet are in Table 1. As shown in the table, with  
 213 the design of linear transform noise to be positive noise (PN), ResNet improves the classification  
 214 accuracy by a large margin. While the salt-and-pepper, which is theoretically harmful noise (HN),  
 215 degrades the models. Note we did not utilize data augmentation techniques for ResNet experiments  
 216 except for normalization. The significant results show that positive noise can effectively improve  
 217 classification accuracy by reducing task complexity.

### 218 4.2.2 ViT Family

219 The results of ViT with different noises on ImageNet are in Table 2. Since the ViT-L is overfitting on  
 220 the ImageNet [13] [34], the positive noise did not work well on the ViT-L. As shown in the table, the  
 221 existence of positive noise improves the classification accuracy of ViT by a large margin compared to  
 222 vanilla ViT. The comparisons with previously published works, such as DeiT [38], SwinTransformer  
 223 [20], DaViT [12], and MaxViT [39], are shown in Table 3, and our positive noise-empowered ViT  
 224 achieved the new state-of-the-art result. Note that the JFT-300M and JFT-4B datasets are private and  
 225 not publicly available [35], and we believe that ViT large and above will benefit from positive noise  
 226 significantly if trained on larger JFT-300M or JFT-4B, which is theoretically supported in section 4.4.

## 227 4.3 Ablation Study

228 We also proactively inject noise into variants of ViT, such as DeiT [38], Swin Transformer [20],  
 229 BEiT [3], and ConViT [9], and the results show that positive noise could benefit various variants  
 230 of ViT by improving classification accuracy significantly. The results of injecting noise to variants  
 231 of ViT are reported in the supplementary. We also did ablation studies on the strength of linear  
 232 transform noise and the injected layer. The results are shown in Fig. 2. We can observe that the  
 233 deeper layer the positive noise injects, the better prediction performance the model can obtain. There  
 234 are reasons behind this phenomenon. First, the latent features of input in the deeper layer have better

Table 3: Comparison between Positive Noise Empowered ViT with other ViT variants. Top 1 Accuracy is shown in percentage. Here PN is the positive noise, i.e., linear transform noise.

Model	Top1 Acc.	Params.	Image Res.	Pretrained Dataset
ViT-B [13]	84.33	86M	224 × 224	ImageNet 21k
DeiT-B [38]	85.70	86M	224 × 224	ImageNet 21k
SwinTransformer-B [20]	86.40	88M	384 × 384	ImageNet 21k
DaViT-B [12]	86.90	88M	384 × 384	ImageNet 21k
MaxViT-B [39]	88.82	119M	512 × 512	JFT-300M (Private)
ViT-22B [10]	89.51	21743M	224 × 224	JFT-4B (Private)
ViT-B+PN	<b>89.99</b>	86M	224 × 224	ImageNet 21k
ViT-B+PN	<b>91.37</b>	86M	384 × 384	ImageNet 21k

235 representations than those in shallow layers; second, injection to shallow layers obtain less mutual  
 236 information gain because of trendy replacing Equation 8 with Equation 7. More results on the small  
 237 dataset TinyImageNet can be found in the supplementary.

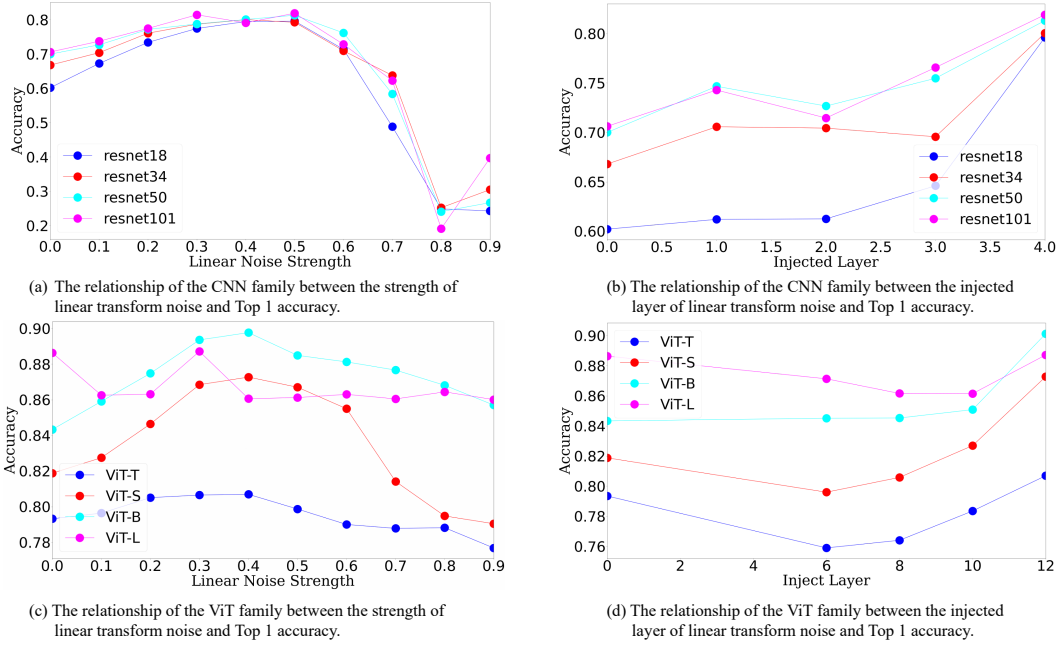


Figure 2: The relationship between the linear transform noise strength and the top 1 accuracy, and between the injected layer and top 1 accuracy. Parts (a) and (b) are the results of the CNN family, while parts (c) and (d) are the results of the ViT family. For parts (a) and (c) the linear transform noise is injected at the last layer. For parts (b) and (d), the influence of positive noise on different layers is shown. Layers 6, 8, 10, and 12 in the ViT family are chosen for the ablation study.

#### 238 4.4 Optimal Quality Matrix

239 As shown in Equation 14, it is interesting to learn about the optimal quality matrix of  $Q$  that maximizes  
 240 the mutual information while satisfying the constraints. This equals minimizing the determinant of  
 241 the matrix sum of  $I$  and  $Q$ . Here, we directly give out the optimal quality matrix of  $Q$  as:

$$Q_{optimal} = \text{diag} \left( \frac{1}{k+1} - 1, \dots, \frac{1}{k+1} - 1 \right) + \frac{1}{k+1} \mathbf{1}_{k \times k} \quad (19)$$

242 where  $k$  is the number of data samples. And the corresponding upper boundary of the mutual  
 243 information as:

$$MI(\mathcal{T}, Q_{optimal} \mathbf{X}) = (k-1) \log(k+1) \quad (20)$$



Table 4: Top 1 accuracy on ImageNet with the optimal quality matrix of linear transform noise.

Model	Top1 Acc.	Params.	Image Res.	Pretrained Dataset
ViT-B+Optimal Q	<b>93.87</b>	86M	224 × 224	ImageNet 21k
ViT-B+Optimal Q	<b>95.12</b>	86M	384 × 384	ImageNet 21k

Table 5: Comparison with various ViT-based methods on **Office-Home**.

Method	Ar2Cl	Ar2Pr	Ar2Re	Cl2Ar	Cl2Pr	Cl2Re	Pr2Ar	Pr2Cl	Pr2Re	Re2Ar	Re2Cl	Re2Pr	Avg.
ViT-B[13]	54.7	83.0	87.2	77.3	83.4	85.6	74.4	50.9	87.2	79.6	54.8	88.8	75.5
TVT-B[44]	74.9	86.8	89.5	82.8	88.0	88.3	79.8	71.9	90.1	85.5	74.6	90.6	83.6
CDTrans-B[43]	68.8	85.0	86.9	81.5	87.1	87.3	79.6	63.3	88.2	82.0	66.0	90.6	80.5
SSRT-B [36]	75.2	89.0	91.1	85.1	88.3	90.0	85.0	74.2	91.3	85.7	78.6	91.8	85.4
ViT-B+PN	<b>78.3</b>	<b>90.6</b>	<b>91.9</b>	<b>87.8</b>	<b>92.1</b>	<b>91.9</b>	<b>85.8</b>	<b>78.7</b>	<b>93.0</b>	<b>88.6</b>	<b>80.6</b>	<b>93.5</b>	<b>87.7</b>

244 The details are provided in the supplementary. We find that the upper boundary of the mutual  
 245 information of injecting positive noise is determined by the number of data samples, i.e., the scale of  
 246 the dataset. Therefore, the larger the dataset, the better effect of injecting positive noise into deep  
 247 models. With the optimal quality matrix and the top 1 accuracy of ViT-B on ImageNet can be further  
 248 improved to 95%, which is shown in Table 4.

Table 6: Comparison with various ViT-based methods on **Visda2017**.

Method	plane	bcycl	bus	car	horse	knife	mcycl	person	plant	sktbrd	train	truck	Avg.
ViT-B[13]	97.7	48.1	86.6	61.6	78.1	63.4	94.7	10.3	87.7	47.7	94.4	35.5	67.1
TVT-B[44]	92.9	85.6	77.5	60.5	93.6	98.2	89.4	76.4	93.6	92.0	91.7	55.7	83.9
CDTrans-B[43]	97.1	90.5	82.4	77.5	96.6	96.1	93.6	<b>88.6</b>	<b>97.9</b>	86.9	90.3	62.8	88.4
SSRT-B [36]	<b>98.9</b>	87.6	<b>89.1</b>	<b>84.8</b>	98.3	<b>98.7</b>	<b>96.3</b>	81.1	94.9	97.9	94.5	43.1	88.8
ViT-B+PN	98.8	<b>95.5</b>	84.8	73.7	<b>98.5</b>	97.2	95.1	76.5	95.9	<b>98.4</b>	<b>98.3</b>	<b>67.2</b>	<b>90.0</b>

## 249 4.5 Domain Adaption Results

250 Unsupervised domain adaptation (UDA) aims to learn transferable knowledge across the source and  
 251 target domains with different distributions [25] [42]. Recently, transformer-based methods achieved  
 252 SOTA results on UDA, therefore, we evaluate the ViT-B with the positive noise on widely used  
 253 UDA benchmarks. Here the positive noise is the linear transform noise identical to that used in the  
 254 classification task. The positive noise is injected into the last layer of the model, the same as the  
 255 classification task. The datasets include **Office Home** [41] and **VisDA2017** [26]. Detailed datasets  
 256 introduction and experiments training settings are in the supplementary. The objective function  
 257 is borrowed from TVT [44], which is the first work that adopts Transformer-based architecture  
 258 for UDA. The results are shown in Table 5 and 6. The ViT-B with positive noise achieves better  
 259 performance than the existing works. These results show that positive noise can improve model  
 260 generality, therefore, benefit deep models in domain adaptation tasks.

## 261 5 Conclusion

262 This study presents a comprehensive investigation into the influence of various common noise types  
 263 on deep learning models, including Gaussian noise, linear transform noise, and salt-and-pepper noise.  
 264 We demonstrate that, under certain conditions, linear transform noise can have a positive effect on  
 265 deep models. Our experiments show that injecting the positive noise in latent space can significantly  
 266 enhance the prediction performance of deep models on image classification tasks, leading to new  
 267 state-of-the-art results on ImageNet. The findings of this study have a broad impact on future research  
 268 and may contribute to the development of more accurate models and their improved performance in  
 269 real-world applications. Moreover, we are excited to explore the potential of positive noise in more  
 270 deep learning tasks.

271 **References**

- 272 [1] Osama K. Al-Shaykh and Russell M. Mersereau. Lossy compression of noisy images. *IEEE*  
273 *Transactions on Image Processing*, 7(12):1641–1652, 1998.
- 274 [2] Wissam A. Albukhanajer, Johann A. Briffa, and Yaochu Jin. Evolutionary multiobjective image  
275 feature extraction in the presence of noise. *IEEE Transactions on Cybernetics*, 45(9):1757–1768,  
276 2014.
- 277 [3] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT pre-training of image transformers. *arXiv*  
278 *preprint arXiv:2106.08254*, 2021.
- 279 [4] Roberto Benzi, Alfonso Sutera, and Angelo Vulpiani. The mechanism of stochastic resonance.  
280 *Journal of Physics A: mathematical and general*, 14(11):L453, 1981.
- 281 [5] George EP Box and David R. Cox. An analysis of transformations. *Journal of the Royal*  
282 *Statistical Society: Series B (Methodological)*, 26(2):211–243, 1964.
- 283 [6] Sebastian Braun, Hannes Gamper, Chandan KA Reddy, and Ivan Tashev. Towards efficient mod-  
284 els for real-time deep noise suppression. In *ICASSP 2021-2021 IEEE International Conference*  
285 *on Acoustics, Speech and Signal Processing (ICASSP)*, pages 656–660, 2021.
- 286 [7] Raymond H. Chan, Chung-Wa Ho, and Mila Nikolova. Salt-and-pepper noise removal by  
287 median-type noise detectors and detail-preserving regularization. *IEEE Transactions on image*  
288 *processing*, 14(10):1479–1485, 2005.
- 289 [8] Thomas M. Cover. Elements of information theory. *John Wiley & Sons*, 1999.
- 290 [9] Stéphane d’Ascoli, Hugo Touvron, Matthew Leavitt, Ari Morcos, Giulio Biroli, and Levent  
291 Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. *arXiv*  
292 *preprint arXiv:2103.10697*, 2021.
- 293 [10] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin  
294 Gilmer, Andreas Steiner, and et al. Scaling vision transformers to 22 billion parameters. *arXiv*  
295 *preprint arXiv:2302.05442 (2023)*, 2023.
- 296 [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Feifei Li. Imagenet: A large-scale  
297 hierarchical image database. In *IEEE conference on computer vision and pattern recognition*,  
298 pages 248–255, 2009.
- 299 [12] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jindong Wang, and Lu Yuan. Davit: Dual  
300 attention vision transformers. In *In Computer Vision–ECCV 2022: 17th European Conference*,  
301 pages 74–92, 2022.
- 302 [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
303 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
304 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
305 recognition at scale. In *arXiv preprint arXiv:2010.11929*, 2020.
- 306 [14] Changyong Feng, Hongyue Wang, Naiji Lu, Tian Chen, Hua He, Ying Lu, and Xin M. Tu.  
307 Log-transformation and its implications for data analysis. *Shanghai archives of psychiatry*,  
308 26(2):105, 2014.
- 309 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
310 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
311 pages 770–778, 2016.
- 312 [16] Solomon Kullback and Richard A. Leibler. On information and sufficiency. *The annals of*  
313 *mathematical statistics*, 22(1):79–86, 1951.
- 314 [17] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N 7*, (7), 2015.
- 315 [18] Yann LeCun and Yoshua Bengio. Convolutional networks for images, speech, and time series.  
316 *The handbook of brain theory and neural networks*, 3361(10), 1995.

- 317 [19] Xuelong Li. Positive-incentive noise. *IEEE Transactions on Neural Networks and Learning*  
318 *Systems*, 2022.
- 319 [20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
320 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*  
321 *of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- 322 [21] Peter McClintock. Can noise actually boost brain power? *Physics World*, 15(7), 2002.
- 323 [22] Toshio Mori and Shoichi Kai. Noise-induced entrainment and stochastic resonance in human  
324 brain waves. *Physical review letters*, 88(21), 2002.
- 325 [23] Rich Ormiston, Tri Nguyen, Michael Coughlin, Rana X. Adhikari, and Erik Katsavounidis.  
326 Noise reduction in gravitational-wave data via deep learning. *Physical Review Research*,  
327 2(3):033066, 2020.
- 328 [24] J. S. Owotogbe, T. S. Ibiyemi, and B. A. Adu. A comprehensive review on various types of  
329 noise in image processing. *int. J. Sci. eng. res*, 10(10):388–393, 2019.
- 330 [25] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on*  
331 *knowledge and data engineering*, 22(10):1345–1359, 2009.
- 332 [26] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko.  
333 Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- 334 [27] Lis Kanashiro Pereira, Yuki Taya, and Ichiro Kobayashi. Multi-layer random perturbation  
335 training for improving model generalization efficiently. *Proceedings of the Fourth BlackboxNLP*  
336 *Workshop on Analyzing and Interpreting Neural Networks for NLP*, 2021.
- 337 [28] Kamiar Radnosrati, Gustaf Hendeby, and Fredrik Gustafsson. Crackling noise. *IEEE Transac-*  
338 *tions on Signal Processing*, 68:3590–3602, 2020.
- 339 [29] Fabrizio Russo. A method for estimation and filtering of gaussian noise in images. *IEEE*  
340 *Transactions on Instrumentation and Measurement*, 52(4):1148–1154, 2003.
- 341 [30] James P. Sethna, Karin A. Dahmen, and Christopher R. Myers. Crackling noise. *Nature*,  
342 410(6825):242–250, 2001.
- 343 [31] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*  
344 *algorithms*. Cambridge university press, Cambridge, 2014.
- 345 [32] Claude Elwood Shannon. A mathematical theory of communication. *ACM SIGMOBILE mobile*  
346 *computing and communications review*, 5(1):3–55, 2001.
- 347 [33] Jan Sijbers, Paul Scheunders, Noel Bonnet, Dirk Van Dyck, and Erik Raman. Quantification and  
348 improvement of the signal-to-noise ratio in a magnetic resonance image acquisition procedure.  
349 *Magnetic resonance imaging*, 14(10):1157–1163, 1996.
- 350 [34] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit,  
351 and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision  
352 transformers. In *arXiv preprint arXiv:2106.10270*, 2021.
- 353 [35] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable  
354 effectiveness of data in deep learning era. In *In Proceedings of the IEEE international conference*  
355 *on computer vision*, pages 843–852, 2017.
- 356 [36] Tao Sun, Cheng Lu, Tianshuo Zhang, and Harbin Ling. Safe self-refinement for transformer-  
357 based domain adaptation. *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
358 *Pattern Recognition*, pages 7191–7200, 2022.
- 359 [37] Sunil Thulasidasan, Tanmoy Bhattacharya, Jeff Bilmes, Gopinath Chennupati, and Jamal  
360 Mohd-Yusof. Combating label noise in deep learning using abstention. In *arXiv preprint*  
361 *arXiv:1905.10964*, 2019.

- 362 [38] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and  
363 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In  
364 *International conference on machine learning*, pages 10347–10357, 2021.
- 365 [39] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and  
366 Yinxiao Li. Maxvit: Multi-axis vision transformer. In *In Computer Vision–ECCV 2022: 17th*  
367 *European Conference*, pages 459–479, 2022.
- 368 [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
369 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*  
370 *processing systems*, 2017.
- 371 [41] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.  
372 Deep hashing network for unsupervised domain adaptation. *CVPR*, pages 5018–5027, 2017.
- 373 [42] Ying Wei, Yu Zhang, Junzhou Huang, and Qiang Yang. Transfer learning via learning to transfer.  
374 *ICML*, pages 5085–5094, 2018.
- 375 [43] Tongkun Xu, Weihua Chen, Fan Wang, Hao Li, and Rong Jin. Cdtrans: Cross-domain trans-  
376 former for unsupervised domain adaptation. *ICLR*, pages 520–530, 2022.
- 377 [44] Jinyu Yang, Jingjing Liu, Ning Xu, and Junzhou Huang. Tvt: Transferable vision transformer  
378 for unsupervised domain adaptation. *WACV*, pages 520–530, 2023.
- 379 [45] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transform-  
380 ers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
381 pages 12104–12113, 2022.

## 382 Checklist

383 The checklist follows the references. Please read the checklist guidelines carefully for information on  
384 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
385 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
386 the appropriate section of your paper or providing a brief inline description. For example:

- 387 • Did you include the license to the code and datasets? **[Yes]** Yes  
388 • Did you include the license to the code and datasets? **[No]** The code and the data are  
389 proprietary.  
390 • Did you include the license to the code and datasets? **[N/A]**

391 Please do not modify the questions and only use the provided macros for your answers. Note that the  
392 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
393 block and only keep the Checklist section heading above along with the questions/answers below.

- 394 1. For all authors...
- 395 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
396 contributions and scope? **[TODO]** Yes  
397 (b) Did you describe the limitations of your work? **[TODO]** Yes  
398 (c) Did you discuss any potential negative societal impacts of your work? **[TODO]** No  
399 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
400 them? **[TODO]** Yes
- 401 2. If you are including theoretical results...
- 402 (a) Did you state the full set of assumptions of all theoretical results? **[TODO]** Yes  
403 (b) Did you include complete proofs of all theoretical results? **[TODO]** Yes
- 404 3. If you ran experiments...
- 405 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
406 mental results (either in the supplemental material or as a URL)? **[TODO]** Yes

- 407 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
408 were chosen)? **[TODO]** Yes
- 409 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
410 ments multiple times)? **[TODO]** No
- 411 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
412 of GPUs, internal cluster, or cloud provider)? No **[TODO]**
- 413 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 414 (a) If your work uses existing assets, did you cite the creators? **[TODO]** N/A
- 415 (b) Did you mention the license of the assets? **[TODO]** N/A
- 416 (c) Did you include any new assets either in the supplemental material or as a URL?  
417 **[TODO]** N/A
- 418 (d) Did you discuss whether and how consent was obtained from people whose data you're  
419 using/curating? **[TODO]** N/A
- 420 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
421 information or offensive content? **[TODO]** N/A
- 422 5. If you used crowdsourcing or conducted research with human subjects...
- 423 (a) Did you include the full text of instructions given to participants and screenshots, if  
424 applicable? **[TODO]** N/A
- 425 (b) Did you describe any potential participant risks, with links to Institutional Review  
426 Board (IRB) approvals, if applicable? **[TODO]** N/A
- 427 (c) Did you include the estimated hourly wage paid to participants and the total amount  
428 spent on participant compensation? **[TODO]** N/A