

---

# Estimating Transition Matrix with Diffusion Models for Instance-Dependent Label Noise

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Learning with noisy labels is a common problem in weakly supervised learning,  
2 where the transition matrix approach is a prevalent method for dealing with label  
3 noise. It estimates the transition probabilities from a clean label distribution to a  
4 noisy label distribution and has garnered continuous attention. However, existing  
5 transition matrix methods predominantly focus on class-dependent noise, making  
6 it challenging to incorporate feature information for learning instance-dependent  
7 label noise. This paper proposes the idea of using diffusion models for estimating  
8 transition matrix in the context of instance-dependent label noise. Specifically, we  
9 first estimate grouped transition matrices through clustering. Then, we introduce  
10 a process of adding noise and denoising with the transition matrix, incorporating  
11 features extracted by unsupervised pre-trained models. The proposed method  
12 enables the estimation of instance-dependent transition matrix and extends the  
13 application of transition matrix method to a broader range of noisy label data.  
14 Experimental results demonstrate the significant effectiveness of our approach on  
15 both synthetic and real-world datasets with instance-dependent noise. The code  
16 will be open sourced upon acceptance of the paper.

## 17 1 Introduction

18 For classification problems with given labels, deep neural networks have demonstrated significant  
19 improvements compared to traditional methods in recent years [25]. The efficacy of deep neural  
20 networks heavily relies on the accuracy of the labels. Directly incorporating polluted erroneous labels  
21 into network learning can result in the network fitting the noise, potentially severely impacting the  
22 predictive performance of the network [8]. However, in reality, obtaining accurate annotated data can  
23 be prohibitively expensive, and a substantial amount of data comes from the Internet or is annotated  
24 by non-expert annotators, inevitably containing noisy labels. Therefore, researching and promoting  
25 methods to mitigate the damage to models and make them more robust in the face of label noise data  
26 is a highly worthwhile problem to investigate, known as the problem of learning with noisy labels  
27 [23, 10, 34, 1].

28 Different approaches have been proposed to address the problem of label noise. One category  
29 [31, 22] involves the design of specialized loss functions or network structures to enhance the model's  
30 robustness against noisy labels. Another major category focuses on sample selection [2, 10, 14],  
31 where samples are partitioned into a set of clean samples and a set of contaminated noisy samples  
32 based on the magnitude of the loss or the similarity of extracted features. The labels of the noisy  
33 samples are then modified or their weights are reduced, followed by learning using semi-supervised  
34 methods. Sample selection methods are currently mainstream and have achieved promising results.  
35 However, the selection process relies heavily on intuition and lacks theoretical support. Additionally,  
36 the sample selection procedure is often complex and computationally intensive. In contrast, another

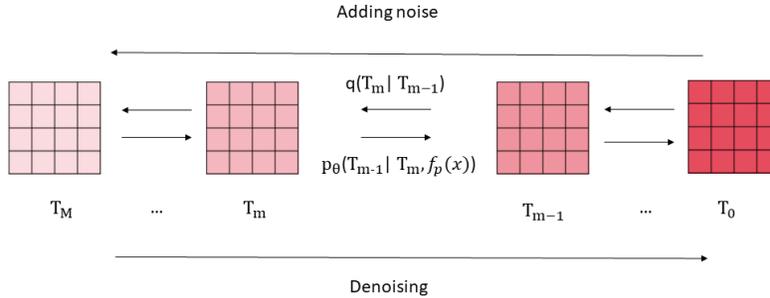


Figure 1: Diffusion Model for Transition Matrix.

37 significant category of methods is the transition matrix method [34, 17, 12, 42], which estimates the  
 38 transition probabilities from the clean label distribution to the noisy label distribution. This class  
 39 of methods reveals the generation process of noisy labels and exhibits statistical consistency, often  
 40 accompanied by theoretical analyses as methodological support. As a result, they have garnered  
 41 continuous attention and occupy an important position in various algorithms for learning with noisy  
 42 labels.

43 In transition matrix methods, accurate estimation of the transition matrix is crucial. If an accurate  
 44 estimation of the transition matrix can be obtained, along with the observed data for estimating the  
 45 posterior distribution of the noisy labels, it is possible to infer the distribution of clean labels for  
 46 neural network learning. Previous transition matrix methods [34, 17, 39] have mainly focused on  
 47 class-dependent label noise, where a single transition matrix is estimated for all samples, which is  
 48 typically straightforward. However, for instance-dependent label noise and complex real-world data,  
 49 the label transition probabilities for each sample are not entirely identical. The transition matrix often  
 50 depends on the specific features of individual samples, requiring the estimation of a separate transition  
 51 matrix for each sample. However, in most cases, a single observed label corresponds to each sample  
 52 in the dataset, making it an identifiability problem to estimate a separate transition matrix for each  
 53 sample [20]. Although some methods [33, 41, 15] have utilized separate small networks to generate  
 54 the transition matrix or divided the data into groups to transform it into a grouped class-dependent  
 55 scenario, there still exist significant estimation errors and a lack of incorporating features effectively  
 56 into the estimation of the transition matrix.

57 To better incorporate the feature information of images into the estimation of the transition matrix,  
 58 this work employs conditional diffusion models. The diffusion model originates from generative  
 59 models and has been widely applied in various computer vision tasks in recent years [36, 7], showing  
 60 remarkable results. The proposed method revolves around the core idea of replacing image samples in  
 61 the original diffusion process with a transition matrix. The matrix undergoes a process of adding noise  
 62 and denoising, where the denoising step incorporates the sample features extracted by a pre-trained  
 63 model as conditions. This generates a feature-dependent transition matrix. The constructed diffusion  
 64 module is illustrated in Figure 1. Additionally, considering the assumption that instance-dependent  
 65 label noise is usually correlated with features [6], clustering methods are utilized at the feature level  
 66 to group samples. Preliminary estimations of the transition matrices are obtained for each group,  
 67 which are then incorporated into the diffusion module for learning. The overall framework of the  
 68 method is depicted in Figure 2.

69 The subsequent sections are organized as follows. Section 2 presents an in-depth review of the  
 70 relevant works. In Section 3, we introduce our proposed model framework. Section 4 outlines  
 71 the experimental analysis conducted on diverse synthetic and real-world noisy datasets, along with  
 72 comparisons against other existing methods. Finally, we provide concluding in Section 5. The  
 73 primary contributions of this paper can be summarized as follows:

- 74 • We propose a method that utilizes diffusion models to add noise and denoise on the transition  
 75 matrix, incorporating image features extracted through pre-trained encoder.
- 76 • By combining the transition matrix-based diffusion model with feature-based clustering, we  
 77 establish a framework capable of addressing instance-dependent label noise problems.

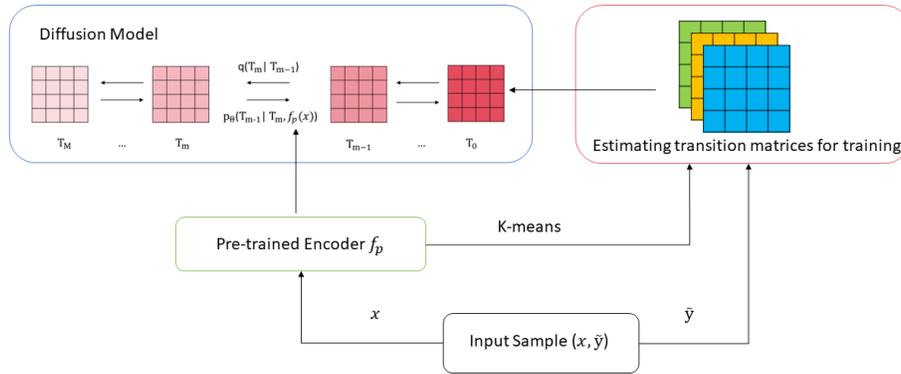


Figure 2: The overall framework of DTM.

- 78 • Our method demonstrates significant improvements over other transition matrix methods on  
 79 both synthetic and real-world noisy datasets, and it achieves comparable performance to  
 80 state-of-the-art methods.

## 81 2 Related Works

### 82 2.1 Transition Matrix Methods

83 Most previous methods for estimating transition matrix in the presence of label noise have primarily  
 84 focused on class-dependent noise scenarios, simplifying the estimation process. Methods such as  
 85 [24, 34] assume the existence of anchor points to identify the transition matrix. [17] and [39]  
 86 introduce different regularization techniques to relax the anchor point assumption. Additionally,  
 87 [26, 38] apply techniques such as meta-learning to estimate the transition matrix, but these approaches  
 88 may require more clean data and computational resources. While these methods are effective for  
 89 handling class-dependent label noise, they are not suitable for instance-dependent noise or real-world  
 90 noisy data.

91 However, estimating an individual transition matrix for each sample without additional assumptions  
 92 or multiple noisy labels is infeasible [20]. To approximate the estimation of the instance-dependent  
 93 transition matrix, [9] utilize an adaptation layer that estimates the transition matrix based on the  
 94 output of each sample. [37] employs a separate network to estimate the transition matrix based on  
 95 Bayesian labels. Some methods, such as [33, 30, 41], employ clustering to learn part-dependent  
 96 or group-dependent matrices, which can be viewed as a compromise between instance-dependent  
 97 and class-dependent methods. Other approaches, including [6, 12], utilize the similarity in the  
 98 feature space to aid in learning the transition matrix. Although these instance-dependent transition  
 99 matrix methods achieve identifiability through specialized treatments, they have not effectively  
 100 utilized feature information in the learning process, resulting in errors in estimating feature-dependent  
 101 transition matrices.

### 102 2.2 Diffusion Models

103 Diffusion models, as generative models, have played a significant role in computer vision [36, 7].  
 104 Prominent examples include DDPM [11], DDIM [27], score matching methods [28], and methods  
 105 based on stochastic differential equations [29]. Diffusion models and their variants have been applied  
 106 to various computer vision tasks such as image generation, image-to-image translation, text-to-image  
 107 generation, among others. However, their application to the problem of label noise is relatively novel.  
 108 To the best of our knowledge, only one existing work [3] has utilized diffusion models for addressing  
 109 this problem. However, this work treats labels as the output of the diffusion model, which limits  
 110 their expressive power due to the low dimension of the labels. Moreover, it overly relies on directly  
 111 incorporating image features as conditions in the label generation process, which depends heavily on

112 pre-trained models and may not be as reasonable as incorporating them into the transition matrix that  
 113 reveals the process of noise generation. Experimental results also support this perspective.

### 114 3 Method

115 In this section, we present the definitions of symbols and introduce our method of using **D**iffusion  
 116 models to construct the **T**ransition **M**atrix (DTM).

#### 117 3.1 Preliminaries

118 Let  $\mathcal{X} \subset \mathbb{R}^d$  be the input image space,  $\mathcal{Y} = \{1, 2, \dots, C\}$  be the label space, where  $C$  is the number  
 119 of classes. Random variables  $(X, Y), (X, \tilde{Y}) \in \mathcal{X} \times \mathcal{Y}$  denote the underlying data distributions  
 120 with true and noisy labels respectively. In general, we can not observe the latent true data samples  
 121  $\mathbb{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , but can only obtain the corrupted data  $\tilde{\mathbb{D}} = \{(\mathbf{x}_i, \tilde{y}_i)\}_{i=1}^N$ , where  $\tilde{y} \in \mathcal{Y}$  is the  
 122 noisy label corrupted from the true label  $y$ , while denote corresponding one-hot label as  $\mathbf{y}$  and  $\tilde{\mathbf{y}}$ .

123 Transition matrix methods use a matrix  $\mathbf{T}(\mathbf{x}) \in [0, 1]^{C \times C}$  to represent the probability from clean  
 124 label to noisy label, where the  $ij$ -th entry of the transition matrix is the probability that the instance  
 125  $\mathbf{x}$  with the clean label  $i$  corrupted to a noisy label  $j$ . The matrix satisfies the requirement that the  
 126 sum of each row  $\sum_{j=1}^C \mathbf{T}_{ij}(\mathbf{x})$  is 1, and usually has the requirement for  $\mathbf{T}_{ii}(\mathbf{x}) > \mathbf{T}_{ij}(\mathbf{x}), \forall j \neq i$ .  
 127 Let  $P(\mathbf{Y}|X = \mathbf{x}) = [P(Y = 1|X = \mathbf{x}), \dots, P(Y = C|X = \mathbf{x})]^\top$  be the clean class-posterior  
 128 probability and  $P(\tilde{\mathbf{Y}}|X = \mathbf{x}) = [P(\tilde{Y} = 1|X = \mathbf{x}), \dots, P(\tilde{Y} = C|X = \mathbf{x})]^\top$  be the noisy  
 129 class-posterior probability, the formula can be write as:

$$P(\tilde{\mathbf{Y}}|X = \mathbf{x}) = \mathbf{T}(\mathbf{x})^\top P(\mathbf{Y}|X = \mathbf{x}). \quad (1)$$

130 By estimating the transition matrix and the noisy class-posterior probability, the clean class-posterior  
 131 probability can be inferred by

$$P(\mathbf{Y}|X = \mathbf{x}) = \mathbf{T}(\mathbf{x})^{-\top} P(\tilde{\mathbf{Y}}|X = \mathbf{x}), \quad (2)$$

132 where the symbol  $-\top$  denotes the transpose of the inverse matrix.

133 The majority of existing methods [24, 10, 17] focus on studying the class-dependent and instance-  
 134 independent transition matrix, i.e.,  $\mathbf{T}(\mathbf{x}) \equiv \mathbf{T}$  for  $\forall \mathbf{x}$ . However, these methods are not applicable to  
 135 instance-dependent noise scenarios where the transition matrix  $\mathbf{T}(\mathbf{x})$  varies with respect to the input  
 136  $X$ . The main focus of our work is to utilize the feature information from input images to construct a  
 137 instance-dependent transition matrix  $\mathbf{T}(\mathbf{x})$ .

#### 138 3.2 Diffusion Model for Transition Matrix

139 We adopt the classic DDPM model [11] from diffusion models as a reference to perform noise  
 140 addition and denoising on the transition matrix. The diagram is illustrated in Figure 1.

141 For the forward diffusion process beginning with transition matrix  $\mathbf{T}_0 \sim q(\mathbf{T})$ , the process of  
 142 gradually adding noise is obtained according to the following Markov process:

$$q(\mathbf{T}_m | \mathbf{T}_{m-1}) = \mathcal{N}\left(\mathbf{T}_m; \sqrt{1 - \beta_m} \mathbf{T}_{m-1}, \beta_m \mathbf{I}\right), \quad (3)$$

143 for  $m = 1, 2, \dots, M$ , where we use  $M$  to replace  $T$ , which is usually used in other diffusion models,  
 144 in above equation for distinguishing from the symbol of transition matrix  $\mathbf{T}$ .

145 We aim to make the distribution of  $q(\mathbf{T}_M)$  approach a standard normal distribution  $\mathcal{N}(0, \mathbf{I})$  and  
 146 through  $\mathbf{T}_M$  to conduct the reverse denoising process by fitting a neural network  $\mu_\theta$  to fit the  
 147 distttribution:

$$p_\theta(\mathbf{T}_{m-1} | \mathbf{T}_m) = \mathcal{N}\left(\mathbf{T}_{m-1}; \mu_\theta(\mathbf{T}_m, \mathbf{x}, f_p, m), \tilde{\beta}_m \mathbf{I}\right), \quad (4)$$

148 where define  $\tilde{\beta}_m = \frac{1 - \bar{\alpha}_{m-1}}{1 - \bar{\alpha}_m} \beta_m, \alpha_m = 1 - \beta_m, \bar{\alpha}_m = \prod_{i=1}^m \alpha_i$ . The  $f_p$  in equation (4) denotes the  
 149 pre-trained encoder for feature extraction.

150 The diffusion model can be learned by optimizing the evidence lower bound:

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_q \left[ \mathcal{L}_M + \sum_{m>1}^M \mathcal{L}_{m-1} + \mathcal{L}_0 \right], \quad (5)$$

151 where

$$\begin{aligned} \mathcal{L}_0 &= -\log p_\theta(\mathbf{T}_0 | \mathbf{T}_1), \\ \mathcal{L}_{m-1} &= D_{\text{KL}}(q(\mathbf{T}_{m-1} | \mathbf{T}_m, \mathbf{T}_0) \| p_\theta(\mathbf{T}_{m-1} | \mathbf{T}_m)), \\ \mathcal{L}_M &= D_{\text{KL}}(q(\mathbf{T}_M | \mathbf{T}_0) \| p_\theta(\mathbf{T}_M)). \end{aligned} \quad (6)$$

152 Similar to the derivation and simplification process of DDPM, when a pre-trained encoder  $f_p$  is  
 153 provided along with the training data incorporating the initial transition matrix  $\mathbf{T}$ , the learning  
 154 algorithm for the diffusion model is presented in Algorithm 1.

---

**Algorithm 1** Diffusion Model for Transition Matrix

---

**Input:** Training data  $\{\mathbf{x}_i, \mathbf{T}_i\}_{i=1}^N$ , pre-trained encoder  $f_p$ .

**while** not converged **do**

  Sample  $(\mathbf{x}_0, \mathbf{T}_0)$  from data

  Sample  $m \sim \{1, \dots, M\}$

  Sample noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$

  Take gradient descent step on the loss:

$$\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_m} \mathbf{T}_0 + \sqrt{1 - \bar{\alpha}_m} \epsilon, \mathbf{x}_0, f_p, m) \right\|^2$$

**end while**

---

155 Next, for each image  $\mathbf{x}$ , we can sample the corresponding transition matrix  $\mathbf{T}(\mathbf{x})$  as shown in  
 156 Algorithm 2.

---

**Algorithm 2** Sample for Transition Matrix

---

  Sample  $\mathbf{T}_M \sim \mathcal{N}(0, \mathbf{I})$

**for**  $m = M, \dots, 1$  **do**

$z \sim \mathcal{N}(0, \mathbf{I})$  if  $t > 1$ , else  $z = \mathbf{0}$

$$\mathbf{T}_{m-1} = \frac{1}{\sqrt{\bar{\alpha}_m}} \left( \mathbf{T}_m - \frac{1 - \bar{\alpha}_m}{\sqrt{1 - \bar{\alpha}_m}} \epsilon_\theta(\mathbf{T}_m, \mathbf{x}, f_p, m) \right) + \sigma_m z$$

**end for**

**Output:**  $\mathbf{T}_0$

---

### 157 3.3 Feature-Dependent Framework

158 From Algorithm 1, it can be observed that there are two components of the diffusion process that  
 159 need to be provided in advance: the pre-trained encoder  $f_p$  and the initial input  $\mathbf{T}(\mathbf{x})$ .

160 The pre-trained encoder  $f_p$  can be obtained through self-supervised learning or directly using the  
 161 large model like CLIP. In our experiments, we employ the commonly used SimCLR [4] method in  
 162 contrastive learning as the feature extraction model.

163 On the other hand, the part involving the transition matrix  $\mathbf{T}(\mathbf{x})$  used for learning the diffusion  
 164 model is also related to the pre-trained encoder  $f_p$ . Based on the assumption that the noise transition  
 165 probability depends on image features, we adopt a group-dependent transition matrix as the initial  
 166 input. We perform clustering algorithms at the feature extraction level  $f_p(\mathbf{x})$ , using the K-means  
 167 method in our experiments, to group the image data. Then, based on the method VolMinNet [17], we  
 168 train class-dependent transition matrices for each group and obtain the initial transition matrix  $\mathbf{T}(\mathbf{x})$   
 169 for each image  $\mathbf{x}$ , which is then used as input in Algorithm 1. It is worth to note that the initial  $\mathbf{T}(\mathbf{x})$   
 170 used as input for the diffusion process does not require different for each  $\mathbf{x}$ . However, the denoising  
 171 process of the diffusion model will further incorporate the feature information into the learning of the  
 172 transition matrix.

173 After obtaining the instance-dependent estimated transition matrix  $\mathbf{T}(\mathbf{x})$ , the neural network can be  
 174 learned to fit the clean label distribution by the loss function:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{T}(\mathbf{x}_i)^\top f_\phi(\mathbf{x}_i), \tilde{\mathbf{y}}_i), \quad (7)$$

175 where  $f_\phi(\cdot) : \mathcal{X} \rightarrow \Delta^{C-1}$  ( $\Delta^{C-1} \subset [0, 1]^C$  is the  $C$ -dimensional simplex) is a differentiable  
 176 function represented by a neural network with parameters  $\phi$  and  $\ell$  is a loss function usually using  
 177 cross-entropy (CE) loss.

178 The schematic diagram of the proposed framework is shown in Figure 2, and the pseudocode is  
 179 presented in Algorithm 3.

---

**Algorithm 3** A framework of DTM

---

**Input:** Training set  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ , pre-trained encoder  $f_p$ , diffusion model  $\epsilon_\theta$ , classification neural network  $f_\phi$ .

- 1: Utilize input data to train  $f_p$  or directly utilizing  $f_p$  to extract features.
- 2: Perform K-means on feature space and estimate the transition matrix for each group to get data  $\{\mathbf{x}_i, \mathbf{T}_i\}_{i=1}^N$ .
- 3: Train the diffusion model  $\epsilon_\theta$  with Algorithm 1.
- 4: Sample instance-dependent train matrix  $\mathbf{T}(\mathbf{x})$  for any input image  $\mathbf{x}_i$  with Algorithm 2.
- 5: Update the parameters of the classification network by incorporating the transition matrix  $\mathbf{T}(\mathbf{x}_i)$  into equation (7).

**Output:** Network parameters  $\phi$ .

---

180 **3.4 Matrix Transformation**

181 Considering that the transition matrix typically require the sum of each row  $\sum_{j=1}^C \mathbf{T}_{ij}(\mathbf{x})$  is 1, and  
 182 for  $\mathbf{T}_{ii}(\mathbf{x}) > \mathbf{T}_{ij}(\mathbf{x}), \forall j \neq i$ , we employ a transformation during the update learning process in our  
 183 practical experiments.

184 We utilize a  $C \times C$  weight matrix  $\mathbf{W} = (w_{ij})$  to assist in the process. Denote matrix  $\mathbf{A}$  as  
 185  $\mathbf{A}_{ii} = 1 + \sigma(w_{ii})$  for all  $i \in \{1, 2, \dots, C\}$  and  $\mathbf{A}_{ij} = \sigma(w_{ij})$  for all  $i \neq j$  where  $\sigma$  is the sigmoid  
 186 function. Then we do the normalization  $\mathbf{T}_{ij} = \frac{\mathbf{A}_{ij}}{\sum_{k=1}^C \mathbf{A}_{kj}}$  to get the transition matrix  $\mathbf{T}$ .

187 Through this transformation, we ensure that the learned transition matrix has row sums equal to 1 and  
 188 that the diagonal elements are the largest in each row. In practical experiments, we apply the diffusion  
 189 modeling discussed in subsection 3.2 to the matrix  $\mathbf{W}$ , and then transform it into the transition matrix  
 190  $\mathbf{T}$  for application. To simplify the notation, we uniformly use the term of transition matrix  $\mathbf{W}$  to  
 191 represent it, unless it leads to singularity.

192 **4 Experiments**

193 In this section, we present experimental findings to showcase the effectiveness of our proposed  
 194 method compared to other methods. We evaluate our approach on both synthetic instance-dependent  
 195 noisy datasets and real-world noisy datasets.

196 **4.1 Datasets**

197 We conduct experiments on following image classification datasets: CIFAR-10 and CIFAR-100 [13],  
 198 CIFAR-10N and CIFAR-100N [32], Clothing1M [35], Webvision and ILSVRC12 [16]. Among  
 199 them, CIFAR-10 and CIFAR-100 both have  $32 \times 32 \times 3$  color images including 50,000 training  
 200 images and 10,000 test images. CIFAR-10 has 10 classes while CIFAR-100 has 100 classes. We  
 201 generate instance-dependent noisy data on CIFAR-10 and CIFAR-100 with noise rates ranging from  
 202 10% to 50%, following the same generation method as in [33]. CIFAR-10N has three annotated  
 203 labels, namely Random1, Random 2 and Random 3. The "Aggregate" is the aggregation of three noisy  
 204 labels by majority voting, and the "Worst" is the dataset with the worst case. For CIFAR-100N, each

205 image contains a coarse label and a fine label given by a human annotator. Clothing1M is a real-world  
 206 dataset consisting of 1 million training images, consisting of 14 categories. WebVision contains 2.4  
 207 million images crawled from the websites using the 1,000 concepts in ImageNet ILSVRC12, but only  
 208 the first 50 classes of the Google image subset are used in our experiments. For the validation set  
 209 selection in our BTR method, we randomly sampled 10 samples from each observed class for each  
 210 dataset to form the validation set, while the remaining samples were used for the training set.

## 211 4.2 Experimental Setup

212 For the pre-trained model, we employ the commonly used SimCLR model [4] from contrastive  
 213 learning, which directly performs self-supervised learning on input images without utilizing additional  
 214 datasets. For the diffusion model, we follow the setup similar to DDPM [11] to set  $\beta_1 = 10^{-4}$ ,  $\beta_M =$   
 215 0.02 and utilize a similar U-Net network architecture but we reduce the  $M$  from 1000 to 10 to  
 216 accelerate the learning process. As for the classification network, it may vary depending on the  
 217 specific dataset. More specifically, for CIFAR-10/10N, we use ResNet-18 as the backbone network  
 218 with batch size 128 and learning rate 0.05. For CIFAR-100/100N, we use ResNet-34 network  
 219 with batch size 128, learning rate 0.02. For clothing1M, we use a ResNet-50 pre-trained with 10  
 220 epochs, batch size 64, learning rate 0.002 for network and divided by 10 after the 5th epoch. We use  
 221 InceptionResNetV2 network on Webvision, with 100 epochs, batch size 32, learning rate 0.02 for  
 222 network and divided by 10 after the 30th and 60th epoch. For clustering, we utilize the K-means  
 223 method, where the number of clusters is set to 10 times the number of classes in the datasets. For  
 224 the initialization of transition matrix, the update method and setting are consistent with [17]. While  
 225 the updates for other parameters are performed using the stochastic gradient descent optimization  
 226 method.

Table 1: Test accuracy with instance-dependent noise on CIFAR-10/100.

	CIFAR-10				
	IDN-10%	IDN-20%	IDN-30%	IDN-40%	IDN-50%
CE	88.86±0.23	86.93±0.17	82.42±0.44	76.68±0.23	58.93±1.54
VolMinNet	89.97±0.57	87.01±0.64	83.80±0.67	79.52±0.83	61.90±1.06
PeerLoss	90.89±0.07	89.21±0.63	85.70±0.56	78.51±1.23	59.08±1.05
BLTM	90.45±0.72	88.14±0.66	84.55±0.48	79.71±0.95	63.33±2.75
PartT	90.32±0.15	89.33±0.70	85.33±1.86	80.59±0.41	64.58±2.86
MEIDTM	92.91±0.07	92.26±0.25	90.73±0.34	85.94±0.92	73.77±0.82
SOP	93.58±0.31	93.07±0.45	92.42±0.43	89.83±0.77	82.52±0.97
CC	95.24±0.20	93.68±0.12	93.31±0.46	<b>94.97±0.09</b>	91.19±0.34
LRA	95.87±0.42	94.70±0.28	93.79±0.40	92.72±0.29	90.95±0.43
DTM	<b>96.45±0.17</b>	<b>95.90±0.21</b>	<b>95.14±0.20</b>	94.82±0.31	<b>92.04±0.42</b>
	CIFAR-100				
	IDN-10%	IDN-20%	IDN-30%	IDN-40%	IDN-50%
CE	66.55±0.23	63.94±0.51	61.97±1.16	58.70±0.56	56.63±0.69
VolMinNet	67.78±0.62	66.13±0.47	61.08±0.90	57.35±0.83	52.60±1.31
PeerLoss	65.64±1.07	63.83±0.48	61.64±0.67	58.30±0.80	55.41±0.28
BLTM	68.42±0.42	66.62±0.85	64.72±0.64	59.38±0.65	55.68±1.43
PartT	67.33±0.33	65.33±0.59	64.56±1.55	59.73±0.76	56.80±1.32
MEIDTM	69.88±0.45	69.16±0.16	66.76±0.30	63.46±0.48	59.18±0.16
SOP	74.09±0.52	73.13±0.46	72.14±0.46	68.98±0.58	64.24±0.86
CC	80.52±0.22	79.61±0.19	77.34±0.31	76.58±0.25	72.68±0.36
LRA	81.20±0.16	80.53±0.29	78.22±0.19	76.55±0.31	72.97±0.51
DTM	<b>82.96±0.25</b>	<b>82.04±0.32</b>	<b>80.87±0.45</b>	<b>78.56±0.60</b>	<b>74.85±0.56</b>

## 227 4.3 Comparison Methods

228 In our experiments, we included the following common transition matrix and baseline methods as  
 229 comparison: (1) VolMinNet [17], (2) PeerLoss [21] (3) BLTM [37], (4) PartT [33], (5) MEIDTM  
 230 [6], as well as state-of-the-art methods for learning with noisy labels: (6) Co-teaching [10], (7) ELR+  
 231 [18], (8) DivideMix [14], (9) SOP and SOP+ [19], (10) PGDF [5], (11) CC [40], (12) LRA [3]  
 232 with SimCLR as encoder similarly.

Table 2: Test accuracy on CIFAR-10N and CIFAR-100N.

	CIFAR-10N					CIFAR-100N
	Aggregate	Random 1	Random 2	Random 3	Worst	Noisy
Co-teaching	91.20±0.13	90.33±0.13	90.30±0.17	90.15±0.18	83.83±0.13	60.37±0.27
ELR+	94.83±0.10	94.43±0.41	94.20±0.24	94.34±0.22	91.09±1.60	66.72±0.07
DivideMix	95.01±0.71	95.16±0.19	94.89±0.23	95.03±0.20	92.56±0.42	71.13±0.48
SOP+	95.61±0.13	95.28±0.13	95.31±0.10	95.39±0.11	93.24±0.21	67.81±0.23
PGDF	95.35±0.12	94.95±0.21	94.78±0.34	94.92±0.28	94.22±0.29	67.76±0.35
CC	95.63±0.21	95.11±0.31	94.93±0.37	95.09±0.21	94.24±0.40	71.21±0.22
LRA	94.57±0.23	94.19±0.17	94.38±0.42	94.02±0.32	93.20±0.59	70.96±0.53
DTM	<b>96.13±0.17</b>	<b>95.98±0.22</b>	<b>96.01±0.28</b>	<b>95.78±0.34</b>	<b>94.93±0.21</b>	<b>72.51±0.30</b>

#### 233 4.4 Experimental Results on Synthetic Datasets

234 We primarily validated our proposed method DTM against previous instance-based transition matrix  
 235 methods on synthetic CIFAR-10/100 noise datasets. These methods mainly focus on estimating the  
 236 transition matrix and some methods applicable to instance-dependent label noise. We performed 5  
 237 independent runs for each experimental configuration, and the average values and standard deviations  
 238 of each experiment are presented in Table 1.

239 The results demonstrate that our proposed DTR method outperforms other methods of the same  
 240 category across various noise rates. It is evident that traditional transition matrix methods for class-  
 241 dependent noise as VolMinNet exhibit subpar performance when handling instance-dependent noise.  
 242 While even advanced transition matrix methods for instance-dependent label noise such as BLTM,  
 243 ParT and MEIDTM, still show significant gaps compared to our method.

244 Furthermore, as the noise rates increase, the test accuracy of existing transition matrix methods  
 245 significantly decline. This is particularly pronounced in the case of CIFAR-100 with 50% instance-  
 246 dependent noise (IDN) data, where all transition matrix methods achieve test accuracy below 60%.  
 247 In contrast, our proposed DTR method achieves a remarkable test accuracy of 74.85%, showcasing  
 248 its exceptional performance. That demonstrates relatively robust performance of DTM with only a  
 249 slight decrease as the noise rate increases.

250 This experiment clearly demonstrates that there is a significant performance gap between previous  
 251 transition matrix methods and other advanced techniques, such as CC and LRA, when dealing with  
 252 instance-dependent noise problems. However, the experimental results indicate that our proposed  
 253 method DTM, which incorporates the diffusion model into the estimation of the transition matrix,  
 254 outperforms these advanced techniques, except for the case of 40% noise in CIFAR-100, where  
 255 our method slightly underperforms CC. It is evident that by leveraging the diffusion modeling to  
 256 estimate the transition matrix, we effectively incorporate the image’s feature information, leading to a  
 257 substantial improvement in the effectiveness of the transition matrix.

#### 258 4.5 Experimental Results on Real-World Datasets

259 In addition to synthetic datasets, we also applied our method to real-world datasets and compared it  
 260 with other state-of-the-art techniques for handling label noise problems. The results are presented in  
 261 Table 2 and Table 3.

Table 3: Test accuracy on Clothing1M, Webvision and ILSVRC12.

	Clothing1M	Webvision	ILSVRC12
Co-teaching	69.2	63.6	61.5
ELR+	74.81	77.78	70.29
DivideMix	74.76	77.32	75.20
SOP+	74.98	77.60	75.29
PGDF	75.19	81.47	75.45
CC	75.40	79.36	76.08
LRA	75.32	80.05	76.64
DTM	<b>75.57</b>	<b>81.95</b>	<b>77.55</b>

262 The results demonstrate that regardless of the type of noise labels, whether it is aggregated, random,  
 263 or the worst-case scenario in CIFAR-10N, as well as in CIFAR-100N with more label categories,  
 264 our method consistently achieves the best results in handling real-world noise. When dealing with  
 265 large datasets like Clothing1M and complex image datasets like Webvision, DTM also performs  
 266 comparably to other state-of-the-art methods.

267 Through extensive experiments on five real-world datasets and the results on synthetic datasets above,  
 268 our method outperforms the LRA method, which also utilizes the diffusion model for label noise  
 269 problems. The LRA method models label diffusion with fewer dimensional information and lacks the  
 270 rationale of our method, which considers noise generation from a transfer probability distribution  
 271 perspective. The experiments demonstrate that our method achieves better learning performance by  
 272 effectively integrating the transition matrix with the diffusion model.

#### 273 4.6 Ablation Study

274 Besides the aforementioned experiments, we conducted ablation studies on proposed DTM method  
 275 to assess the importance of each component. Table 4 presents the comparative results under 20%  
 276 and 40% instance-dependent noise rates, where "w/o" denotes "without". We conducted ablation  
 277 experiments on three components of our method, they are diffusion module, pre-trained encoder  
 278 module, and clustering module respectively. "w/o diffusion" indicates directly using the features  
 279 extracted by the pre-trained model for the classification task with the transition matrix. "w/o pre-train"  
 280 means not extracting features through self-supervised learning and directly utilizing the classification  
 281 network with the diffusion model. "w/o clustering" indicates that the initial transition matrix used for  
 282 the diffusion model is the same for all samples.

Table 4: Ablation study of DTR. The data in the table represents the test accuracy.

	CIFAR-10		CIFAR-100	
	IDN-0.2	IDN-0.4	IDN-0.2	IDN-0.4
w/o pre-train	90.52	83.61	66.17	61.79
w/o clustering	92.25	88.35	71.93	66.47
w/o diffusion	93.74	91.66	79.82	73.51
DTR	<b>95.90</b>	<b>94.82</b>	<b>82.04</b>	<b>78.56</b>

283 From the results in Table 4, it can be observed that regardless of which component of diffusion  
 284 module, pre-trained encoder module and clustering module is missing, the performance is consistently  
 285 weaker compared to the original DTM. This indicates that each module plays a crucial role in our  
 286 method. Our approach effectively combines the transition matrix, diffusion model, and pre-trained  
 287 feature extraction, leading to significant improvements.

## 288 5 Conclusion

289 In this paper, we propose a method that models the transition matrix using diffusion models, incorpor-  
 290 ating the feature information extracted by a pre-trained encoder into the estimation of the transition  
 291 matrix. This approach enables the model to handle instance-dependent label noise with a wider range  
 292 of applicability. Experimental results on both synthetic and real-world noisy datasets demonstrate the  
 293 effectiveness of our proposed method.

## 294 References

295 [1] Görkem Algan and Ilkay Ulusoy. Image classification with deep learning in the presence of  
 296 noisy labels: A survey. *Knowledge-Based Systems*, 215:106771, 2021.

297 [2] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio,  
 298 Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al.  
 299 A closer look at memorization in deep networks. In *International Conference on Machine*  
 300 *Learning*, pages 233–242. PMLR, 2017.

- 301 [3] Jian Chen, Ruiyi Zhang, Tong Yu, Rohan Sharma, Zhiqiang Xu, Tong Sun, and Changyou Chen.  
 302 Label-retrieval-augmented diffusion models for learning from noisy labels. *arXiv preprint*  
 303 *arXiv:2305.19518*, 2023.
- 304 [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
 305 for contrastive learning of visual representations. In *International Conference on Machine*  
 306 *Learning*, pages 1597–1607. PMLR, 2020.
- 307 [5] Wenkai Chen, Chuang Zhu, and Mengting Li. Sample prior guided robust model learning to  
 308 suppress noisy labels. In *Joint European Conference on Machine Learning and Knowledge*  
 309 *Discovery in Databases*, pages 3–19. Springer, 2023.
- 310 [6] De Cheng, Tongliang Liu, Yixiong Ning, Nannan Wang, Bo Han, Gang Niu, Xinbo Gao,  
 311 and Masashi Sugiyama. Instance-dependent label-noise learning with manifold-regularized  
 312 transition matrix estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
 313 *and Pattern Recognition*, pages 16630–16639, 2022.
- 314 [7] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion  
 315 models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
 316 2023.
- 317 [8] Amit Daniely and Elad Granot. Generalization bounds for neural networks via approximate  
 318 description length. *Advances in Neural Information Processing Systems*, 32, 2019.
- 319 [9] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adapta-  
 320 tion layer. In *International Conference on Learning Representations*, 2016.
- 321 [10] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi  
 322 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.  
 323 *Advances in Neural Information Processing Systems*, 31, 2018.
- 324 [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
 325 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 326 [12] Zhimeng Jiang, Kaixiong Zhou, Zirui Liu, Li Li, Rui Chen, Soo-Hyun Choi, and Xia Hu. An  
 327 information fusion approach to learning with instance-dependent label noise. In *International*  
 328 *Conference on Learning Representations*, 2021.
- 329 [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
 330 2009.
- 331 [14] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as  
 332 semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020.
- 333 [15] Shikun Li, Xiaobo Xia, Hansong Zhang, Yibing Zhan, Shiming Ge, and Tongliang Liu. Esti-  
 334 mating noise transition matrix with label correlations for noisy multi-label learning. *Advances*  
 335 *in Neural Information Processing Systems*, 35:24184–24198, 2022.
- 336 [16] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. Webvision database:  
 337 Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017.
- 338 [17] Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi Sugiyama. Provably end-to-end  
 339 label-noise learning without anchor points. In *International Conference on Machine Learning*,  
 340 pages 6403–6413. PMLR, 2021.
- 341 [18] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-  
 342 learning regularization prevents memorization of noisy labels. *Advances in Neural Information*  
 343 *Processing Systems*, 33:20331–20342, 2020.
- 344 [19] Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. Robust training under label noise by over-  
 345 parameterization. In *International Conference on Machine Learning*, pages 14153–14172.  
 346 PMLR, 2022.
- 347 [20] Yang Liu, Hao Cheng, and Kun Zhang. Identifiability of label noise transition matrix. In  
 348 *International Conference on Machine Learning*, pages 21475–21496. PMLR, 2023.

- 349 [21] Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing  
350 noise rates. In *International Conference on Machine Learning*, pages 6226–6236. PMLR, 2020.
- 351 [22] Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey.  
352 Normalized loss functions for deep learning with noisy labels. In *International Conference on*  
353 *Machine Learning*, pages 6543–6553. PMLR, 2020.
- 354 [23] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning  
355 with noisy labels. *Advances in Neural Information Processing Systems*, 26, 2013.
- 356 [24] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu.  
357 Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings*  
358 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1944–1952, 2017.
- 359 [25] Samira Pouyanfar, Saad Sadiq, Yilin Yan, Haiman Tian, Yudong Tao, Maria Presa Reyes, Mei-  
360 Ling Shyu, Shu-Ching Chen, and Sundaraja S Iyengar. A survey on deep learning: Algorithms,  
361 techniques, and applications. *ACM Computing Surveys (CSUR)*, 51(5):1–36, 2018.
- 362 [26] Jun Shu, Qian Zhao, Zongben Xu, and Deyu Meng. Meta transition adaptation for robust deep  
363 learning with noisy labels. *arXiv preprint arXiv:2006.05697*, 2020.
- 364 [27] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
365 *preprint arXiv:2010.02502*, 2020.
- 366 [28] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data  
367 distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- 368 [29] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and  
369 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*  
370 *preprint arXiv:2011.13456*, 2020.
- 371 [30] Jialu Wang, Yang Liu, and Caleb Levy. Fair classification with group-dependent label noise. In  
372 *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages  
373 526–536, 2021.
- 374 [31] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross  
375 entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International*  
376 *Conference on Computer Vision*, pages 322–330, 2019.
- 377 [32] Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning  
378 with noisy labels revisited: A study using real-world human annotations. *arXiv preprint*  
379 *arXiv:2110.12088*, 2021.
- 380 [33] Xiaobo Xia, Tongliang Liu, Bo Han, Nannan Wang, Mingming Gong, Haifeng Liu, Gang Niu,  
381 Dacheng Tao, and Masashi Sugiyama. Part-dependent label noise: Towards instance-dependent  
382 label noise. *Advances in Neural Information Processing Systems*, 33:7597–7610, 2020.
- 383 [34] Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi  
384 Sugiyama. Are anchor points really indispensable in label-noise learning? *Advances in Neural*  
385 *Information Processing Systems*, 32, 2019.
- 386 [35] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive  
387 noisy labeled data for image classification. In *Proceedings of the IEEE Conference on Computer*  
388 *Vision and Pattern Recognition*, pages 2691–2699, 2015.
- 389 [36] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang,  
390 Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and  
391 applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- 392 [37] Shuo Yang, Erkun Yang, Bo Han, Yang Liu, Min Xu, Gang Niu, and Tongliang Liu. Estimating  
393 instance-dependent bayes-label transition matrix using a deep neural network. In *International*  
394 *Conference on Machine Learning*, pages 25302–25312. PMLR, 2022.

- 395 [38] LIN Yong, Renjie Pi, Weizhong Zhang, Xiaobo Xia, Jiahui Gao, Xiao Zhou, Tongliang Liu,  
396 and Bo Han. A holistic view of label noise transition matrix in deep learning and beyond. In  
397 *The Eleventh International Conference on Learning Representations, 2022*.
- 398 [39] Yivan Zhang, Gang Niu, and Masashi Sugiyama. Learning noise transition matrix from only  
399 noisy labels via total variation regularization. In *International Conference on Machine Learning*,  
400 pages 12501–12512. PMLR, 2021.
- 401 [40] Ganlong Zhao, Guanbin Li, Yipeng Qin, Feng Liu, and Yizhou Yu. Centrality and consistency:  
402 two-stage clean samples identification for learning with instance-dependent noisy labels. In  
403 *European Conference on Computer Vision*, pages 21–37. Springer, 2022.
- 404 [41] Zhaowei Zhu, Yiwen Song, and Yang Liu. Clusterability as an alternative to anchor points  
405 when learning with noisy labels. In *International Conference on Machine Learning*, pages  
406 12912–12923. PMLR, 2021.
- 407 [42] Zhaowei Zhu, Jialu Wang, and Yang Liu. Beyond images: Label noise transition matrix  
408 estimation for tasks with lower-quality features. In *International Conference on Machine*  
409 *Learning*, pages 27633–27653. PMLR, 2022.

## 410 **NeurIPS Paper Checklist**

### 411 **1. Claims**

412 Question: Do the main claims made in the abstract and introduction accurately reflect the  
413 paper's contributions and scope?

414 Answer: [\[Yes\]](#)

415 Justification: The main content and contributions of the work are reflected in the abstract  
416 and introduction.

417 Guidelines:

- 418 • The answer NA means that the abstract and introduction do not include the claims  
419 made in the paper.
- 420 • The abstract and/or introduction should clearly state the claims made, including the  
421 contributions made in the paper and important assumptions and limitations. A No or  
422 NA answer to this question will not be perceived well by the reviewers.
- 423 • The claims made should match theoretical and experimental results, and reflect how  
424 much the results can be expected to generalize to other settings.
- 425 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
426 are not attained by the paper.

### 427 **2. Limitations**

428 Question: Does the paper discuss the limitations of the work performed by the authors?

429 Answer: [\[Yes\]](#)

430 Justification: In the experimental section, we analyze the applicability and limitations of our  
431 method.

432 Guidelines:

- 433 • The answer NA means that the paper has no limitation while the answer No means that  
434 the paper has limitations, but those are not discussed in the paper.
- 435 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 436 • The paper should point out any strong assumptions and how robust the results are to  
437 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
438 model well-specification, asymptotic approximations only holding locally). The authors  
439 should reflect on how these assumptions might be violated in practice and what the  
440 implications would be.
- 441 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
442 only tested on a few datasets or with a few runs. In general, empirical results often  
443 depend on implicit assumptions, which should be articulated.
- 444 • The authors should reflect on the factors that influence the performance of the approach.  
445 For example, a facial recognition algorithm may perform poorly when image resolution  
446 is low or images are taken in low lighting. Or a speech-to-text system might not be  
447 used reliably to provide closed captions for online lectures because it fails to handle  
448 technical jargon.
- 449 • The authors should discuss the computational efficiency of the proposed algorithms  
450 and how they scale with dataset size.
- 451 • If applicable, the authors should discuss possible limitations of their approach to  
452 address problems of privacy and fairness.
- 453 • While the authors might fear that complete honesty about limitations might be used by  
454 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
455 limitations that aren't acknowledged in the paper. The authors should use their best  
456 judgment and recognize that individual actions in favor of transparency play an impor-  
457 tant role in developing norms that preserve the integrity of the community. Reviewers  
458 will be specifically instructed to not penalize honesty concerning limitations.

### 459 **3. Theory Assumptions and Proofs**

460 Question: For each theoretical result, does the paper provide the full set of assumptions and  
461 a complete (and correct) proof?

462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515

Answer: [NA]

Justification: The focus of the work is on application and does not include a theoretical proof component.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper provides a detailed description of the model construction and the specifics of the experimental data.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

**5. Open access to data and code**

516 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
517 tions to faithfully reproduce the main experimental results, as described in supplemental  
518 material?

519 Answer: [No]

520 Justification: Upon acceptance of our paper, we will provide open-source code. The data we  
521 used is from commonly available open-source datasets.

522 Guidelines:

- 523 • The answer NA means that paper does not include experiments requiring code.
- 524 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
525 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 526 • While we encourage the release of code and data, we understand that this might not be  
527 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
528 including code, unless this is central to the contribution (e.g., for a new open-source  
529 benchmark).
- 530 • The instructions should contain the exact command and environment needed to run to  
531 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
532 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 533 • The authors should provide instructions on data access and preparation, including how  
534 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 535 • The authors should provide scripts to reproduce all experimental results for the new  
536 proposed method and baselines. If only a subset of experiments are reproducible, they  
537 should state which ones are omitted from the script and why.
- 538 • At submission time, to preserve anonymity, the authors should release anonymized  
539 versions (if applicable).
- 540 • Providing as much information as possible in supplemental material (appended to the  
541 paper) is recommended, but including URLs to data and code is permitted.

## 542 6. Experimental Setting/Details

543 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
544 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
545 results?

546 Answer: [Yes]

547 Justification: The experimental section of the paper provides details of the model and data.

548 Guidelines:

- 549 • The answer NA means that the paper does not include experiments.
- 550 • The experimental setting should be presented in the core of the paper to a level of detail  
551 that is necessary to appreciate the results and make sense of them.
- 552 • The full details can be provided either with the code, in appendix, or as supplemental  
553 material.

## 554 7. Experiment Statistical Significance

555 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
556 information about the statistical significance of the experiments?

557 Answer: [Yes]

558 Justification: We conducted multiple repeated experiments to validate our approach and  
559 performed ablation experiments.

560 Guidelines:

- 561 • The answer NA means that the paper does not include experiments.
- 562 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
563 dence intervals, or statistical significance tests, at least for the experiments that support  
564 the main claims of the paper.
- 565 • The factors of variability that the error bars are capturing should be clearly stated (for  
566 example, train/test split, initialization, random drawing of some parameter, or overall  
567 run with given experimental conditions).

- 568 • The method for calculating the error bars should be explained (closed form formula,  
569 call to a library function, bootstrap, etc.)
- 570 • The assumptions made should be given (e.g., Normally distributed errors).
- 571 • It should be clear whether the error bar is the standard deviation or the standard error  
572 of the mean.
- 573 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
574 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
575 of Normality of errors is not verified.
- 576 • For asymmetric distributions, the authors should be careful not to show in tables or  
577 figures symmetric error bars that would yield results that are out of range (e.g. negative  
578 error rates).
- 579 • If error bars are reported in tables or plots, The authors should explain in the text how  
580 they were calculated and reference the corresponding figures or tables in the text.

## 581 8. Experiments Compute Resources

582 Question: For each experiment, does the paper provide sufficient information on the com-  
583 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
584 the experiments?

585 Answer: [Yes]

586 Justification: We list the relevant details in the experimental section.

587 Guidelines:

- 588 • The answer NA means that the paper does not include experiments.
- 589 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
590 or cloud provider, including relevant memory and storage.
- 591 • The paper should provide the amount of compute required for each of the individual  
592 experimental runs as well as estimate the total compute.
- 593 • The paper should disclose whether the full research project required more compute  
594 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
595 didn't make it into the paper).

## 596 9. Code Of Ethics

597 Question: Does the research conducted in the paper conform, in every respect, with the  
598 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

599 Answer: [Yes]

600 Justification: We submitted the paper following the NeurIPS Code of Ethics.

601 Guidelines:

- 602 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 603 • If the authors answer No, they should explain the special circumstances that require a  
604 deviation from the Code of Ethics.
- 605 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
606 eration due to laws or regulations in their jurisdiction).

## 607 10. Broader Impacts

608 Question: Does the paper discuss both potential positive societal impacts and negative  
609 societal impacts of the work performed?

610 Answer: [Yes]

611 Justification: We discuss the positive implications of our work and ensure it does not have  
612 any negative societal impact.

613 Guidelines:

- 614 • The answer NA means that there is no societal impact of the work performed.
- 615 • If the authors answer NA or No, they should explain why their work has no societal  
616 impact or why the paper does not address societal impact.

- 617
- 618
- 619
- 620
- 621
- 622
- 623
- 624
- 625
- 626
- 627
- 628
- 629
- 630
- 631
- 632
- 633
- 634
- 635
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 636 11. Safeguards

637 Question: Does the paper describe safeguards that have been put in place for responsible  
638 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
639 image generators, or scraped datasets)?

640 Answer: [NA]

641 Justification: There are no concerns in this regard regarding this work.

642 Guidelines:

- 643
- 644
- 645
- 646
- 647
- 648
- 649
- 650
- 651
- 652
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 653 12. Licenses for existing assets

654 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
655 the paper, properly credited and are the license and terms of use explicitly mentioned and  
656 properly respected?

657 Answer: [Yes]

658 Justification: The data and code used in our work are all publicly available and open-source.

659 Guidelines:

- 660
- 661
- 662
- 663
- 664
- 665
- 666
- 667
- 668
- 669
- 670
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 671
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 672  
673  
674

675 **13. New Assets**

676 Question: Are new assets introduced in the paper well documented and is the documentation  
677 provided alongside the assets?

678 Answer: [NA]

679 Justification: The paper currently does not include any new assets.

680 Guidelines:

- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 681  
682  
683  
684  
685  
686  
687  
688

689 **14. Crowdsourcing and Research with Human Subjects**

690 Question: For crowdsourcing experiments and research with human subjects, does the paper  
691 include the full text of instructions given to participants and screenshots, if applicable, as  
692 well as details about compensation (if any)?

693 Answer: [NA]

694 Justification: The paper does not involve crowdsourcing nor research with human subjects.

695 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 696  
697  
698  
699  
700  
701  
702  
703

704 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
705 Subjects**

706 Question: Does the paper describe potential risks incurred by study participants, whether  
707 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
708 approvals (or an equivalent approval/review based on the requirements of your country or  
709 institution) were obtained?

710 Answer: [NA]

711 Justification: The paper does not involve crowdsourcing nor research with human subjects.

712 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
  - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 713  
714  
715  
716  
717  
718  
719  
720  
721  
722