

INFORMATIVE DATA REWEIGHTING FOR IMAGE CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep Neural Networks (DNNs) have achieved remarkable success in image classification tasks. However, their training typically requires large-scale, high-quality labeled datasets, which may be scarce or infeasible to obtain in certain computer vision tasks. To alleviate this challenge, Generative Data Augmentation (GDA) has been introduced to improve model performance by increasing the number of training samples with synthetic data generated by models such as Diffusion Models (DMs). Despite its benefits, GDA-generated synthetic samples often contain noise, which can negatively impact the performance of image classification models when incorporated into training. Prior approaches, including data selection and reweighting techniques, aim to address this issue but often rely on external expert models or clean metadata. In this work, we introduce Informative Data Reweighting (IDR), a principled sample reweighting framework based on the Information Bottleneck (IB) principle, to enhance the performance of DNNs for image classification using GDA. Through extensive experiments, we demonstrate that IDR effectively prioritizes more informative training samples in the augmented training set comprising original real training samples and synthetic training samples, resulting in substantial improvements over existing data selection and reweighting strategies for GDA in image classification. The code for IDR is available at <https://anonymous.4open.science/r/IDR-3BE0/>.

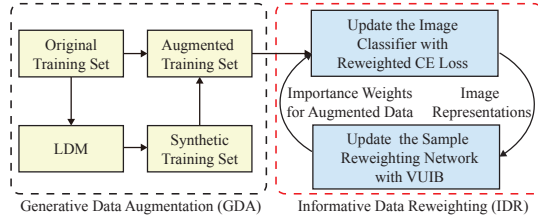
1 INTRODUCTION

Deep Neural Networks (DNNs) have achieved strong performance across computer vision tasks such as image classification, semantic segmentation, and object detection, but their success critically depends on large-scale, high-quality annotated datasets (Feng et al., 2020). However, obtaining such datasets poses substantial challenges for many computer vision tasks (El Jiani et al., 2022; Xiao et al., 2023) due to constraints such as resource limitations and privacy concerns (Esteva et al., 2021; Price & Cohen, 2019; Ali et al., 2023; Ramudu et al., 2023). To mitigate this issue, researchers have explored Generative Data Augmentation (GDA) (Saryıldız et al., 2023; Lei et al., 2023; Azizi et al., 2023; Trabucco et al., 2024), leveraging techniques such as Generative Adversarial Networks (GANs) (Zhang et al., 2021; Li et al., 2022) and Diffusion Models (DMs) (He et al., 2023; Tian et al., 2023; Yuan et al., 2022; Bansal & Grover, 2023; Vendrow et al., 2023) to generate synthetic training samples that can be combined with the original training set for data augmentation, leading to an augmented training set.

Limitations in Existing GDA Methods. Despite its potential, synthetic samples generated by GDA often contain noise, such as local visual artifacts and semantic inconsistencies, which can degrade downstream image classification performance (Corvi et al., 2023; Azizi et al., 2023; He et al., 2023). To mitigate this issue, prior work has proposed data selection and sample reweighting strategies that train classifiers on curated or weighted synthetic samples (Chhabra et al., 2024; He et al., 2023), typically relying on meta-networks with clean metadata to assign higher importance to informative samples, while assuming the availability of such clean metadata (Shu et al., 2019; Guo et al., 2022; Jain et al., 2024). A closely related study, CBF (He et al., 2023), filters low-quality synthetic samples using CLIP zero-shot confidence (Radford et al., 2021) in GDA, but its effectiveness relies heavily on CLIP, an expert vision-language model pretrained on large-scale, non-public image-text datasets. GenDataAgent (Li et al., 2025) filters out-of-distribution samples from an expert generative model, Stable Diffusion, using VoG-based logit gradients, while it cannot be employed to select informative in-domain synthetic samples as in GDA (He et al., 2023; Azizi et al., 2023), where the generative

054 model is trained on the same domain as the original training data. In summary, the current machine
 055 learning literature lacks a principled sample reweighting method for GDA, which does not rely on
 056 clean metadata or external expert models. We further note that the upper bound for mutual information
 057 between the learned features and the input training features in the IB loss in IDS (Wang et al., 2025)
 058 and DCS-Transformer (Wang et al., 2024) is vacuous in terms of the learned features, with details in
 059 Section E.2 of the appendix. Such drawback of IDS results in its inferior performance compared to
 060 our IDR as evidenced in Tables 1, 5, 6, and 7 of Sections 4.1, D.4, D.5, and D.6.

061 To address this issue, we propose a novel In-
 062 formative Data Reweighting (IDR) method for
 063 reweighting samples in the augmented training
 064 set obtained from GDA, inspired by the In-
 065 formation Bottleneck (IB) principle. As illus-
 066 trated in Figure 1, IDR starts with GDA to gen-
 067 erate a synthetic training set using the LDM
 068 trained on the original training set, which is
 069 then combined with the original training set
 070 to form an augmented training set. IDR aims
 071 to assign higher importance weights to more
 072 informative training samples in the augmented



073 Figure 1: The pipeline for Generative Data Aug-
 074 mentation (GDA) and Informative Data Reweighting
 075 (IDR).

076 training set that contribute more to the image classification task. To achieve this goal, IDR trains a
 077 sample reweighting network that assigns importance weights to the samples in the augmented training
 078 set by reducing the IB loss on them. Let \tilde{X} be the input features, and \tilde{Z} be the learned features by
 079 the classification network. Let \tilde{Y} be the ground truth training class labels. The principle of IB is to
 080 increase the mutual information between \tilde{Z} and \tilde{Y} while reducing the mutual information between
 081 \tilde{Z} and \tilde{X} . That is, the IB principle encourages reduction of the IB loss, $I(\tilde{Z}, \tilde{X}) - I(\tilde{Z}, \tilde{Y})$, where
 082 $I(\cdot, \cdot)$ denotes mutual information modeling the correlation of its input variables. The reduction of
 083 the IB loss is achieved by optimizing a novel variational upper bound for the IB loss, termed VUIB,
 084 which can be optimized by standard SGD algorithms. As illustrated in Figure 1, the cross-entropy
 085 loss reweighted by the importance weights and the VUIB are iteratively optimized to update the
 086 weights in the classification network and the sample reweighting network, inspired by the bi-level
 087 optimization method in (Shu et al., 2019).

088 **Our Contributions.** The contributions of this paper are presented as follows.

089 First, to the best of our knowledge, IDR is among the first to perform sample reweighting in GDA
 090 by employing a principled IB framework in the sample reweighting process, where each sample in
 091 the augmented training set receives an importance weight, aiming to improve the accuracy of the
 092 classifier trained on the reweighted training samples. A sample reweighting network is optimized
 093 for reducing the IB loss on the augmented training set, such that the IB principle, learning features
 094 more strongly correlated with class labels while decreasing their correlation with the inputs, is better
 095 adhered to. In contrast to existing works in sample reweighting (Shu et al., 2019; Guo et al., 2022;
 096 Jain et al., 2024) and data selection (Chhabra et al., 2024; He et al., 2023), IDR does not require
 097 either clean metadata or external expert models.

098 Second, to facilitate the reduction of the IB loss under mini-batch-based SGD optimization, we
 099 derive a novel and theoretically justified variational upper bound for the IB loss, termed VUIB, which
 100 is amenable to optimization using standard SGD algorithms. Accordingly, the IB loss is reduced by
 101 optimizing its corresponding upper bound, namely VUIB. In contrast to existing upper bounds for the
 102 IB loss, such as VIB (Dai et al., 2018; Srivastava et al., 2021), which relies on an impractical Gaussian
 103 assumption on the hidden features of deep neural networks, and APIB (Guo et al., 2023), which only
 104 reduces an approximation of the IB loss, IDR directly reduces the variational upper bound VUIB
 105 without introducing any distributional assumptions on the hidden representations or approximation
 106 to the IB loss. Furthermore, the proposed VUIB is computationally efficient, with a complexity of
 107 $\Theta(N'T_0 + N'C)$, where C denotes the number of classes, N' represents the number of training
 samples, and T_0 is the computational cost of a forward and backward pass of the neural network for
 each sample. By comparison, the mutual-information upper bound employed in CLUB (Cheng et al.,
 2020), while avoiding distributional assumptions on hidden features, incurs a significantly higher
 computational complexity of $\Theta(N'^2T_0)$, since $N' \gg C$. As shown in Table 2, models trained using
 VUIB achieve substantially better performance than those based on CLUB, VIB, and APIB. During

the training of IDR, the importance-weighted cross-entropy loss and the VUIB are optimized in an iterative manner to update the parameters of the classification network and the sample reweighting network. As demonstrated by the results in Section 4.1, IDR significantly outperforms state-of-the-art data reweighting and data selection methods for image classification on ImageNet-1K.

Recent works propose improved conventional data augmentation strategies for image classification based on transformations of the original training samples (Kim & Kim, 2025; Kang et al., 2025). In comparison, IDR leverages GDA to expand the training set and employs informative reweighting to mitigate the noise introduced by synthetic samples, making it orthogonal and complementary to conventional data augmentation approaches. Experiments in Section D.2 demonstrate that IDR outperforms existing conventional augmentation methods and achieves further performance improvements when combined with them.

It is also worthwhile to mention that while the reweighting is originally developed to render importance scores of synthetic data, it is revealed in this work that reweighting can also be applied to the original real training samples so that more informative real samples contributing more to the classification task receive higher importance weights. Section D.3 presents an ablation study showing that jointly reweighting real and synthetic samples in IDR is more effective than reweighting synthetic samples alone. Throughout this paper we use $[N]$ to denote all the natural numbers between 1 and N .

2 RELATED WORKS

Information Bottleneck Principle. The Information Bottleneck (IB) principle (Tishby et al., 2000) states that an optimal DNN compresses input data, retaining only information necessary for predicting targets, thus increasing mutual information with outputs while reducing it with inputs. Deep VIB (Alemi et al., 2017) incorporates IB as a training objective. Empirical (Lai et al., 2021; Zhou et al., 2022) and theoretical (Amjad & Geiger, 2020; Kawaguchi et al., 2023) studies show that closer adherence to IB improves performance.

Generative Data Augmentation, Data Selection and Sample Reweighting. Generative data augmentation (GDA) creates synthetic training data but remains challenging. Although it can boost image classification, synthetic samples often add noise (Azizi et al., 2023; Trabucco et al., 2024; Na et al., 2024), motivating quality control. Prior work follows three paths: improving synthesis quality (Sarıyıldız et al., 2023; Lei et al., 2023; Zhou et al., 2023), selecting high-quality subsets (Song et al., 2023; Lin et al., 2023; He et al., 2023; Chhabra et al., 2024; Li et al., 2025), and reweighting samples (Mo et al., 2019; Shu et al., 2019; Guo et al., 2022; Jain et al., 2024). For example, CBF (He et al., 2023) uses CLIP zero-shot confidence (Radford et al., 2021), while GenDataAgent (Li et al., 2025) filters synthetic data via VoG-based logit gradients (Agarwal et al., 2022).

3 INFORMATIVE DATA REWEIGHTING

Given the original training set $\mathcal{D}_{\text{real}} = \{x_i, y_i\}_{i=1}^N$ for image classification where x_i is the input feature with class label $y_i \in [C]$ and C is the number of classes in the training set, we aim to generate synthetic training set $\mathcal{D}_{\text{syn}} = \{\hat{x}_i, \hat{y}_i\}_{i=1}^M$ with diffusion models and train a classifier on the augmented training set $\mathcal{D}_{\text{aug}} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^{N'} = \{x_i, y_i\}_{i=1}^N \cup \{\hat{x}_i, \hat{y}_i\}_{i=1}^M$ with $N' = N + M$. To mitigate the negative impact of noise in the synthetic training samples, we propose Informative Data Reweighting (IDR) to reweight the training samples in \mathcal{D}_{aug} with a sample reweighting network. The sample reweighting network is trained by optimizing the variational upper bound for the Information Bottleneck (IB) loss on the augmented training set in the hope that more informative training samples receive higher weights, thus improving the performance of the classifier trained on the augmented training set. In Section 3.1, we describe the details for generating the synthetic training samples with diffusion models. We then introduce our novel variational upper bound for the IB loss in Section 3.2. In Section 3.3, we describe the training of the reweighting network and the classifier network in IDR.

3.1 GENERATING SYNTHETIC TRAINING SAMPLES WITH DIFFUSION MODELS

To generate labeled synthetic training samples, we train a conditional Latent Diffusion Model (LDM) (Rombach et al., 2022) with Classifier-Free Guidance (CFG) (Ho & Salimans, 2022) on the latent features of the images in the training set, which are generated by an off-the-shelf pretrained variational autoencoder (VAE) model from Stable Diffusion (Rombach et al., 2022). Detailed formulations of the training and inference of diffusion models, LDM, and CFG are deferred to

Section B.1 of the appendix. We use Diffusion Transformers (DiTs) (Peebles & Xie, 2023) as the backbones of the LDMs in our work. Let v_e and v_d be the fixed pretrained encoder and decoder. The encoder of the VAE is first applied to generate the latent features $\{h_i\}_{i=1}^N$ of $\mathcal{D}_{\text{real}}$, where $h_i = v_e(x_i)$ is the latent feature of the i -th image. The LDM is then trained on $\{h_i, y_i\}_{i=1}^N$ by minimizing the loss in Equation 10 in Section B.1 of the appendix. Algorithm 2 in Section B.1 of the appendix describes the training algorithm of the LDM. Once the training of the LDM is finished, the latent features $\{\hat{h}_i\}_{i=1}^M$ are generated for a set of pre-defined synthetic labels $\{\hat{y}_i\}_{i=1}^M$. The synthetic labels are obtained by uniformly sampling class indices from all the classes following the convention in existing GDA methods (Azizi et al., 2023; Trabucco et al., 2024), which ensures balanced coverage across all classes. The number of synthetic samples is decided by cross-validation as detailed in Section 4.1. The synthetic training images $\{\hat{x}_i\}_{i=1}^M$ are then generated by applying the pretrained decoder on the latent features $\{\hat{h}_i\}_{i=1}^M$, where $\hat{x}_i = v_d(\hat{h}_i)$. Algorithm 3 in Section B.1 of the appendix describes the generation process of the synthetic training set. After obtaining the synthetic training set \mathcal{D}_{syn} with the LDM, we combine it with the original training set $\mathcal{D}_{\text{real}}$ to obtain the augmented training set \mathcal{D}_{aug} . Next, the classifier network in IDR can be trained together with the sample reweighting network on \mathcal{D}_{aug} to be described in Section 3.3.

3.2 VARIATIONAL UPPER BOUND FOR THE IB LOSS

Although data reweighting methods have achieved remarkable success in image classification, existing data reweighting methods do not explicitly ensure that more informative training samples which contribute more to classification tasks receive higher importance weights when applied to GDA. As a result, synthetic training samples with noise introduced by the generative models, such as local visual artifacts (unnatural textures, color bleeding, or over-smoothed regions) (Corvi et al., 2023) or semantic inconsistencies (Azizi et al., 2023; He et al., 2023), could potentially degrade the performance of the classification model trained on the augmented training set. To address this issue, we propose a novel information-theoretic data reweighting method for GDA, termed Informative Data Reweighting (IDR), which encourages the sample reweighting network to assign higher importance weights to more informative training samples by explicitly reducing the IB loss on the augmented training set. Again, we note that our reweighting is applied to both real and synthetic training samples for maximal benefit of our reweighting scheme. IDR aims to reduce the IB loss, $\text{IB}(\tilde{Z}, \tilde{X}, \tilde{Y}) = I(\tilde{Z}, \tilde{X}) - I(\tilde{Z}, \tilde{Y})$, where $I(\cdot, \cdot)$ stands for the mutual information. \tilde{X} , \tilde{Z} , and \tilde{Y} denote the random variables representing the input feature, learned feature, and ground truth training class label of the samples in the augmented training set, respectively. Reduction of the IB loss ensures that the learned features of the augmented training set are more correlated with their class labels and less correlated with their input features. Lacking such an information-theoretic mechanism, existing sample reweighting methods, which do not explicitly learn informative importance weights, may introduce noisy information to classification tasks by assigning high importance weights to noisy training samples. By reducing the IB loss, IDR explicitly encourages more informative training samples in the augmented training set to receive higher importance weights by increasing the correlation of the learned features with class labels while decreasing their correlation with the noise in the training samples, ultimately leading to more discriminative features.

In order to assign higher importance weights to more informative training samples, we propose to train the reweighting network by explicitly reducing the IB loss on the augmented training set. To this end, we first derive a variational upper bound for the IB loss, which can be optimized by standard SGD algorithms. Given the augmented training set \mathcal{D}_{aug} , we first specify how to compute the IB loss, $\text{IB}(\Theta) = I(\tilde{Z}(\Theta), \tilde{X}) - I(\tilde{Z}(\Theta), \tilde{Y})$, where Θ is the weights of the classification network, \tilde{X} is a random variable representing the input feature which takes values in $\{\tilde{x}_i\}_{i=1}^{N'}$, $\tilde{Z}(\Theta)$ is a random variable representing the learned feature which takes values in $\{\tilde{z}_i(\Theta)\}_{i=1}^{N'}$ with $\tilde{z}_i(\Theta)$ being the learned feature for the i -th training sample in \mathcal{D}_{aug} . \tilde{Y} is a random variable representing the class label, which takes values in $\{\tilde{y}_i\}_{i=1}^{N'}$. We define $\mathcal{C}(\theta, \Theta) = \left\{ \left\{ c_k^{(\text{input})}(\theta) \right\}_{k=1}^C, \left\{ c_k^{(\text{feat})}(\theta, \Theta) \right\}_{k=1}^C \right\}$ as the class centroids of the input features and the learned features on \mathcal{D}_{aug} , where θ denotes the parameters of the sample reweighting network. The formulas for the computation of $\mathcal{C}(\theta, \Theta)$ can be found in Equation (2). We abbreviate $\tilde{Z}(\Theta)$ as \tilde{Z} , $c_k^{(\text{input})}(\theta)$ as $c_k^{(\text{input})}$, and $c_k^{(\text{feat})}(\theta, \Theta)$ as $c_k^{(\text{feat})}$.

Then we define the probability that \tilde{z}_i belongs to class a as $\Pr[\tilde{Z} \in a] = \frac{1}{N'} \sum_{i=1}^{N'} \phi(\tilde{z}_i, c_a^{(\text{feat})})$

with $\phi(\tilde{z}_i, c_a^{(\text{feat})}) = \frac{\exp(-\|\tilde{z}_i - c_a^{(\text{feat})}\|_2^2)}{\sum_{a'=1}^C \exp(-\|\tilde{z}_i - c_{a'}^{(\text{feat})}\|_2^2)}$ following (Wang et al., 2024; 2025). Similarly,

we define the probability that \tilde{x}_i belongs to class y as $\Pr[\tilde{X} \in y] = \frac{1}{N'} \sum_{i=1}^{N'} \phi(\tilde{x}_i, c_y^{(\text{input})})$.

Similarly, we define the joint probabilities $\Pr[\tilde{Z} \in a, \tilde{X} \in y] = \frac{1}{N'} \sum_{i=1}^{N'} \phi(\tilde{z}_i \in a, \tilde{x}_i \in y)$, where $\phi(\tilde{z}_i \in a, \tilde{x}_i \in y)$ is defined in (12) in the appendix. As a result, we can compute the mutual information by $I(\tilde{Z}, \tilde{X}) = \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{X} \in y] \log \frac{\Pr[\tilde{Z} \in a, \tilde{X} \in y]}{\Pr[\tilde{Z} \in a] \Pr[\tilde{X} \in y]}$,

$I(\tilde{Z}, \tilde{X}) = \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{X} \in y] \log \frac{\Pr[\tilde{Z} \in a, \tilde{X} \in y]}{\Pr[\tilde{Z} \in a] \Pr[\tilde{X} \in y]}$, and $I(\tilde{Z}, \tilde{Y}) =$

$\sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{\Pr[\tilde{Z} \in a, \tilde{Y} = y]}{\Pr[\tilde{Z} \in a] \Pr[\tilde{Y} = y]}$, and then compute the IB loss $\text{IB}(\Theta) =$

$\text{IB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}) = I(\tilde{Z}, \tilde{X}) - I(\tilde{Z}, \tilde{Y})$. Given a variational distribution $Q(\tilde{Z} \in a | \tilde{Y} = y)$ for

$y, a \in [C]$, the following theorem gives our new variational upper bound, $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}})$, for the IB loss $\text{IB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}})$.

Theorem 3.1. Let $\Pr[\tilde{X} \in y] = \sum_{i=1}^{N'} \mathbb{1}_{\{\tilde{y}_i = y\}} / N' := p_y$ be the prior probability for every $y \in [C]$.

Then

$$\text{IB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}) \leq \text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}), \quad (1)$$

where

$$\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}) := \frac{1}{N'} \sum_{i=1}^{N'} \text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \tilde{x}_i), \quad \text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \tilde{x}_i) := U_i - V_i.$$

Here

$$U_i := \sum_{a=1}^C \sum_{y=1}^C \phi(i, a, y) \log \left(\frac{\phi(i, a, y)}{p_y \phi(\tilde{z}_i, c_a^{(\text{feat})})} \right), \quad V_i := \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i = y\}} \log Q(\tilde{Z} \in a | \tilde{Y} = y),$$

where $\phi(i, a, y) = \phi(\tilde{z}_i \in a, \tilde{x}_i \in y)$. $Q(\tilde{Z} \in a | \tilde{Y} = y)$ is the variational conditional probability that \tilde{Z} belongs to class a given its label \tilde{Y} being y , which is computed efficiently by Algorithm 5 in the appendix. $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \tilde{x}_i)$ can be interpreted as the information bottleneck upper bound for the i -th sample in \mathcal{D}_{aug} . The proof of this theorem follows by applying Lemma E.1 and Lemma E.2 in Section E of the appendix. We remark that $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \tilde{x}_i)$ is ready to be optimized by standard SGD algorithms because it is the summation of losses on individual training points.

Efficient and Distribution-Free Variational Upper Bound. Since our novel variational upper bound is distribution-free, we demonstrate the computational efficiency of VUIB over the existing distribution-free upper bound for the IB loss in CLUB (Cheng et al., 2020). Let T_0 denote the computational cost of a forward and backward pass of the neural network for each individual training sample. The overall computational complexity for VUIB is $\Theta(N'T_0 + N'C)$. By comparison, computing the upper bound for the IB loss introduced in CLUB (Cheng et al., 2020) incurs a significantly higher computational complexity exceeding $\Theta(N'^2 T_0)$ since $N' \gg C$. Notably, the term $\Theta(N'^2 T_0)$ accounts solely for the computation of the upper bound for the mutual information $I(\tilde{Z}, \tilde{X})$. In addition, CLUB (Cheng et al., 2020) requires the evaluation of a lower bound on the mutual information $I(\tilde{Z}, \tilde{Y})$ in order to construct an upper bound on the IB loss $I(\tilde{Z}, \tilde{X}) - I(\tilde{Z}, \tilde{Y})$. A detailed analysis of the computational complexity of both CLUB and VUIB is provided in Section F of the appendix. Unlike existing upper bounds for the IB loss, such as VIB (Dai et al., 2018; Srivastava et al., 2021), which rely on an impractical Gaussian assumption over the hidden features of DNNs, and APIB (Guo et al., 2023), which only reduces an approximation of the IB loss, IBMA directly optimizes a variational upper bound of the IB loss, namely VUIB, without introducing distributional assumptions on the hidden features or approximations to the IB objective. As reported in Table 2 in Section 4.3, VUIB consistently achieves substantially better performance than the competing approaches, and requires much less training time than the distribution-free baseline method, CLUB (Cheng et al., 2020).

3.3 INFORMATIVE DATA REWEIGHTING (IDR)

IDR aims to assign higher importance weights to more informative training samples in the augmented training set that contribute more to the image classification task. To this end, IDR trains a sample reweighting network $g_\theta(\cdot)$ to learn importance weights $\{g_\theta(\tilde{x}_i) \in [0, 1]\}_{i=1}^{N'}$ of the samples in the augmented training set by reducing the IB loss on them, where $g_\theta(\cdot)$ is a DNN and θ denotes its parameters. Following (Shu et al., 2019), we employ a two-layer multilayer perceptron as the sample reweighting network, where a sigmoid activation is applied to the output to produce importance weights. We remark that the reweighting network plays a role similar to that of the meta-networks in (Shu et al., 2019; Jain et al., 2024), which generate the importance weights for training samples. Inspired by existing works (Wang et al., 2025), our reweighting scheme applies to both synthetic and real training samples.

To train the sample reweighting network $g_\theta(\cdot)$ to reduce the IB loss on \mathcal{D}_{aug} , such that more informative samples receive higher importance weights, we optimize $g_\theta(\cdot)$ by reducing the variational upper bound for the IB loss, VUIB. Let $f'_\Theta(\cdot)$ denote the representation learning backbone of the image classifier $f_\Theta(\cdot)$ excluding the last linear layer. To compute the VUIB on the augmented training set \mathcal{D}_{aug} , we first compute the class centroids for the input features and the learned features using all the samples in \mathcal{D}_{aug} . The class centroids for the input features and the learned features can be computed by

$$c_k^{(\text{input})}(\theta) = \frac{\sum_{i=1}^{N'} g_\theta(\tilde{x}_i) \tilde{x}_i \mathbb{I}_{\{\tilde{y}_i=k\}}}{\sum_{i=1}^{N'} g_\theta(\tilde{x}_i) \mathbb{I}_{\{\tilde{y}_i=k\}}}, \quad c_k^{(\text{feat})}(\theta, \Theta) = \frac{\sum_{i=1}^{N'} g_\theta(\tilde{x}_i) f'_\Theta(\tilde{x}_i) \mathbb{I}_{\{\tilde{y}_i=k\}}}{\sum_{i=1}^{N'} g_\theta(\tilde{x}_i) \mathbb{I}_{\{\tilde{y}_i=k\}}}, \quad (2)$$

where $k \in [C]$ is the class index and C is the number of classes, and $\mathbb{I}_{\{\cdot\}}$ is an indicator function. Next, the VUIB on \mathcal{D}_{aug} can be computed using Equation (1). With the sample reweighting network $g_\theta(\cdot)$, the overall training loss for the classifier $f_\Theta(\cdot)$ on the augmented training set \mathcal{D}_{aug} is $\mathcal{L}_{\text{train}}(\theta, \Theta, \mathcal{D}_{\text{aug}}) = \frac{1}{N'} \sum_{i=1}^{N'} g_\theta(\tilde{x}_i) \text{CE}(f_\Theta(\tilde{x}_i), \tilde{y}_i)$, where $\text{CE}(\cdot, \cdot)$ is the cross-entropy function. To train the classifier $f_\Theta(\cdot)$ by optimizing $\mathcal{L}_{\text{train}}(\theta, \Theta, \mathcal{D}_{\text{aug}})$ while training the sample reweighting network g_θ by optimizing VUIB($\theta, \Theta, \mathcal{D}_{\text{aug}}$), we formulate a bi-level optimization objective for IDR as

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}_{\text{train}}(\theta^*, \Theta, \mathcal{D}_{\text{aug}}), \quad \text{s.t. } \theta^* = \arg \min_{\theta} \text{VUIB}(\mathcal{C}(\theta, \Theta^*), \Theta^*, \mathcal{D}_{\text{aug}}), \quad (3)$$

where Θ^* and θ^* are the optimal parameters for the classifier $f_\Theta(\cdot)$ and the sample reweighting network $g_\theta(\cdot)$.

Optimization of IDR. To train the classifier $f_\Theta(\cdot)$ and the sample reweighting network $g_\theta(\cdot)$ with the optimization objective in Equation (3), we adopt an alternating stochastic gradient descent update strategy commonly used for solving bi-level optimization problems (Shu et al., 2019; Algan & Ulusoy, 2021; Jain et al., 2024). In the bi-level optimization framework, the lower level learns sample importance weights via a reweighting network, while the upper level trains the classifier using these weights to improve generalization. At the t -th epoch, the parameters of the sample reweighting network are first updated by $\theta^{(t)} = \theta^{(t-1)} - \eta_\theta \nabla_{\theta} \text{VUIB}(\mathcal{C}(\theta, \Theta^{(t-1)}), \Theta^{(t-1)}, \mathcal{D}_{\text{aug}})$, where η_θ is the learning rate of θ . $\theta^{(t)}$ and $\Theta^{(t)}$ are the parameters of the sample reweighting network and the classifier network at the t -th epoch. Next, the parameters of the classifier are updated by $\Theta^{(t)} = \Theta^{(t-1)} - \eta_\Theta \nabla_{\Theta} \mathcal{L}_{\text{train}}(\theta^{(t-1)}, \Theta, \mathcal{D}_{\text{aug}})$, where η_Θ is the learning rate of Θ . Since both VUIB and $\mathcal{L}_{\text{train}}$ are separable and amenable to mini-batch stochastic gradient descent (SGD), the entire optimization process of IDR can be efficiently conducted using mini-batch SGD. Algorithm 1 in Section A of the appendix describes the training process of IDR. $Q^{(t)}(\tilde{Z} \in a | \tilde{Y} = y)$ is the variational conditional probability that a feature \tilde{Z} belongs to class a given class label y at the t -th epoch, which is computed efficiently by Algorithm 5 in the appendix. It is worthwhile to mention that the class centroids in the input feature space can be efficiently pre-computed prior to training using FAISS (Douze et al., 2025), which incurs marginal training overhead as detailed in Section D.8, and it is not necessary to project the input features into a lower-dimensional space using methods such as VAEs (Kingma & Welling, 2014), which require additional training time.

4 EXPERIMENTS

We present a comprehensive evaluation of our proposed IDR. The implementation details of our experiments and the performance of IDR for image classification are presented in Section 4.1.

The impact of varying the number of synthetic training samples is investigated in Section 4.2. A comparison between the proposed variational upper bound for the IB loss and existing IB formulations is presented in Section 4.3. Additional experiment results are presented in Section D of the appendix. A qualitative analysis of the data reweighting with IDR using t-SNE visualization is conducted in Section D.1. Section D.2 compares IDR with recent conventional data augmentation methods, which do not rely on generative models, and shows that IDR is complementary to the conventional data augmentation methods. Section D.3 presents an ablation study demonstrating the advantages of jointly reweighting real and synthetic samples in IDR over reweighting synthetic samples alone. The transferability of IDR models pretrained on ImageNet-1K is further evaluated for object detection, instance segmentation, and fine-grained image classification tasks in Sections D.4, D.5, and D.6. The effect of different diffusion models used for synthetic data generation in IDR is examined in Section D.7. Ablation studies and training time analysis of the individual components of IDR, including the VUIB and the reweighting mechanism, are provided in Section D.8. The comparison between IDR and representative active learning methods is reported in Section D.9. The statistical significance of the performance improvements achieved by IDR is analyzed in Section D.10.

4.1 IMAGE CLASSIFICATION ON IMAGENET-1K

Implementation Details. We evaluate IDR for image classification on ImageNet-1K (Deng et al., 2009). ViT-S, ViT-B (Dosovitskiy et al., 2021), Swin-T, and Swin-B (Liu et al., 2021) are used as the base classification networks. All models are fine-tuned on the augmented training set from checkpoints pretrained on the original ImageNet-1K training set. Detailed settings for the fine-tuning process are presented in Section C of the appendix. DiT-B (Peebles & Xie, 2023) is employed for synthetic data generation, with detailed settings deferred to Section C of the appendix. The same set of synthetic images generated by DiT-B on ImageNet-1K is used for IDR and all baseline models in this section. The number of synthetic images for IDR and baselines is selected via five-fold cross-validation from $\{N, 2N, \dots, 10N\}$, where N is the size of the ImageNet-1K training set.

In this section, we compare the performance of competing data selection and data reweighting methods with our IDR for GDA on ImageNet-1K (Deng et al., 2009). In particular, we compare our IDR models with competitive data selection and sample reweighting methods, including IE (Chhabra et al., 2024), CBF (He et al., 2023), MW-Net (Shu et al., 2019), OTR (Guo et al., 2022), REVAR (Jain et al., 2024), IDS (Wang et al., 2025), and GenDataAgent (Li et al., 2025). Similar to IDR, all baseline methods reweight both the real and synthetic samples in the augmented training set. The number of synthetic samples added to the augmented training set for the baseline methods is determined via cross-validation, following the same setting as IDR. The results for IDR averaged over 10 runs with different random initializations (\pm standard deviation) are shown in Table 1, and statistical significance against the best baseline is reported in Appendix Section D.10. It is observed in Table 1 that IDR consistently outperforms all competing methods across different backbones. For instance, IDR achieves a top-1 accuracy of 85.7%, outperforming the strongest baseline, GenDataAgent, by 1.1%, when using Swin-B as the classification network, demonstrating the effectiveness of the IB-based data reweighting method on the augmented training set.

Table 1: Performance comparisons between IDR and existing data selection and data reweighting methods on ImageNet-1K. N denotes the size of the original ImageNet-1K training set.

Method	Backbone	Synthetic Data Size ($\times N$)	Top-1
ViT-S/16	ViT-S/16	-	81.2
MW-Net (Shu et al., 2019)		2	81.8
OTR (Guo et al., 2022)		3	81.6
IE (Chhabra et al., 2024)		2	81.7
CBF (He et al., 2023)		2	81.9
REVAR (Jain et al., 2024)		3	81.9
IDS (Wang et al., 2025)		3	82.1
GenDataAgent (Li et al., 2025)		3	82.1
IDR (Ours)		4	83.2 \pm 0.18
ViT-B/16		ViT-B/16	-
MW-Net (Shu et al., 2019)	2		84.1
OTR (Guo et al., 2022)	2		84.0
IE (Chhabra et al., 2024)	3		84.0
CBF (He et al., 2023)	2		84.4
REVAR (Jain et al., 2024)	2		84.2
IDS (Wang et al., 2025)	3		84.4
GenDataAgent (Li et al., 2025)	3		84.5
IDR (Ours)	5		85.6 \pm 0.22
Swin-T	Swin-T		-
MW-Net (Shu et al., 2019)		2	81.9
OTR (Guo et al., 2022)		2	81.6
IE (Chhabra et al., 2024)		2	81.9
CBF (He et al., 2023)		2	82.2
REVAR (Jain et al., 2024)		2	82.0
IDS (Wang et al., 2025)		3	82.2
GenDataAgent (Li et al., 2025)		3	82.3
IDR (Ours)		4	83.4 \pm 0.27
Swin-B		Swin-B	-
MW-Net (Shu et al., 2019)	2		83.9
OTR (Guo et al., 2022)	2		84.0
IE (Chhabra et al., 2024)	2		83.9
CBF (He et al., 2023)	2		84.4
REVAR (Jain et al., 2024)	3		84.3
IDS (Wang et al., 2025)	3		84.5
GenDataAgent (Li et al., 2025)	3		84.6
IDR (Ours)	5		85.7 \pm 0.19

4.2 IMPACT OF SYNTHETIC SAMPLE QUANTITY

We conduct an ablation study by progressively increasing the number of synthetic samples in the augmented training set, ranging from N to $10N$, where N denotes the size of the original training set. ViT-B is used as the classification network. It can be observed from Figure 2 that although a moderate amount of synthetic samples benefits the classification accuracy, the performance of all the competing baselines and IDR eventually decreases with excessive synthetic samples due to the inherent noise in the synthetic samples (Azizi et al., 2023; Trabucco et al., 2024; Na et al., 2024). However, our IDR still exhibits a clear advantage over competing baselines. As illustrated in Figure 2, both CBF and GenDataAgent achieve their best performance with a relatively small amount of synthetic samples. In particular, CBF and GenDataAgent perform best with $2N$ and $3N$ synthetic samples, respectively. Both CBF and GenDataAgent exhibit a significant decrease in top-1 accuracy once the number of synthetic samples exceeds $3N$, indicating that noise in synthetic samples significantly degrades the classification network. In a strong contrast, IDR is more robust to noise in the generated synthetic samples compared with the baselines, and benefits from a larger amount of synthetic samples. In particular, the performance of IDR continues to improve until the number of synthetic samples reaches $5N$. Importantly, the top-1 accuracy of IDR with excessive synthetic samples ($> 5N$) still significantly surpasses the best accuracy of all baselines, with only a marginal decrease from the peak performance at $5N$. It can also be observed from Figure 2 that IDR consistently outperforms CBF and GenDataAgent when the same amount of synthetic samples is used, ranging from N to $10N$.

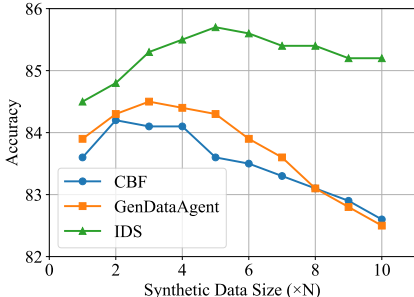


Figure 2: Performance comparison on ImageNet-1K with different amounts of synthetic training data in the augmented training set.

4.3 COMPARISON BETWEEN VUIB AND EXISTING UPPER BOUNDS FOR THE IB LOSS

We compare the proposed variational upper bound for the IB loss, VUIB, with existing works deriving the upper bound for the IB loss (Cheng et al., 2020; Dai et al., 2018; Srivastava et al., 2021). To compare the performance of VUIB with VIB, APIB, and CLUB, we conduct an ablation study by replacing VUIB with VIB, APIB, and CLUB in IDR for image classification on ImageNet-1K. The training time is evaluated on the augmented training set of ImageNet-1K with $5N$ synthetic samples, where N is the size of the original training set. It is observed in Table 2 that the unrealistic Gaussian distribution assumption on the hidden features imposed by VIB (Dai et al., 2018) and the IB approximation in APIB (Guo et al., 2023) lead to degraded performance compared with our reduction of VUIB. For instance, the model based on VUIB outperforms the model based on APIB by 1.2% in top-1 classification accuracy. In addition, the model based on VUIB outperforms the model based on CLUB by 1.0% in accuracy while using only 49.7% of the training time of the CLUB-based model.

Table 2: Comparison of different methods for the training overhead incurred by IB bound. ViT-B is used as the base model. The training time is evaluated on four NVIDIA A100 GPUs.

Methods	Top-1	Training Time (Hours / Epoch)
VIB (Dai et al., 2018)	84.2	1.69
APIB (Guo et al., 2023)	84.4	1.70
CLUB (Cheng et al., 2020)	84.6	3.47
VUIB (Ours)	85.6	1.72

5 CONCLUSION

In this paper, we propose Informative Data Reweighting (IDR), a novel method designed to reweight samples in the augmented training set for Generative Data Augmentation (GDA) based on an information-theoretic measure, the Information Bottleneck (IB). IDR trains a sample reweighting network to reduce the IB loss on the augmented training set, such that the IB principle, learning features more correlated with the outputs and less correlated with the inputs, is better adhered to. Extensive experiments and ablation studies demonstrate that IDR successfully assigns higher weights to more informative samples in the augmented training set for image classification, and significantly outperforms existing data selection and data reweighting methods for GDA.

REFERENCES

- 432
433
434 Chirag Agarwal, Daniel D’souza, and Sara Hooker. Estimating example difficulty using variance of
435 gradients. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022,*
436 *New Orleans, LA, USA, June 18-24, 2022*, pp. 10358–10368. IEEE, 2022.
- 437 Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information
438 bottleneck. In *5th International Conference on Learning Representations, ICLR 2017, Toulon,*
439 *France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 440 Görkem Algan and Ilkay Ulusoy. Meta soft label generation for noisy labels. In *2020 25th Interna-*
441 *tional Conference on Pattern Recognition (ICPR)*, pp. 7142–7148. IEEE, 2021.
- 442 Omar Ali, Wiem Abdelbaki, Anup Shrestha, Ersin Elbasi, Mohammad Abdallah Ali Alryalat, and
443 Yogesh K Dwivedi. A systematic literature review of artificial intelligence in the healthcare sector:
444 Benefits, challenges, methodologies, and functionalities. *Journal of Innovation & Knowledge*, 8
445 (1):100333, 2023.
- 446 Rana Ali Amjad and Bernhard C. Geiger. Learning representations for neural network-based
447 classification using the information bottleneck principle. *IEEE Trans. Pattern Anal. Mach. Intell.*,
448 42(9):2225–2239, 2020.
- 449 Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J. Fleet.
450 Synthetic data from diffusion models improves imagenet classification. *Trans. Mach. Learn. Res.*,
451 2023, 2023.
- 452 Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated
453 datasets. *arXiv preprint arXiv:2302.02503*, 2023.
- 454 Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: high quality object detection and instance
455 segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(5):1483–1498, 2021.
- 456 Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang Liu, Zhe Gan, and Lawrence Carin. CLUB: A
457 contrastive log-ratio upper bound of mutual information. In *Proceedings of the 37th International*
458 *Conference on Machine Learning, ICML 2020*, volume 119, pp. 1779–1788. PMLR, 2020.
- 459 Anshuman Chhabra, Peizhao Li, Prasant Mohapatra, and Hongfu Liu. ”what data benefits my
460 classifier?” enhancing model performance and interpretability through influence-based data
461 selection. In *The Twelfth International Conference on Learning Representations*, 2024.
- 462 Riccardo Corvi, Davide Cozzolino, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. Intriguing
463 properties of synthetic images: from generative adversarial networks to diffusion models. In
464 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023 - Workshops,*
465 *Vancouver, BC, Canada, June 17-24, 2023*, pp. 973–982. IEEE, 2023.
- 466 Bin Dai, Chen Zhu, Baining Guo, and David P. Wipf. Compressing neural networks using the
467 variational information bottleneck. In *Proceedings of the 35th International Conference on*
468 *Machine Learning, ICML 2018*, volume 80, pp. 1143–1152. PMLR, 2018.
- 469 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
470 hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and*
471 *Pattern Recognition*, pp. 248–255. IEEE, 2009.
- 472 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
473 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
474 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
475 In *International Conference on Learning Representations*, 2021.
- 476 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel
477 Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *IEEE Transactions on*
478 *Big Data*, 2025.
- 479 Laila El Jiani, Sanaa El Filali, et al. Overcome medical image data scarcity by data augmentation
480 techniques: A review. In *2022 International Conference on Microelectronics (ICM)*, pp. 21–24.
481 IEEE, 2022.
- 482
483
484
485

- 486 Andre Esteva, Katherine Chou, Serena Yeung, Nikhil Naik, Ali Madani, Ali Mottaghi, Yun Liu, Eric
487 Topol, Jeff Dean, and Richard Socher. Deep learning-enabled medical computer vision. *NPJ*
488 *digital medicine*, 4(1):5, 2021.
- 489 Ruibin Feng, Zongwei Zhou, Michael B Gotway, and Jianming Liang. Parts2whole: Self-supervised
490 contrastive learning via reconstruction. In *MICCAI Workshop on Domain Adaptation and Repre-*
491 *sentation Transfer*, pp. 85–95. Springer, 2020.
- 493 Mingfei Gao, Zizhao Zhang, Guo Yu, Sercan Ö Arik, Larry S Davis, and Tomas Pfister. Consistency-
494 based semi-supervised active learning: Towards minimizing labeling cost. In *European Conference*
495 *on Computer Vision*, pp. 510–526. Springer, 2020.
- 496 Dandan Guo, Zhuo Li, He Zhao, Mingyuan Zhou, Hongyuan Zha, et al. Learning to re-weight
497 examples with optimal transport for imbalanced classification. *Advances in Neural Information*
498 *Processing Systems*, 35:25517–25530, 2022.
- 500 Song Guo, Lei Zhang, Xiawu Zheng, Yan Wang, Yuchao Li, Fei Chao, Chenglin Wu, Shengchuan
501 Zhang, and Rongrong Ji. Automatic network pruning via hilbert-schmidt independence criterion
502 lasso under information bottleneck principle. In *IEEE/CVF International Conference on Computer*
503 *Vision, ICCV 2023*, pp. 17412–17423. IEEE, 2023.
- 504 Yeho Gwon, Sehyun Hwang, Hoyoung Kim, Jungseul Ok, and Suha Kwak. Enhancing cost efficiency
505 in active learning with candidate set query. *Trans. Mach. Learn. Res.*, 2025, 2025.
- 506 Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip H. S. Torr, Song Bai, and
507 Xiaojuan Qi. Is synthetic data from generative models ready for image recognition? In *The*
508 *Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May*
509 *1-5, 2023*, 2023.
- 511 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
512 2022.
- 513 Nishant Jain, Karthikeyan Shanmugam, and Pradeep Shenoy. Learning model uncertainty as variance-
514 minimizing instance weights. In *The Twelfth International Conference on Learning Representations*,
515 2024.
- 516 Minsoo Kang, Seong-Wan Lee, and Suhyun Kim. Model-agnostic and efficient mixup augmentation
517 guided by saliency maps. *IEEE Trans. Pattern Anal. Mach. Intell.*, 47(8):6946–6958, 2025.
- 519 Kenji Kawaguchi, Zhun Deng, Xu Ji, and Jiaoyang Huang. How does information bottleneck help
520 deep learning? In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan
521 Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023,*
522 *23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning*
523 *Research*, pp. 16049–16096. PMLR, 2023.
- 524 Bum Jun Kim and Sang Woo Kim. Configuring data augmentations to reduce variance shift in
525 positional embedding of vision transformers. In Toby Walsh, Julie Shah, and Zico Kolter (eds.),
526 *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25*
527 *- March 4, 2025, Philadelphia, PA, USA*, pp. 17840–17849. AAAI Press, 2025.
- 528 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann
529 LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB,*
530 *Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- 532 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly,
533 and Neil Houlsby. Big transfer (bit): General visual representation learning. In Andrea Vedaldi,
534 Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *Computer Vision - ECCV 2020 - 16th*
535 *European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part V*, volume 12350 of
536 *Lecture Notes in Computer Science*, pp. 491–507. Springer, 2020.
- 537 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
538 categorization. In *2013 IEEE International Conference on Computer Vision Workshops, ICCV*
539 *Workshops 2013, Sydney, Australia, December 1-8, 2013*, pp. 554–561. IEEE Computer Society,
2013.

- 540 Dan Kushnir and Luca Venturi. Diffusion-based sampling for deep active learning. In *2023 International Conference on Sampling Theory and Applications (SampTA)*, pp. 1–9, 2023. doi: 10.1109/SampTA59647.2023.10301392.
- 541
542
543 Qiuxia Lai, Yu Li, Ailing Zeng, Minhao Liu, Hanqiu Sun, and Qiang Xu. Information bottleneck approach to spatial attention learning. In Zhi-Hua Zhou (ed.), *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pp. 779–785. ijcai.org, 2021.
- 544
545
546
547 Shiye Lei, Hao Chen, Sen Zhang, Bo Zhao, and Dacheng Tao. Image captions are natural prompts for text-to-image models. *arXiv preprint arXiv:2307.08526*, 2023.
- 548
549
550 Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Big-datasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21330–21340, 2022.
- 551
552
553 Zhiteng Li, Lele Chen, Jerone T. A. Andrews, Yunhao Ba, Yulun Zhang, and Alice Xiang. Gen-dataagent: On-the-fly dataset augmentation with synthetic data. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025.
- 554
555
556
557 Shaobo Lin, Kun Wang, Xingyu Zeng, and Rui Zhao. Explore the power of synthetic data on few-shot object detection. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023 - Workshops, Vancouver, BC, Canada, June 17-24, 2023*, pp. 638–647. IEEE, 2023.
- 558
559
560
561 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 2014.
- 562
563
564 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 565
566
567 Sangwoo Mo, Chiheon Kim, Sungwoong Kim, Minsu Cho, and Jinwoo Shin. Mining gold samples for conditional gans. *Advances in Neural Information Processing Systems*, 32, 2019.
- 568
569
570
571 Byeonghu Na, Yeongmin Kim, HeeSun Bae, Jung Hyun Lee, Se Jung Kwon, Wanmo Kang, and Il chul Moon. Label-noise robust diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.
- 572
573
574 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Sixth Indian Conference on Computer Vision, Graphics & Image Processing, ICVGIP 2008, Bhubaneswar, India, 16-19 December 2008*, pp. 722–729. IEEE Computer Society, 2008.
- 575
576
577
578 Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, pp. 3498–3505. IEEE Computer Society, 2012.
- 579
580
581 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 582
583
584
585 W Nicholson Price and I Glenn Cohen. Privacy in the age of medical big data. *Nature medicine*, 25(1):37–43, 2019.
- 586
587
588
589 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 2021.
- 590
591
592
593 Kama Ramudu, V Murali Mohan, D Jyothirmai, DVSSSV Prasad, Ruchi Agrawal, and Sampath Boopathi. Machine learning and artificial intelligence in disease prediction: Applications, challenges, limitations, case studies, and future directions. In *Contemporary Applications of Data Fusion for Advanced Healthcare Informatics*, pp. 297–318. IGI Global, 2023.

- 594 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
595 resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer*
596 *Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp.
597 10674–10685. IEEE, 2022.
- 598
599 Mert Bülent Saryıldız, Karteek Alahari, Diane Larlus, and Yannis Kalantidis. Fake it till you make
600 it: Learning transferable representations from synthetic imagenet clones. In *Proceedings of the*
601 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8011–8021, June
602 2023.
- 603 Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-
604 weight-net: Learning an explicit mapping for sample weighting. *Advances in neural information*
605 *processing systems*, 32, 2019.
- 606 Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In
607 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5972–5981, 2019.
608
- 609 Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy
610 labels with deep neural networks: A survey. *IEEE Trans. Neural Networks Learn. Syst.*, 34(11):
611 8135–8153, 2023.
- 612 Ayush Srivastava, Oshin Dutta, Jigyasa Gupta, Sumeet Agarwal, and Prathosh AP. A variational
613 information bottleneck based method to compress sequential networks for human action recognition.
614 In *IEEE Winter Conference on Applications of Computer Vision, WACV 2021*, pp. 2744–2753.
615 IEEE, 2021.
- 616 Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic
617 images from text-to-image models make strong visual representation learners. *arXiv preprint*
618 *arXiv:2306.00984*, 2023.
- 619 Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*
620 *preprint physics/0004057*, 2000.
- 621
622 Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective data augmenta-
623 tion with diffusion models. In *The Twelfth International Conference on Learning Representations,*
624 *ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- 625 Joshua Vendrow, Saachi Jain, Logan Engstrom, and Aleksander Madry. Dataset interfaces: Diagnos-
626 ing model failures using controllable counterfactual generation. *arXiv preprint arXiv:2302.07865*,
627 2023.
- 628
629 Yancheng Wang, Ping Li, and Yingzhen Yang. Visual transformer with differentiable channel
630 selection: An information bottleneck inspired approach. In *Forty-first International Conference on*
631 *Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- 632 Yancheng Wang, Rajeev Goel, Marko Jojic, Alvin C. Silva, Teresa Wu, and Yingzhen Yang. Informa-
633 tive synthetic data generation for thorax disease classification. In *Conference on Uncertainty in*
634 *Artificial Intelligence, Rio Othon Palace, Rio de Janeiro, Brazil, 21-25 July 2025*, volume 286 of
635 *Proceedings of Machine Learning Research*, pp. 4489–4514. PMLR, 2025.
- 636
637 Junfei Xiao, Yutong Bai, Alan Yuille, and Zongwei Zhou. Delving into masked autoencoders for
638 multi-label thorax disease classification. In *Proceedings of the IEEE/CVF Winter Conference on*
639 *Applications of Computer Vision*, pp. 3588–3600, 2023.
- 640 Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing
641 for scene understanding. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair
642 Weiss (eds.), *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany,*
643 *September 8-14, 2018, Proceedings, Part V*, volume 11209 of *Lecture Notes in Computer Science*,
644 pp. 432–448. Springer, 2018.
- 645
646 Jianan Yang, Haobo Wang, Sai Wu, Gang Chen, and Junbo Zhao. Towards controlled data augmenta-
647 tions for active learning. In *International Conference on Machine Learning*, pp. 39524–39542.
PMLR, 2023.

648 Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF*
649 *conference on computer vision and pattern recognition*, pp. 93–102, 2019.

650
651 Jianhao Yuan, Francesco Pinto, Adam Davies, Aarushi Gupta, and Philip Torr. Not just pretty
652 pictures: Text-to-image generators enable interpretable interventions for robust representations.
653 *arXiv preprint arXiv:2212.11237*, 2022.

654
655 Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio
656 Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort.
657 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.

658
659 Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba.
660 Semantic understanding of scenes through the ADE20K dataset. *Int. J. Comput. Vis.*, 127(3):
661 302–321, 2019.

662
663 Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and
664 Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference*
665 *on Machine Learning*, pp. 27378–27394. PMLR, 2022.

666
667 Yongchao Zhou, Hshmat Sahak, and Jimmy Ba. Training on thin air: Improve image classification
668 with generated data. *arXiv preprint arXiv:2305.15316*, 2023.

669 A ALGORITHM OF IDR

670
671 Algorithm 1 describes the training process of IDR.

672 **Algorithm 1** Algorithm of IDR

673 **Input:** The augmented training set \mathcal{D}_{aug} , epoch number t_{max} , the learning rates η_{θ} and η_{Θ} .

- 674 1: Initialize the classifier network parameters $\Theta^{(0)}$ and the sample reweighting network parameters $\theta^{(0)}$.
 - 675 2: **for** $t = 1, 2, \dots, t_{\text{max}}$ **do**
 - 676 3: Compute the class centroids of the input features and image representations $\mathcal{C}(\theta^{(t-1)}, \Theta^{(t-1)})$ by
677 Equation (2).
 - 678 4: Update $\theta^{(t)}$ by applying mini-batch gradient descent on \mathcal{D}_{aug} using $\theta^{(t)} = \theta^{(t-1)} -$
679 $\eta_{\theta} \nabla_{\theta} \text{VUIB}(\mathcal{C}(\theta, \Theta^{(t-1)}), \Theta^{(t-1)}, \mathcal{D}_{\text{aug}})$.
 - 680 5: Update $\Theta^{(t)}$ by applying mini-batch gradient descent on \mathcal{D}_{aug} using $\Theta^{(t)} = \Theta^{(t-1)} -$
681 $\eta_{\Theta} \nabla_{\Theta} \mathcal{L}_{\text{train}}(\theta^{(t-1)}, \Theta, \mathcal{D}_{\text{aug}})$.
 - 682 6: Compute $Q^{(t)}(\tilde{Z} \in a | \tilde{Y} = y)$ by Algorithm 5 in Section F of the appendix.
 - 683 7: **end for**
 - 684 8: **return** The trained weights Θ of the classifier network and the trained weights θ of the sample reweighting
685 network.
-

687 B INFORMATION ON DIFFUSION MODELS

688 B.1 FORMULATIONS OF DIFFUSION MODELS

689
690 **Diffusion models (DMs)** are probabilistic latent-variable frameworks that represent data x^0 through
691 a Markovian transformation from a highly noisy variable x_T back to the original signal, while
692 preserving identical dimensionality across all intermediate states. They are characterized by two
693 coupled Markov chains: a forward noising procedure given by $q(x^{(1:T)} | x^0) = \prod_{t=1}^T q(x^{(t)} | x^{(t-1)})$
694 and a learnable reverse-time denoising chain formulated as $p_{\omega}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\omega}(x^{(t-1)} |$
695 $x^{(t)})$. In the forward direction, Gaussian perturbations are incrementally injected into the data:
696

$$697 \quad q(x^{(t)} | x^{(t-1)}) = \mathcal{N}(x^{(t)}; \sqrt{1 - \beta^{(t)}}x^{(t-1)}, \beta^{(t)}\mathbf{I}), \quad (4)$$

698
699 where the noise schedule $\beta^{(1:T)}$ controls the variance added at each diffusion step t . This schedule is
700 selected such that the terminal distribution x_T closely follows a standard normal distribution, i.e.,
701 $q(x_T) \approx \mathcal{N}(0, \mathbf{I})$. Once specified, the forward diffusion process q remains fixed and is not learned.

Sample synthesis in DMs relies on training a parametric reverse process that progressively removes noise from $x_{T:1}$ to recover the clean sample x^0 :

$$p_\omega(x^{(t-1)} | x^{(t)}) = \mathcal{N}(x^{(t-1)}; \mu_\omega(x^{(t)}, t), (\rho^{(t)})^2 \mathbf{I}), \quad (5)$$

with the prior $p(x_T)$ defined as $\mathcal{N}(0, \mathbf{I})$. In practice, neural architectures such as U-Nets or Transformers are employed to model the mean function μ_ω , while the variance $\rho^{(t)}$ is typically fixed in advance.

For learning, the forward chain $q(x^{(1:T)} | x^0)$ is treated as a known posterior, and the parameters of the reverse process $p_\omega(x_{0:T})$ are optimized to maximize a variational lower bound on the data likelihood. Since directly optimizing the likelihood objective often causes numerical instability, a commonly adopted and more stable surrogate loss is:

$$\mathcal{L}_{\text{DM}} = \mathbb{E}_{x^0, \varepsilon \sim \mathcal{N}(0, \mathbf{I}), t} \left\| \varepsilon - \varepsilon_\omega(x^{(t)}, t) \right\|_2^2, \quad (6)$$

where ε_ω denotes a neural predictor of the injected noise ε , enabling the model to infer $x^{(t-1)}$ from $x^{(t)}$ at each timestep. After training, new samples are obtained via iterative ancestral sampling:

$$x^{(t-1)} = \frac{1}{\sqrt{1 - \beta^{(t)}}} \left(x^{(t)} - \frac{\beta^{(t)}}{\sqrt{1 - (\alpha^{(t)})^2}} \varepsilon_\omega(x^{(t)}, t) + \rho^{(t)} \varepsilon \right), \quad (7)$$

initialized from the Gaussian prior $x_T \sim p(x_T) = \mathcal{N}(x_T; 0, \mathbf{I})$.

Latent Diffusion Models (LDMs) extend conventional diffusion models by performing the diffusion process in a compressed latent space, thereby significantly reducing computational cost. Specifically, the original data x^0 is first mapped into a latent representation h^0 using an encoder. Diffusion in the latent domain follows:

$$q(h^{(t)} | h^{(t-1)}) = \mathcal{N}(h^{(t)}; \sqrt{1 - \beta^{(t)}} h^{(t-1)}, \beta^{(t)} I), \quad (8)$$

while the corresponding reverse-time dynamics aim to reconstruct the clean latent code h^0 from h_T :

$$p_\omega(h^{(t-1)} | h^{(t)}) = \mathcal{N}(h^{(t-1)}; \mu_\omega(h^{(t)}, t), (\rho^{(t)})^2 I), \quad (9)$$

after which h^0 is decoded back into the original data space. The optimization objective for LDMs is defined as

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{h_e(x), \varepsilon \sim \mathcal{N}(0, I), t} \left\| \varepsilon - \varepsilon_\omega(h^{(t)}, t, y) \right\|_2^2, \quad (10)$$

which mirrors the noise-prediction formulation used in standard diffusion models.

Classifier-Free Guidance (CFG) enhances conditional generation by combining conditional and unconditional noise estimators during sampling, allowing explicit control over class-conditional behavior without an auxiliary classifier. When applied to LDMs, the guided sampling update is expressed as:

$$h^{(t-1)} = \frac{1}{\sqrt{1 - \beta^{(t)}}} \left(h^{(t)} - \frac{\beta^{(t)}}{\sqrt{1 - (\alpha^{(t)})^2}} \tilde{\varepsilon}^{(t)} + \rho^{(t)} \varepsilon \right), \quad (11)$$

where $\tilde{\varepsilon}^{(t)} = (1 + \omega) \varepsilon_\omega(h^{(t)}, y, t) - \gamma \varepsilon_\omega(h^{(t)}, t)$, and γ controls the strength of guidance, steering the sampling trajectory toward desired conditions.

Algorithm 2 outlines the training procedure for the LDM, while Algorithm 3 details the generation process used to construct the synthetic training set.

C DETAILED EXPERIMENT SETTINGS

We evaluate IDR for general image classification on ImageNet-1K (Deng et al., 2009), which is a large-scale general image classification benchmark containing approximately 1.28 million training images and 50,000 validation images across 1,000 classes. We fine-tune ViT-S, ViT-B (Dosovitskiy et al., 2021), Swin-T, and Swin-B (Liu et al., 2021) on the augmented ImageNet-1K training set, starting from checkpoints pretrained on the original ImageNet-1K training set. The fine-tuning is

Algorithm 2 Training Procedure for the LDM

Input: The original training dataset $\mathcal{D}_{\text{real}} = \{x_i, y_i\}_{i=1}^N$, the encoder v_e from a fixed pretrained VAE, and the total number of LDM training epochs t_{LDM} .

Output: The learned parameters ω of the LDM.

- 1: Initialize the LDM parameters ω .
- 2: Map the input samples $\{x_i\}_{i=1}^N$ into latent representations $\{h_i\}_{i=1}^N$ using the encoder v_e , i.e., $h_i = v_e(x_i)$.
- 3: **for** $t = 1, 2, \dots, t_{\text{LDM}}$ **do**
- 4: Optimize ω via mini-batch stochastic gradient descent over $\{h_i\}_{i=1}^N$ by minimizing the loss \mathcal{L}_{LDM} defined in Equation (10).
- 5: **end for**
- 6: **return** The optimized LDM parameters ω .

Algorithm 3 Synthetic Training Set Generation

Input: The target labels of the synthetic dataset $\{\hat{y}_i\}_{i=1}^M$, the trained LDM parameters ω , and the decoder v_d from the fixed pretrained VAE.

Output: The generated synthetic dataset $\mathcal{D}_{\text{syn}} = \{\hat{x}_i, \hat{y}_i\}_{i=1}^M$.

- 1: **for** $i = 1, 2, \dots, M$ **do**
- 2: Draw a noise vector $\epsilon \sim \mathcal{N}(0, I)$.
- 3: Produce a synthetic latent representation \hat{h}_i from ϵ using the LDM according to Equation (9) in Section B of the appendix.
- 4: Recover the synthetic input feature by decoding the latent code, $\hat{x}_i = v_d(\hat{h}_i)$.
- 5: **end for**
- 6: **return** The synthetic training dataset $\mathcal{D}_{\text{syn}} = \{\hat{x}_i, \hat{y}_i\}_{i=1}^M$.

performed for 100 epochs using the AdamW optimizer, with an initial learning rate of 1×10^{-4} , weight decay of 1×10^{-4} , and momentum parameters $(\beta_1, \beta_2) = (0.9, 0.999)$. We adopt a cosine learning-rate decay schedule with a five-epoch linear warm-up. All experiments are conducted on four NVIDIA A100 GPUs with a total batch size of 1024, corresponding to 256 images per GPU. Input images are resized to 224×224 , and standard ImageNet data augmentation is applied, including random resized cropping, horizontal flipping, and color jittering. Label smoothing with a factor of 0.1 is used during training. We apply the Exponential Moving Average (EMA) to model parameters with a decay rate of 0.9999 for evaluation. The parameters of the sample reweighting network are optimized by the AdamW optimizer with the learning rate fixed to 1×10^{-5} . The hidden dimension of the sample reweighting network is set to 512.

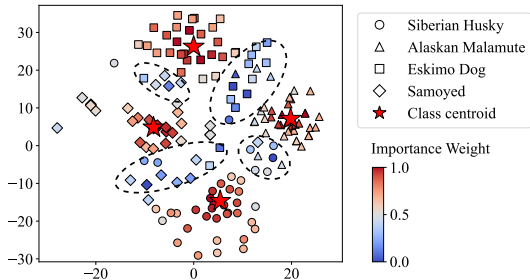
For synthetic data generation, we adopt the DiT-B backbone (Peebles & Xie, 2023). The diffusion model is trained on the latent representations encoded by the pretrained VAE from Stable Diffusion using classifier-free guidance, optimized with a mean-squared error noise prediction objective for 2,800 epochs using AdamW with a learning rate of 1×10^{-4} following (Peebles & Xie, 2023). During inference, synthetic samples are generated with a guidance scale of 4.0 and 128 sampling steps, and the resulting latent features are decoded back to the image space using the fixed VAE decoder. The number of synthetic images added to the augmented training set for IDR and the competing baseline methods is selected by five-fold cross-validation from $\{N, 2N, \dots, 10N\}$, where N denotes the number of training samples in the original training set of ImageNet-1K.

D ADDITIONAL EXPERIMENT RESULTS

D.1 T-SNE VISUALIZATION ANALYSIS

To qualitatively examine how IDR reweights synthetic samples in the augmented training set, we illustrate the importance weights of sampled synthetic images from four visually similar dog classes in ImageNet-1K using t-SNE in Figure 3, where marker shapes indicate classes and colors denote importance weights. Equation (2) shows that the class centroids are largely decided by samples with high importance weights. This is reflected by the t-SNE visualization, which shows that synthetic images located closer to the class centroids tend to receive higher importance weights, whereas samples located near class boundaries or closer to the centroids of other classes are assigned lower importance weights. For instance, the samples in the regions highlighted by the dashed circles correspond to marginal areas between classes and exhibit reduced importance weights.

810
811
812
813
814
815
816
817
818

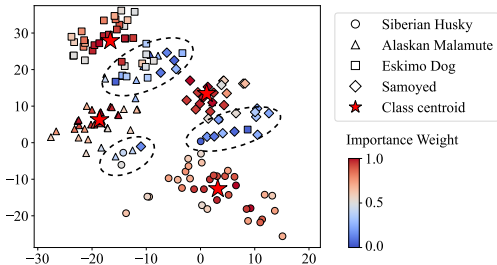


819 Figure 3: Visualization of the importance weights learned by IDR on synthetic images. For each class,
820 40 synthetic images are randomly sampled from the augmented training set. Different marker shapes
821 denote different classes, corresponding to four visually similar dog categories from ImageNet-1K.
822 The color of each marker indicates its importance weight.

824
825
826
827
828
829

To qualitatively analyze how IDR reweights real training samples, we further visualize the importance weights of sampled real images from four visually similar dog classes in ImageNet-1K using t-SNE, as illustrated in Figure 4. The observed distribution is similar to that for synthetic samples in Figure 3 in Section D.1 of the main paper. Real samples located closer to their class centroids usually receive higher importance weights, whereas samples near class boundaries or closer to the centroids of other classes are assigned lower importance weights.

830
831
832
833
834
835
836
837
838
839



841 Figure 4: Visualization of the importance weights learned by IDR on real images. For each class,
842 40 real images are randomly sampled from the augmented training set and projected into a two-
843 dimensional space using t-SNE. Different marker shapes denote different classes, corresponding to
844 four visually similar dog categories from ImageNet-1K, including Siberian Husky, Alaskan Malamute,
845 Eskimo Dog, and Samoyed, which are challenging to distinguish. The color of each point indicates
846 its importance weight, with redder colors indicating higher importance and bluer colors indicating
847 lower importance.

848
849
850

D.2 COMPARISON WITH CONVENTIONAL DATA AUGMENTATION METHODS

851
852
853
854
855
856
857
858
859
860
861
862
863

Recent works have proposed better conventional data augmentation strategies to improve image classification performance (Kim & Kim, 2025; Kang et al., 2025) without using generative models. In particular, Position-Aware Augmentation (PAA) proposes to adapt existing data augmentation methods, such as random resize crop, Mixup, CutMix, and random erasing, to explicitly account for the sensitivity of Vision Transformer positional embeddings to spatial transformations. In parallel, GuidedMixup (Kang et al., 2025) proposes saliency-guided mixup by optimizing sample pairings and pixel-wise mixing ratios to preserve semantically salient regions and reduce label noise during data mixing. These conventional data augmentation strategies do not rely on GDA and are orthogonal to IDR, as they operate directly on the original training samples without selection or reweighting. As shown in Table 3, IDR outperforms both PAA and GuidedMixup, demonstrating the effectiveness of GDA with the proposed informative data reweighting. For instance, IDR outperforms PAA by 1.1% in top-1 classification accuracy on ImageNet-1K. Moreover, combining IDR with either PAA or GuidedMixup further improves performance. Specifically, IDR+PAA achieves a 0.5% improvement over IDR alone and a 1.6% improvement over PAA alone, demonstrating that IDR is complementary

864 to conventional data augmentation strategies and can be seamlessly integrated with them to yield
865 additive performance gains.
866

867
868 Table 3: Comparison with conventional data augmentation methods and their combinations with IDR
869 on ImageNet-1K using ViT-B.

Method	Top-1 (%)
ViT-B	83.7
PAA (Kang et al., 2025)	84.5
GuidedMixup (Kang et al., 2025)	84.3
IDR	85.6
IDR+PAA	86.1
IDR+GuidedMixup	86.0

875 876 D.3 ABLATION STUDY ON REWEIGHTING REAL AND SYNTHETIC SAMPLES 877

878 We conduct an ablation study to examine the effect of jointly reweighting real and synthetic samples
879 within the proposed IDR framework, in comparison to reweighting only synthetic training samples.
880 To this end, we build an ablation model, IDR (Synthetic Only), in which the reweighting network $g_{\theta}(\cdot)$
881 is applied exclusively to reweight synthetic samples. All experiments are conducted on ImageNet-1K
882 using ViT-B as the classification backbone under identical training configurations. The number of
883 synthetic samples is determined via cross-validation. All baseline methods compared in this study
884 reweight both the real and synthetic training samples. As shown in Table 4, reweighting synthetic
885 samples alone already yields a substantial performance improvement over the baseline methods,
886 highlighting the effectiveness of importance weighting in mitigating noise introduced by synthetic
887 samples. Nevertheless, jointly reweighting both real and synthetic samples results in additional
888 performance gains. For instance, IDR (Synthetic Only) outperforms the best baseline method,
889 GenDataAgent, by 0.7% in top-1 accuracy, while the IDR model further improves performance by
890 0.4%, highlighting the benefits of jointly reweighting real and synthetic samples.

891
892 Table 4: Ablation study on reweighting real and synthetic samples on ImageNet-1K using ViT-B.

Method	Synthetic Data Size ($\times N$)	Top-1
ViT-B/16	-	83.7
MW-Net (Shu et al., 2019)	2	84.1
OTR (Guo et al., 2022)	2	84.0
IE (Chhabra et al., 2024)	3	84.0
CBF (He et al., 2023)	2	84.4
REVAR (Jain et al., 2024)	2	84.2
IDS (Wang et al., 2025)	3	84.4
GenDataAgent (Li et al., 2025)	3	84.5
IDR (Synthetic Only)	5	<u>85.2</u>
IDR	5	85.6

902 D.4 OBJECT DETECTION 903

904 To evaluate the transferability of the features learned by IDR, we assess models pretrained by IDR on
905 ImageNet-1K for the object detection task. We incorporate the ImageNet pre-trained models based on
906 Swin-B in Table 1 into the Cascade Mask R-CNN framework (Cai & Vasconcelos, 2021) for object
907 detection on the MS-COCO dataset (Lin et al., 2014). We follow the Swin training setup (Liu et al.,
908 2021) with standard image resizing, AdamW optimization with learning rate 1×10^{-4} , weight decay
909 0.05, batch size 16, and a 3×36 -epoch schedule, and report COCO box and mask mAP metrics
910 following (Cai & Vasconcelos, 2021). It is observed in Table 5 that IDR consistently outperforms all
911 competing methods across mAP^{box} and mAP^{m} . For instance, IDR outperforms the strongest baseline
912 by 1.0% in mAP^{m} , demonstrating better transferability of learned features to dense prediction tasks.

913 D.5 INSTANCE SEGMENTATION 914

915 In this section, we evaluate the performance of IDR models pre-trained on ImageNet for segmentation
916 on the ADE20K (Zhou et al., 2019) using the UperNet (Xiao et al., 2018) segmentation framework
917 following (Liu et al., 2021). We follow the training and evaluation protocol in (Liu et al., 2021), where
both our model and the baselines are trained on the training split and evaluated on the validation

Table 5: Object Detection Results on COCO.

Methods	mAP ^{box}	mAP ^m
MW-Net (Shu et al., 2019)	52.0	45.2
OTR (Guo et al., 2022)	52.4	45.4
IE (Chhabra et al., 2024)	52.1	45.3
CBF (He et al., 2023)	52.5	45.5
REVAR (Jain et al., 2024)	52.5	45.6
IDS (Wang et al., 2025)	52.7	45.6
GenDataAgent (Li et al., 2025)	52.8	45.8
IDR (Ours)	53.6	46.8

split of the dataset. All models are optimized using AdamW for a total of 160000 iterations with a batch size of 16, an initial learning rate of 6×10^{-5} , and a weight decay of 0.01. The learning rate follows a linear decay schedule after a warm-up phase of 1500 iterations. To enhance generalization, we employ data augmentation techniques including random horizontal flipping, random rescaling with a ratio range of [0.5, 2.0], and random photometric distortions. Stochastic depth regularization is applied with a drop rate of 0.2. For inference, we use multi-scale testing with scale factors varying from 0.5 to 1.75. It is observed in Table 6 that IDR consistently outperforms all competing methods. For instance, IDR outperforms the strongest baseline by 0.7% in mIoU, demonstrating stronger representation transferability for segmentation.

Table 6: Instance Segmentation Results on ADE20K.

Methods	mIoU
MW-Net (Shu et al., 2019)	51.6
OTR (Guo et al., 2022)	51.9
IE (Chhabra et al., 2024)	51.6
CBF (He et al., 2023)	52.1
REVAR (Jain et al., 2024)	52.2
IDS (Wang et al., 2025)	52.1
GenDataAgent (Li et al., 2025)	52.3
IDR (Ours)	53.0

D.6 TRANSFER LEARNING CAPABILITY OF MODELS PRETRAINED BY IDR

We evaluate the transfer learning capability of IDR-pretrained Swin-B models on three standard fine-grained image classification benchmarks: Oxford Flowers-102 (Nilsback & Zisserman, 2008), Oxford-IIIT Pet (Parkhi et al., 2012), and Stanford Cars (Krause et al., 2013). Following established transfer learning protocols (Kolesnikov et al., 2020), the IDR-pretrained Swin-B model is fine-tuned on the training split of each downstream dataset and evaluated on the corresponding test split. Swin-B is also used as the feature backbone for all compared baseline methods. All models are initialized with ImageNet-pretrained weights and fine-tuned for 50 epochs using the Adam optimizer. The learning rate is set to 1×10^{-5} for all datasets. The top-1 classification accuracy on each dataset is reported in Table 7. The results indicate that IDR-pretrained Swin-B models maintain strong transfer learning performance across diverse fine-grained classification tasks. For instance, on the Stanford Cars dataset, IDR achieves a top-1 accuracy of 93.9%, outperforming the strongest baseline, GenDataAgent, by 1.0%.

Table 7: Top-1 classification accuracy comparison for transfer learning on the Oxford Flowers-102, Oxford-IIIT Pet, and Stanford Cars datasets.

Methods	Flowers	Pet	Cars
MW-Net (Shu et al., 2019)	97.0	95.3	91.6
OTR (Guo et al., 2022)	97.1	95.5	91.8
IE (Chhabra et al., 2024)	97.2	95.6	92.0
CBF (He et al., 2023)	97.3	95.8	92.1
REVAR (Jain et al., 2024)	97.6	96.1	92.5
IDS (Wang et al., 2025)	97.5	96.0	92.3
GenDataAgent (Li et al., 2025)	97.6	96.0	92.4
IDR (Ours)	98.4	97.2	93.9

D.7 ABLATION STUDY ON THE DIFFUSION MODELS FOR THE DATA GENERATION IN IDR

To evaluate the impact of the diffusion model used for the data generation in IDR, we compare the performance of IDR using three different diffusion models for GDA, which are DiT-B, DiT-L, and DiT-XL (Peebles & Xie, 2023). ViT-B is used as the classification network for all the models in this ablation study. The data generation time and the classification accuracy on ImageNet-1K are shown in Table 8. It is observed that the performance of the IDR model is not sensitive to the selection of the diffusion models used for data generation. The IDR model based on the largest DiT model DiT-XL only outperforms the IDR based on the smallest DiT model DiT-B by 0.1% in top-1 classification accuracy on ImageNet-1K, demonstrating that the performance of IDR mainly stems from our informative reweighting instead of the diffusion model. In addition, the results in Table 8 show that the synthetic data generation process with the diffusion models in IDR is efficient, with less than 0.1 seconds/image.

Table 8: Performance comparison between IDR models employing different diffusion models for data generation. The data generation time is evaluated on one NVIDIA A100 GPU.

Methods	Accuracy	Generation Time (seconds/image)
ViT-B	83.7	-
IDR (DiT-B)	85.6	0.095
IDR (DiT-L)	85.7	0.151
IDR (DiT-XL)	85.7	0.176

D.8 ABLATION STUDY AND TRAINING TIME ANALYSIS OF THE IDR

To evaluate the effectiveness and efficiency of different components in IDR, we compare the classification performance and the training time of IDR with two ablation models, which are IDR without VUIB and IDR without the reweighting network. The ablation model, IDR without VUIB, optimizes the sample reweighting network by minimizing cross-entropy loss on a set of clean metadata selected from the ImageNet training set, following (Shu et al., 2019). The ablation model, IDR without reweighting, optimizes the classification network by incorporating the VUIB as a regularization term in the training objective without reweighting the real and synthetic samples in the augmented training set. All experiments are conducted using ViT-B as the classification backbone, and training time is measured on four NVIDIA A100 GPUs. We report the training time in seconds per batch with the same batch size in Table 9 because different GDA methods in Table 1 have augmented training sets of different sizes. It is observed in Table 9 that IDR outperforms the strongest baseline, GenDataAgent (Li et al., 2025), by 1.1% in top-1 accuracy while requiring only a marginal 2.5% increase in training time. In addition, it is observed in Table 9 that removing either VUIB or the reweighting network leads to a clear degradation in performance. For instance, removing VUIB decreases the top-1 accuracy of IDR by 1.3%, which demonstrates the critical role of VUIB in learning robust and discriminative representations from the augmented training set. It is worthwhile to mention that the class centroids in the input feature space can be efficiently pre-computed prior to training using FAISS (Douze et al., 2025), which supports fast, memory-efficient, and highly parallelized centroid and distance computations on GPUs, even for high-dimensional features. In particular, computing the class centroids and the distances between the input features of all samples in the augmented training set and their corresponding class centroids requires only 15 minutes, which happens one time before the training. Computing the class centroids and the distances between the learned features of all samples in the augmented training set and their corresponding class centroids requires only 6.1 seconds in each epoch. As a result, it is not necessary to project the input features into a lower-dimensional space using methods such as VAEs (Kingma & Welling, 2014), which require additional training time.

D.9 COMPARISON WITH ACTIVE LEARNING METHODS

Active learning (AL) methods aim to minimize the effort required for labeling training data by strategically choosing the most informative instances for annotation (Sinha et al., 2019; Yoo & Kweon, 2019; Gao et al., 2020; Kushnir & Venturi, 2023; Yang et al., 2023; Chhabra et al., 2024). The selection of the data for annotation by active learning methods is usually achieved by identifying

Table 9: Ablation study of IDR with training time analysis. The training time is evaluated on four NVIDIA A100 GPUs. The setup of all the methods in this table is the same as that in Table 1.

Methods	Top-1	Training Time (Seconds / Batch)
ViT-B	83.7	0.739
MW-Net (Shu et al., 2019)	84.1	0.801
OTR (Guo et al., 2022)	84.0	0.772
IE (Chhabra et al., 2024)	84.0	0.739
CBF (He et al., 2023)	84.4	0.714
REVAR (Jain et al., 2024)	84.2	0.825
IDS (Wang et al., 2025)	84.4	0.825
GenDataAgent (Li et al., 2025)	84.5	0.805
IDR w/o VUIB	84.3	0.753
IDR w/o Reweighting	85.2	0.801
IDR (Ours)	85.6	0.825

the most informative data points. Such a process works similarly to the data reweighting process in IDR for identifying the most informative synthetic samples. To show the advantage of IDR over active learning methods in selecting the most informative synthetic samples, we compare IDR with state-of-the-art active learning methods, including CAMPAL (Yang et al., 2023), SAAL (Chhabra et al., 2024), and CSQ (Gwon et al., 2025). The active learning methods are adopted to select training data in each epoch from the augmented training set. The results are shown in Table 10. It is observed that IDR significantly outperforms the competing active learning methods. For instance, IDR outperforms the best active learning method, CSQ, by 1.2% in top-1 classification accuracy on ImageNet-1K, demonstrating the superiority of IDR in selecting informative training samples compared to active learning methods.

Table 10: Comparison between IDR and active learning methods.

Methods	Top-1
ViT-B	83.7
CAMPAL (Yang et al., 2023)	84.1
SAAL (Chhabra et al., 2024)	84.2
CSQ (Gwon et al., 2025)	84.4
IDR (Ours)	85.6

D.10 IMPROVEMENT SIGNIFICANCE ANALYSIS.

To verify that the improvement of our proposed IDR on existing methods is statistically significant and outside the range of random error, we train both IDR and the best baseline methods on different datasets for 10 times with different seeds for random initialization of the networks and train/val/test splits. Next, we perform the t-test between the results of IDR and the results of the best baseline methods with different backbones to assess if the improvement of IDR is statistically significant. The mean and standard deviation of the results and the p-values of the t-test are shown in Table 11. It is observed that the largest p-value is 2.6×10^{-6} , which is less than 0.05. The t-test results suggest that the improvement of IDR over the baseline methods is statistically significant with $p \ll 0.05$, and it is not caused by random error.

Table 11: P-values of t-test between IDR and the baseline methods with the best performance.

Backbone	Best Baseline	IDR	p-value
ViT-S/16	82.1±0.25	83.2±0.18	2.6×10^{-6}
ViT-B/16	84.5±0.21	85.6±0.22	5.4×10^{-6}
Swin-T	82.3±0.29	83.4±0.27	1.3×10^{-9}
Swin-B	84.6±0.27	85.7±0.19	7.2×10^{-8}

1080 D.11 T-SNE VISUALIZATION ANALYSIS ON REAL SAMPLES IN THE AUGMENTED TRAINING
 1081 SET

1083 E PROOF OF THEOREM 3.1

1085 We define

$$1086 \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) = \frac{\exp\left(-\|\tilde{z}_i - c_a^{(\text{feat})}\|_2^2 - \|\tilde{x}_i - c_y^{(\text{input})}\|_2^2\right)}{\sum_{a'=1}^C \sum_{y'=1}^C \exp\left(-\|\tilde{z}_i - c_{a'}^{(\text{feat})}\|_2^2 - \|\tilde{x}_i - c_{y'}^{(\text{input})}\|_2^2\right)}, \quad a, y \in [C]. \quad (12)$$

1093 **Lemma E.1.** Let $\Pr[\tilde{X} \in y] = \sum_{i=1}^{N'} \mathbb{1}_{\{\tilde{y}_i=y\}}/N' := p_y$ for every $y \in [C]$, then

$$1094 I(\tilde{Z}, \tilde{X}) \leq \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) \log \left(\frac{\phi(\tilde{z}_i \in a, \tilde{x}_i \in y)}{p_y \phi(\tilde{z}_i, c_a^{(\text{feat})})} \right). \quad (13)$$

1102 *Proof.* Then the joint probability $\Pr[\tilde{Z} \in a, \tilde{X} \in y]$ is specified by

$$1103 \Pr[\tilde{Z} \in a, \tilde{X} \in y] = \frac{1}{N'} \sum_{i=1}^{N'} \phi(\tilde{z}_i \in a, \tilde{x}_i \in y). \quad (14)$$

1109 It then follows from (14) and the log sum inequality that

$$1110 \begin{aligned} 1111 I(\tilde{Z}, \tilde{X}) &= \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{X} \in y] \log \frac{\Pr[\tilde{Z} \in a, \tilde{X} \in y]}{\Pr[\tilde{Z} \in a] \Pr[\tilde{X} \in y]} \\ 1112 &\leq \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) \log(\phi(\tilde{z}_i \in a, \tilde{x}_i \in y)) \\ 1113 &\quad - \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) \log(\phi(\tilde{z}_i, c_a^{(\text{feat})}) p_y) \\ 1114 &= \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) \log \left(\frac{\phi(\tilde{z}_i \in a, \tilde{x}_i \in y)}{p_y \phi(\tilde{z}_i, c_a^{(\text{feat})})} \right). \end{aligned} \quad (15)$$

1126 □

1129 **Lemma E.2.**

$$1130 I(\tilde{Z}, \tilde{Y}) \geq \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i=y\}} \log Q(\tilde{Z} \in a | \tilde{Y} = y) \quad (16)$$

1134 *Proof.* Let $Q(\tilde{Z}|\tilde{Y})$ be a variational distribution. We have

$$\begin{aligned}
1135 & \\
1136 & I(\tilde{Z}, \tilde{Y}) = \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{\Pr[\tilde{Z} \in a, \tilde{Y} = y]}{\Pr[\tilde{Z} \in a] \Pr[\tilde{Y} = y]} \\
1137 & \\
1138 & \\
1139 & = \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{\Pr[\tilde{Z} \in a|\tilde{Y} = y] Q(\tilde{Z} \in a|\tilde{Y} = y)}{\Pr[\tilde{Z} \in a] Q(\tilde{Z} \in a|\tilde{Y} = y)} \\
1140 & \\
1141 & \geq \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{\Pr[\tilde{Z} \in a|\tilde{Y} = y]}{Q(\tilde{Z} \in a|\tilde{Y} = y)} + \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{Q(\tilde{Z} \in a|\tilde{Y} = y)}{\Pr[\tilde{Z} \in a]} \\
1142 & \\
1143 & = \text{KL}\left(P(\tilde{Z}|\tilde{Y}) \parallel Q(\tilde{Z}|\tilde{Y})\right) + \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{Q(\tilde{Z} \in a|\tilde{Y} = y)}{\Pr[\tilde{Z} \in a]} \\
1144 & \\
1145 & \geq \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log \frac{Q(\tilde{Z} \in a|\tilde{Y} = y)}{\Pr[\tilde{Z} \in a]} \\
1146 & \\
1147 & = \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log Q(\tilde{Z} \in a|\tilde{Y} = y) + H\left(P(\tilde{Z})\right) \\
1148 & \\
1149 & \geq \sum_{a=1}^C \sum_{y=1}^C \Pr[\tilde{Z} \in a, \tilde{Y} = y] \log Q(\tilde{Z} \in a|\tilde{Y} = y) \\
1150 & \\
1151 & \geq \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i=y\}} \log Q(\tilde{Z} \in a|\tilde{Y} = y), \tag{17} \\
1152 & \\
1153 & \\
1154 & \\
1155 & \\
1156 & \\
1157 & \\
1158 & \\
1159 & \\
1160 & \\
1161 &
\end{aligned}$$

1162 where the last inequality follows from the joint probability $\Pr[\tilde{Z} \in a, \tilde{Y} = y] =$
1163 $\frac{1}{N'} \sum_{i=1}^{N'} \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i=y\}}$. \square

1166 *Proof of Theorem 3.1.* Equation 1 in Theorem 3.1 of the main paper follows from
1167 $\text{IB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}) = I(\tilde{Z}, \tilde{X}) - I(\tilde{Z}, \tilde{Y})$, the upper bound for $I(\tilde{Z}, \tilde{X})$ in Lemma E.1
1168 and the lower bound for $I(\tilde{Z}, \tilde{Y})$ in Lemma E.2. \square

1170 E.1 COMPUTATION OF $Q^{(t)}(\tilde{Z}|\tilde{Y})$

1171 The variational distribution $Q^{(t)}(\tilde{Z}|\tilde{Y})$ can be computed by

$$\begin{aligned}
1172 & \\
1173 & \\
1174 & Q^{(t)}(\tilde{Z} \in a|\tilde{Y} = y) = \Pr[\tilde{Z} \in a|\tilde{Y} = y] = \frac{\sum_{i=1}^{N'} \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i=y\}}}{\sum_{i=1}^{N'} \mathbb{1}_{\{\tilde{y}_i=y\}}}. \tag{18} \\
1175 & \\
1176 & \\
1177 &
\end{aligned}$$

1178 Algorithm 5 computes $Q^{(t)}(\tilde{Z}|\tilde{Y})$ efficiently with a time complexity of $\Theta(N'C + N'T_0)$, where C
1179 is the number of classes, N' is the number of training samples, and T_0 denotes the computational
1180 complexity of a forward and backward pass of the neural network with respect to each training
1181 sample. $Q^{(t)}(\tilde{Z}|\tilde{Y})$ is the variational distribution $Q(\tilde{Z}|\tilde{Y})$ computed at the t -th epoch in Algorithm 1
1182 of the main paper.

1184 E.2 ISSUE ABOUT THE UPPER BOUND FOR THE IB LOSS IN IDS (WANG ET AL., 2025) AND 1185 DCS-TRANSFORMER (WANG ET AL., 2024)

1186 The derivation of the upper bound for the IB loss in IDS (Wang et al., 2025) and DCS-
1187 Transformer (Wang et al., 2024) follows a similar recipe, with an upper bound for the mutual

information between the learned features \widehat{Z} and the input features \widehat{X} , $I(\widehat{Z}, \widehat{X})$, and a lower bound for the mutual information between the learned features \widehat{Z} and the class label \widehat{Y} , $I(\widehat{Z}, \widehat{Y})$. However, the upper bounds for $I(\widehat{Z}, \widehat{X})$ in both IDS (Wang et al., 2025) and DCS-Transformer (Wang et al., 2024) is

$$\begin{aligned} I(\widehat{Z}, \widehat{X}) &\leq \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\widehat{z}_i, c_a^{(\text{feat})}) \phi(\widehat{x}_i, c_y^{(\text{input})}) \log \phi(\widehat{x}_i, c_y^{(\text{input})}) \\ &= \frac{1}{N'} \sum_{i=1}^{N'} \sum_{y=1}^C \phi(\widehat{x}_i, c_y^{(\text{input})}) \log \phi(\widehat{x}_i, c_y^{(\text{input})}), \end{aligned} \quad (19)$$

where the last equality follows from $\sum_{a=1}^C \phi(\widehat{z}_i, c_a^{(\text{feat})}) = 1$. In (19), $\widehat{x}_i, \widehat{z}_i$ stand for the input training feature and the learned feature, respectively, and $c_a^{(\text{feat})}, c_y^{(\text{input})}$ stand for the class centroids of the input training features and the learned features. Compared to our derived upper bound for $I(\widehat{Z}, \widehat{X})$ on the RHS of (13) of Lemma E.1, the upper bound for the RHS of (19) does not depend on the learned features $\{\widehat{z}_i\}$ and it only depends on the input training data. As a result, the optimization of the upper bound for $I(\widehat{Z}, \widehat{X})$ in the IB loss in (19) cannot optimize the weights of the neural network for reducing the IB loss for data reweighting, making the upper bound for $I(\widehat{Z}, \widehat{X})$ vacuous. In a strong contrast, our upper bound (13) depends on both the real training data and the learned features $\{\widehat{z}_i\}$, leading to improved performance of our IDR over IDS in Tables 1, 5, 6, and 7 of Sections 4.1, D.4, D.5, and D.6.

Algorithm 4 MI Minimization with vCLUB (Algorithm 1 in CLUB (Cheng et al., 2020))

- 1: **for** each training iteration **do**
 - 2: Sample $\{(x_i, z_i)\}_{i=1}^{N'}$ from $p_\sigma(x, z)$
 - 3: Compute log-likelihood $\mathcal{L}(\theta) = \frac{1}{N'} \sum_{i=1}^{N'} \log q_\theta(z_i|x_i)$
 - 4: Update $q_\theta(z|x)$ by maximizing $\mathcal{L}(\theta)$
 - 5: **for** $i = 1$ to N' **do**
 - 6: $L_i = \log q_\theta(z_i|x_i) - \frac{1}{N'} \sum_{j=1}^{N'} \log q_\theta(z_j|x_i)$
 - 7: **end for**
 - 8: Update $p_\sigma(x, z)$ by minimizing $\widehat{I}_{\text{vCLUB}} = \frac{1}{N'} \sum_{i=1}^{N'} L_i$
 - 9: **end for**
-

Algorithm 5 Efficient Computation of $Q^{(t)} (\widehat{Z} \in a \mid \widehat{Y} = y)$

Input: The precomputed soft assignments $\phi(\widehat{z}_i, c_a^{(\text{feat})})$ for $i \in [N']$ and $a \in [C]$, the class labels $\{\widehat{y}_i\}_{i=1}^{N'}$, and the number of classes C . Here $\{\widehat{z}_i\}_{i=1}^{N'}$ are computed at the t -th epoch in Algorithm 1 of the main paper.

Output: Conditional distribution matrix $Q \in \mathbb{R}^{C \times C}$

- 1: Initialize $Q \leftarrow \mathbf{0}^{C \times C}$ and count vector $M \leftarrow \mathbf{0}^C$.
 - 2: **for** $i = 1$ to N' **do**
 - 3: **for** $a = 1$ to C **do**
 - 4: $Q[a, \widehat{y}_i] \leftarrow Q[a, \widehat{y}_i] + \phi(\widehat{z}_i, c_a^{(\text{feat})})$
 - 5: **end for**
 - 6: $M[\widehat{y}_i] \leftarrow M[\widehat{y}_i] + 1$
 - 7: **end for**
 - 8: **for** $y = 1$ to C **do**
 - 9: **for** $a = 1$ to C **do**
 - 10: $Q[a, y] \leftarrow Q[a, y]/M[y]$
 - 11: **end for**
 - 12: **end for**
 - 13: **return** Q
-

F COMPUTATIONAL COMPLEXITY ANALYSIS

Given samples $\{(x_i, z_i)\}_{i=1}^{N'}$ drawn from the joint distribution $p_\sigma(x, z) = p_\sigma(z|x)p(x)$, the goal of CLUB is to optimize the parameters of the predictive neural network $p_\sigma(z|x)$ such that the resulting joint model $p_\sigma(x, z)$ induces minimal mutual information between x and z . Here, $p_\sigma(z|x)$ serves as the main predictive model, such as an encoding model or a classification model, while $p(x)$ denotes the empirical data distribution. In addition, $q_\theta(z|x)$ is a variational conditional distribution, also implemented as a neural network, used to approximate $p(z|x)$ and to provide a differentiable estimator of mutual information. The optimization alternates between updating $q_\theta(z|x)$ by maximizing the log-likelihood $\mathcal{L}(\theta) = \frac{1}{N'} \sum_{i=1}^{N'} \log q_\theta(z_i|x_i)$ to improve the approximation accuracy, and updating $p_\sigma(x, z)$ by reducing the upper bound for the mutual information, $I(x, z)$, $\widehat{I}_{\text{vCLUB}} = \frac{1}{N'} \sum_{i=1}^{N'} L_i$, where $L_i = \log q_\theta(z_i|x_i) - \frac{1}{N'} \sum_{j=1}^{N'} \log q_\theta(z_j|x_i)$. Through this alternating process, the algorithm jointly learns an accurate variational estimator q_θ and a predictive model p_σ . The complete training procedure is summarized in Algorithm 4.

Suppose the computation of $q_\theta(z|x)$ costs T_q and that of $p_\sigma(z|x)$ costs T_p . Computing $\mathcal{L}(\theta) = \frac{1}{N'} \sum_{i=1}^{N'} \log q_\theta(z_i|x_i)$ requires $\Theta(N'T_q)$ time. The computation of vCLUB $\widehat{I}_{\text{vCLUB}} = \frac{1}{N'} \sum_{i=1}^{N'} \left[\log q_\theta(z_i|x_i) - \frac{1}{N'} \sum_{j=1}^{N'} \log q_\theta(z_j|x_i) \right]$ requires computing $[\log q_\theta(z_j|x_i)]$ for $i, j \in [N']$, which requires $\Theta(N'^2T_q)$ time. Updating p_σ for reducing $\widehat{I}_{\text{vCLUB}}$ over the same N' inputs adds $N'T_p$ time. Hence, one training epoch requires $\Theta(N'^2T_q + N'T_q + N'T_p) = \Theta(N'^2T_q + N'T_p)$.

Computational Complexity of VUIB. Herein we analyze the computational complexity for computing the upper bound, VUIB, for the IB loss. Let T_0 denote the complexity of a forward and backward computation of the model predicting \tilde{z}_i for one sample. For each epoch, the computation of $\{\tilde{z}_i\}_{i=1}^{N'}$ using the classification network takes $\Theta(N'T_0)$ time. Once $\{\tilde{z}_i\}_{i=1}^{N'}$ are computed, we pre-compute $\phi(\tilde{z}_i, c_a^{(\text{feat})})$ for $i \in [N']$ and $a \in [C]$, which takes $\Theta(N'C)$ time.

To compute $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}})$, we separately compute $\frac{1}{N'} \sum_{i=1}^{N'} U_i$, which is the upper bound for $I(\tilde{Z}, \tilde{X})$, and $\frac{1}{N'} \sum_{i=1}^{N'} V_i$, which is the lower bound for $I(\tilde{Z}, \tilde{Y})$. Let $A_{ia} := \exp(-\|\tilde{z}_i - c_a^{(\text{feat})}\|_2^2)$ and $B_{iy} := \exp(-\|\tilde{x}_i - c_y^{(\text{input})}\|_2^2)$ for all $a \in [C]$ and $y \in [C]$. Let $D_i := \sum_{a'=1}^C \sum_{y'=1}^C A_{ia'} B_{iy'} = \widehat{A}_i \widehat{B}_i$, where $\widehat{A}_i := (\sum_{a'=1}^C A_{ia'})$, $\widehat{B}_i := (\sum_{y'=1}^C B_{iy'})$. For each $i \in [N']$, we first pre-compute all A_{ia} and B_{iy} for $a \in [C]$ and $y \in [C]$, which takes $\Theta(C)$ time. For each $i \in [N']$, we then pre-compute \widehat{A}_i , \widehat{B}_i , and D_i , which takes $\Theta(C)$ time. Then U_i can be computed by

$$\begin{aligned} U_i &= \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i \in a, \tilde{x}_i \in y) \log \left(\frac{\phi(\tilde{z}_i \in a, \tilde{x}_i \in y)}{p_y \phi(\tilde{z}_i, c_a^{(\text{feat})})} \right) \\ &= \frac{1}{D_i} \sum_{a=1}^C \sum_{y=1}^C A_{ia} B_{iy} \log \left(\frac{A_{ia} B_{iy} / D_i}{p_y \phi(\tilde{z}_i, c_a^{(\text{feat})})} \right) \\ &= -\log D_i + \frac{1}{D_i} \sum_{a=1}^C \sum_{y=1}^C A_{ia} B_{iy} (\log A_{ia} + \log B_{iy} - \log p_y - \log \phi(\tilde{z}_i, c_a^{(\text{feat})})) \\ &= -\log D_i - \frac{\widehat{A}_i}{D_i} \sum_{y=1}^C B_{iy} \log p_y + \frac{1}{D_i} \left(\sum_{a=1}^C A_{ia} \log A_{ia} \right) \widehat{B}_i \\ &\quad + \frac{\widehat{A}_i}{D_i} \left(\sum_{y=1}^C B_{iy} \log B_{iy} \right) - \frac{\widehat{B}_i}{D_i} \left(\sum_{a=1}^C A_{ia} \log \phi(\tilde{z}_i, c_a^{(\text{feat})}) \right). \end{aligned}$$

For each $i \in [N']$, the computation of $\frac{\widehat{A}_i}{D_i} \sum_{y=1}^C B_{iy} \log p_y$, $\frac{1}{D_i} \left(\sum_{a=1}^C A_{ia} \log A_{ia} \right) \widehat{B}_i$, $\frac{\widehat{A}_i}{D_i} \left(\sum_{y=1}^C B_{iy} \log B_{iy} \right)$, and $\frac{\widehat{B}_i}{D_i} \left(\sum_{a=1}^C A_{ia} \log \phi(\tilde{z}_i, c_a^{(\text{feat})}) \right)$ each takes $\Theta(C)$ time. As a result, the computation of U_i for each $i \in [N']$ takes $\Theta(C + C + 1 + C + C + C + C) = \Theta(C)$ time.

1296 Therefore, the computation of $\frac{1}{N'} \sum_{i=1}^{N'} U_i$, which is the upper bound for $I(\tilde{Z}, \tilde{X})$, takes $\Theta(N'C)$
 1297 time.

1298 Given the soft assignment matrix $\phi(\tilde{z}_i, c_a^{(\text{feat})})$ for all $i \in [N']$ and $a \in [C]$, the conditional distribution
 1299 matrix $Q(\tilde{Z} \in a \mid \tilde{Y} = y) \in \mathbb{R}^{C \times C}$ can be efficiently computed following Algorithm 5. We denote
 1300 by $Q[a, y] = Q(\tilde{Z} \in a \mid \tilde{Y} = y)$ the (a, y) -th entry of Q , representing the conditional probability
 1301 that a learned feature \tilde{Z} belongs to class a given class label y . Each entry $Q[a, y]$ is computed by
 1302 aggregating the soft assignment values $\phi(\tilde{z}_i, c_a^{(\text{feat})})$ over all $i \in [N']$ such that $\tilde{y}_i = y$, followed
 1303 by normalization with respect to the total number of samples in that class. The accumulation step
 1304 (lines 2–7 in Algorithm 5) requires $\Theta(N'C)$ time, while the normalization step (lines 8–12 in
 1305 Algorithm 5) requires $\Theta(C^2)$ time. Since $N' \gg C$, the computational complexity of computing Q is
 1306 $\Theta(N'C + C^2) = \Theta(N'C)$.

1307 Let $\mathcal{I}_y := \{i \in [N'] \mid \tilde{y}_i = y\}$ be the index set of the samples from the class y for $y \in [C]$ in the
 1308 augmented training set. The lower bound for the mutual information $I(\tilde{Z}, \tilde{Y})$ can be computed by

$$\begin{aligned}
 1311 \quad \frac{1}{N'} \sum_{i=1}^{N'} V_i &= \frac{1}{N'} \sum_{i=1}^{N'} \sum_{a=1}^C \sum_{y=1}^C \phi(\tilde{z}_i, c_a^{(\text{feat})}) \mathbb{1}_{\{\tilde{y}_i=y\}} \log Q(\tilde{Z} \in a \mid \tilde{Y} = y) \\
 1312 &= \frac{1}{N'} \sum_{a=1}^C \sum_{y=1}^C \sum_{i \in \mathcal{I}_y} \phi(\tilde{z}_i, c_a^{(\text{feat})}) \log Q(\tilde{Z} \in a \mid \tilde{Y} = y). \\
 1313 & \\
 1314 & \\
 1315 & \\
 1316 &
 \end{aligned}$$

1317 For each $a \in [C]$, the computation of $\sum_{y=1}^C \sum_{i \in \mathcal{I}_y} \phi(\tilde{z}_i, c_a^{(\text{feat})}) \log Q(\tilde{Z} \in a \mid \tilde{Y} = y)$ takes
 1318 $\sum_{y=1}^C |\mathcal{I}_y| = \Theta(N')$ time. As a result, the computation of $\frac{1}{N'} \sum_{i=1}^{N'} V_i$, which is the lower bound
 1319 for $I(\tilde{Z}, \tilde{Y})$, takes $\Theta(N'C)$ time.

1320 Since $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}}) = \frac{1}{N'} \sum_{i=1}^{N'} U_i - \frac{1}{N'} \sum_{i=1}^{N'} V_i$, the overall computation cost for
 1321 $\text{VUIB}(\mathcal{C}(\theta, \Theta), \Theta, \mathcal{D}_{\text{aug}})$ is $\Theta(N'T_0 + N'C + N'C + N'C + N'C) = \Theta(N'T_0 + N'C)$.

1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349