

# Supplementary Materials of Semantic-aware Data Augmentation for Text-to-image Synthesis

Anonymous Author(s)  
Submission Id: 890

## ACM Reference Format:

Anonymous Author(s). 2023. Supplementary Materials of Semantic-aware Data Augmentation for Text-to-image Synthesis. In *Proceedings of ACM International Conference on Multimedia (ACM Multimedia '2023)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## A MORE MATHEMATICAL DETAILS

Here we provide more details for our derivations and proofs.

### A.1 Derivation Details of Training Objectives for $G$ with $ITA$

Based on empirical risk minimization (ERM), the empirical risk for generator  $G$  is defined as:

$$R_k(\theta) := \frac{1}{k} \sum_{i=1}^k L(\theta, X_i). \quad (1)$$

Its standard augmented version and corresponding augmented loss are defined as:

$$\hat{R}_k(\theta) := \frac{1}{k} \sum_{i=1}^k \int_{\mathcal{A}} L(\theta, f(X)) dQ_{ITA}(f), \quad (2)$$

where  $Q_{ITA}$  is a probability distribution on a group  $\mathcal{A}$  of  $ITA$  transforms from which  $f$  is sampled. Since only one  $ITA$  will be used, the general sample objective with  $ITA$  is defined as:

$$\min_{\theta} \hat{R}_k(\theta) := \frac{1}{k} \sum_{i=1}^k L(\theta, ITA(X_i)). \quad (3)$$

We then define the solution of Eq. (3) as:

$$\theta_{ITA}^* \in \arg \min_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k L(\theta, ITA(X_i)), \quad (4)$$

where  $\Theta$  is defined as some parameter space.

### A.2 Proof Details

**PROPOSITION A.1 ( $ITA$  INCREASES T2ISYN SEMANTIC CONSISTENCY).**  
Assume exact invariance holds. Consider an unaugmented text-image

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ACM Multimedia '2023, October 28– November 03, 2023, Ottawa, Ontario, Canada

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

generator  $\hat{\theta}(X)$  of  $G$  and its augmented version  $\hat{\theta}_{ITA}$ . For any real-valued convex loss  $S(\theta, \cdot)$  that measures the semantic consistency, we have:

$$\mathbb{E}[S(\theta, \hat{\theta}(X))] \geq \mathbb{E}[S(\theta, \hat{\theta}_{ITA}(X))], \quad (5)$$

which means with  $ITA$ , a model can have lower  $\mathbb{E}[S(\theta, \hat{\theta}_{ITA}(X))]$  then a better text-image consistency.

**PROOF.** From Group-Theoretic Framework for Data Augmentation [2], we obtain a direct consequence that:

$$\text{Cov}[S(\hat{\theta}_{ITA}(X))] \leq \text{Cov}[S(\hat{\theta}(X))], \quad (6)$$

where  $\text{Cov}[\cdot]$  means the covariance matrix decreases in the Loewner order. Therefore for any real-valued convex loss function  $S(\theta, \cdot)$ , we have:

$$\mathbb{E}_{X \sim \mathcal{X}}[S(G(X))] \geq \mathbb{E}_{X \sim \mathcal{X}}[S(\bar{G}(X))], \quad (7)$$

where  $\bar{G}$  is the conditional expectation of  $G$  and  $S(\theta, \cdot)$  is any real-valued convex loss function. Empirically,  $S(\theta, \cdot)$  can be a real-valued convex loss produced by discriminators, perceptual semantic loss produced by pretrained models, and others. It also suggests that  $ITA$  can be considered as an algorithmic regularization like other data augmentations, and augmented  $G$  can obtain better text-image consistency.  $\square$

**PROPOSITION A.2.**  $ITA_C$  can be considered a closed-form solution for general textual semantic preserving augmentation methods of T2Isyn.

**PROOF.** Assume exact invariance holds. Captions offered in the dataset are based on real images, thus:  $e_s|_r \triangleq e_s$ . Assume all mentioned following models are well-trained. We consider two situations:

- (1) For methods that use extra models to generate more textual data based on real images  $r$  (such as I2T2I [5] which uses a pretrained captioning model), we have:

$$e_s|_r \sim \mathcal{N}(m_r, \mathbb{C}_{ss|r}\mathbb{I}) = Q_s|_r, \quad (8)$$

$$e'_s|_r \sim \mathcal{N}(m_r, \mathbb{C}'_{ss|r}\mathbb{I}) = Q'_s|_r. \quad (9)$$

When the extra models are trained on the dataset used for T2Isyn, exact invariance holds. We have:

$$Q_s|_r =_d Q'_s|_r, \quad (10)$$

$$e'_s|_r \sim \mathcal{N}(m_r, \mathbb{C}_{ss|r}\mathbb{I}). \quad (11)$$

- (2) Consider methods that use extra models that generate synonymous texts based on real texts (such as retrieving texts from the dataset and refining the conflicts like RiFeGan [3], using extra pretrained synonymous text generating model, and our proposed  $ITA_T$ ). Assume exact invariance holds. Captions offered in the dataset are based on real images,

thus:  $e_{s|r} \triangleq e_s, e_{s|r} \sim Q_{ss|r}$ . Augmented texts  $e'_s$  are retrieved from the dataset and refine the semantic conflicts between  $e'_{s|r}$  and  $e_{s|r}$  based on the main semantics of real images  $r$ . Therefore:

$$e_{s|r} \sim N(m_r, \mathbb{C}_{ss|r}\mathbb{I}) = Q_{s|r}, \quad (12)$$

$$e'_{s|r} \sim N(m_r, \mathbb{C}_{ss|r}\mathbb{I}) = Q_{s|r}. \quad (13)$$

Due to  $e_{s|r}$  depending on the semantics of  $r$ ,  $e_{s|r}$  should maintain the main semantics of  $r$ . Therefore we have:

$$m_{s|r} \approx m_r, \quad (14)$$

$$e'_{s|r} \sim N(m_{s|r}, \mathbb{C}_{ss|r}\mathbb{I}), \quad (15)$$

where  $ITA_C$  is a closed-form solution. Therefore,  $ITA_C$  can be considered a closed-form solution for general textual semantic preserving augmentation methods of T2Isyn.  $\square$

**PROPOSITION A.3.** *Assume that  $E_I$  is linear. Constraining the distribution  $Q_E$  of  $e_{f|s}$  can additionally constrain the distribution  $\mathcal{F}$  of  $f$ .*

**PROOF.** There are two situations:

- (1) If  $E_I$  is inevitable, Proposition A.3 is obvious.
- (2) If  $E_I$  is not inevitable, constraining  $\mathcal{F}$  can affect  $\mathcal{E}$  in not Nullspace of  $E_I \setminus \text{Null}(E_I)$ :

$$C(\mathcal{E}) \propto C(-\text{Null}(E_I)(\mathcal{F})), \quad (16)$$

where  $C(\cdot)$  is a certain constraint. For Nullspace, there will be no effect. If not all the mass of  $\mathcal{F}$  locates in the  $\text{Null}(E_I)$ , Proposition A.3 holds. If  $\mathcal{F}$  all locates in the  $\text{Null}(E_I)$  while  $E_I$  is well trained, it means  $\mathcal{F}$  does not contain any semantics that matches textual semantics, inferring a total collapse of  $G$ . Since we assume the  $G$  can learn the representation, it is impossible that  $\mathcal{F}$  all locates in the  $\text{Null}(E_I)$ .

Therefore, Proposition A.3 holds.  $\square$

**PROPOSITION A.4.**  *$L_r$  leads to  $|e'_{f|s} - e_{f|s}|$  is less than or equal to a sequence  $\Lambda$  of positive constants, further constrains the semantic manifold of generated embeddings to meet the Lipschitz condition.*

**PROOF.** From Eq. (??), we have following constrain for  $e'_{f|s}$  and  $e_{f|s}$ :

$$|e'_{f|s} - e_{f|s}| \leq |e''_{f|s} - e_{f|s}| = |\epsilon^*| \odot \beta \cdot d(\mathbb{C}_{rr|s}). \quad (17)$$

For each dimension of semantic embeddings, we have:

$$\begin{aligned} |e'_{f|s}^d - e_{f|s}^d| &= \beta \cdot \mathbb{E}[(e'_{s|r}^d - e_{s|r}^d)^2], \\ &\leq \beta \cdot \max[(e'_{s|r}^d - e_{s|r}^d)^2] \\ &= \beta \cdot [(e''_{s|r}^d - e_{s|r}^d)^2], \\ &= |\epsilon^{*d}| \cdot \beta \cdot d(\mathbb{C}_{rr|s})^d \\ &= \beta \cdot d(\mathbb{C}_{rr|s})^d, \end{aligned} \quad (18)$$

$$|e'_{f|s}^d - e_{f|s}^d| \leq \beta \cdot d(\mathbb{C}_{rr|s})^d, \quad (19)$$

where  $d = \{1, \dots, n\}$  and  $n$  is the dimension of the semantic embedding;  $d(\cdot)$  represents diagonal part of a matrix;  $\beta$  is a positive constant. Due to the fact of the many-to-many relationship between

texts and images, we have  $d(\mathbb{C}_{rr|s})^d > 0$ . Assume exact invariance holds,  $|\epsilon^{*d}| = 1; \beta \cdot d(\mathbb{C}_{rr|s})^d > 0$  is a constant. Thus:

$$|e'_{f|s} - e_{f|s}| \leq \Lambda. \quad (20)$$

If we use  $e''_{s|r}$  to generate images, we can alter Eq. (20) to:

$$|e''_{f|s} - e_{f|s}| = \Lambda. \quad (21)$$

Similar to Eq. 21, we can have:

$$|e''_{s|r} - e_{s|r}| = \Lambda_s, \quad (22)$$

where  $\Lambda_s$  is also a sequence of positive constants. Then we have:

$$\frac{|e''_{f|s} - e_{f|s}|}{|e''_{s|r} - e_{s|r}|} = \frac{\Lambda}{\Lambda_s} = M. \quad (23)$$

Due to the findings that semantic features in deep feature space are usually linearized [1, 11, 13], we assume semantic features for texts and images are linearized. Following Eq. (??), we can further have that:

$$|\delta| \frac{|e'_{f|s} - e_{f|s}|}{|e'_{s|r} - e_{s|r}|} = \frac{|e''_{f|s} - e_{f|s}|}{|e''_{s|r} - e_{s|r}|} = \frac{M}{|\delta|} \leq K, \text{ s.t. } e'_{s|r} \neq e_{s|r}, \quad (24)$$

where  $\delta$  is a non-zero coefficient. Finally,  $e'_{f|s} = E_I(G(e'_{s|r}))$ ,  $e_{f|s} = E_I(G(e_{s|r}))$  where  $E_I$  is the image encoder, we have:

$$\frac{|E_I(G(e'_{s|r})) - E_I(G(e_{s|r}))|}{|e'_{s|r} - e_{s|r}|} \leq K, \text{ s.t. } e'_{s|r} \neq e_{s|r}, \quad (25)$$

where it meets Lipschitz condition.  $\square$

**PROPOSITION A.5.**  *$L_r$  provides tighter image semantic constraints than  $L_{db}$  [6] which is defined as:*

$$L_{db} = 1 - \frac{(e'_{s|r} - e_{s|r}) \cdot (e'_{f|s} - e_{f|s})}{\|(e'_{s|r} - e_{s|r})\|^2 \cdot \|(e'_{f|s} - e_{f|s})\|^2}, \quad (26)$$

**PROOF.** For Eq. (??), assume  $L_{db} = 0$  and use  $\epsilon^*$ , combining with Eq. (??), we have:

$$\frac{(e''_{f|s} - e_{f|s})}{\|e''_{f|s} - e_{f|s}\|^2} = \frac{\|e''_{s|r} - e_{s|r}\|^2}{(e''_{s|r} - e_{s|r})} \quad (27)$$

$$= \frac{\|\beta \cdot \epsilon^* \odot d(\mathbb{C}_{ss|r})\|^2}{\beta \cdot \epsilon^* \odot d(\mathbb{C}_{ss|r})}. \quad (28)$$

Therefore:

$$|e''_{f|s} - e_{f|s}| = \|e''_{f|s} - e_{f|s}\|^2 \cdot \frac{\|\beta \cdot \epsilon^* \odot d(\mathbb{C}_{ss|r})\|^2}{\|\beta \cdot \epsilon^* \odot d(\mathbb{C}_{ss|r})\|} \geq 0. \quad (29)$$

where preservation of semantic collapse is not guaranteed due to the distance between  $e'_{f|s}$  and  $e_{f|s}$  is not contained. This infers that when two slightly semantic distinct textual embeddings are given, the generated images' semantics can also be the same.

Assume  $L_r = 0$ , we have:

$$|e''_{f|s} - e_{f|s}| = |\epsilon^*| \odot \beta \cdot d(\mathbb{C}_{rr|s}) \quad (30)$$

$$= \beta \cdot d(\mathbb{C}_{rr|s}) \quad (31)$$

$$> 0, \quad (32)$$

where provides tighter constraints than  $L_{db}$ . Note that the proof details of  $\beta \cdot d(\mathbb{C}_{rr|s}) > 0$  can refer to Appendix Proposition A.5.

□

## B MORE DETAILS OF APPLICATION

If there is no specification, for all frameworks, we use their original losses  $L(\theta, X)$  and  $L(\theta, X')$  with  $GisC$ :  $L_{db}$  or  $L_r$ . See specified parameter settings in Table 1. We then demonstrate detailed implementations for tested frameworks. Note that since  $ITA_C$  needs no more training, thus model with  $ITA_C$  requires no more implementation of  $L_S$  in Eq. (??) and  $L_{ITA_T}$  in Eq. (??).

Section B.1 are implementations of  $L_S$  in Eq. (??) and  $L_{ITA_T}$  in Eq. (??) with different backbone: AttnGAN [15], DF-GAN [10] and VQ-GAN + CLIP [14].

### B.1 Applying $ITA_T$

**B.1.1 Applying  $ITA_T$  to DF-GAN.** DF-GAN [10] is a currently proposed one-way output T2Isyn backbone. For  $L_{DF_D} = (\theta_{DF}, \cdot)$  of DF-GAN's Discriminator  $D_{DF}$ , we use it as  $L_S$  for DF-GAN:

$$\begin{aligned} L_{S-DF_D} = & L_{DF_D}(\theta_{DF_D}, (e_s|_r, \mathcal{G}_{DF}(e_s|_r)) + \\ & L_{DF_D}(\theta_{DF_D}, (e'_s|_r, \mathcal{G}_{DF}(e'_s|_r)) + \\ & L_{DF_D}(\theta_{DF_D}, (e_s|_r, \mathcal{G}_{DF}(e'_s|_r))), \end{aligned} \quad (33)$$

where  $L_{DF}$  is the simplified representation for DF-GAN's original Discriminator losses;  $\mathcal{G} = h_{DF}(\mathcal{G}_{DF}(\cdot))$  where  $(\cdot)$  takes a textual embedding,  $h_{DF}$  maps generated images of ( $\mathcal{G}_{DF}$  on the textual embedding). Notations in the following frameworks are similar. All embeddings used in DF-GAN are gained from DAMSM images and text encoders.

Then for generator  $G$  loss  $L_{DF_G}(\theta_{DF_G}, \cdot)$ , we have loss:

$$\begin{aligned} L_{S-DF_G} = & L_{DF_G}(\theta_{DF_G}, (e_s|_r, \mathcal{G}_{DF}(e_s|_r))) + \\ & L_{DF_G}(\theta_{DF_G}, (e'_s|_r, \mathcal{G}_{DF}(e'_s|_r))) + \\ & L_{DF_G}(\theta_{DF_G}, (e_s|_r, \mathcal{G}_{DF}(e'_s|_r))). \end{aligned} \quad (34)$$

Since DF-GAN only uses one discriminator  $D_{DF}$  for both semantic matching and image quality supervision. Therefore, we can use  $L_{DF_D}(\theta_{DF}, \cdot)$  to force  $\mathcal{G}_{DF}(t'_s)$  be consistent with  $t_s$  by optimizing parameters  $r$  of  $ITA_T$ :

$$\begin{aligned} L_{ITA_T-DF} = & r \cdot L_{iemse}(e_s|_r, e'_s|_r) + \\ & (1 - r) \cdot L_{DF_G}(\alpha, (e_s|_r, \mathcal{G}_{DF}(e'_s|_r))). \end{aligned} \quad (35)$$

**B.1.2 Applying  $ITA_T$  to AttnGAN.** AttnGAN [15] is a widely used backbone for GAN-based text-to-image generation baseline. Since AttnGAN uses both sentence-level  $e_s|_r$  and word-level semantics  $e_w$  embeddings, we implement augmented sentence and words  $e'_s|_r, e'_w$  as  $e'_s|_r = ITA_T(e_s|_r), e'_w = e_w + (e'_s|_r - e_s|_r)$ . Other implementations refer to Section B.1.1.

All embeddings used in AttnGAN are gained from DAMSM images and text encoders.

**B.1.3 Applying  $ITA_T$  to VQ-GAN + CLIP.** We use the released checkpoint and code of [14] for tuning. Notice the [14] is originally trained on the clip embeddings of images; we directly altered it by using textual CLIP embeddings. We only tune the transformer part for organizing the discrete code, while the image-generating

decoder part is fixed. Due to the long training time, we only tune the model for 20 epochs with  $L_r$  and use its original  $L_{vqclip}(\theta, X)$  for our augmented  $X'$  as  $L_{vqclip}(\theta, X')$ . Other settings follow the original settings. All embeddings used in VQ-GAN + CLIP are gained from CLIP images and text encoders. We only test  $ITA_C + L_r$  with VQ-GAN + CLIP due to its long training time.

**B.1.4  $ITA_T$  Implementation Suggestions.** Training  $ITA_T$  with adversarial models needs concern about how to avoid exploding gradient. Because in the early stage, the discriminators may not provide meaningful semantic bounding on  $ITA_T$ , causing the augmented  $e'_s|_r$  located too far from  $e_s|_r$  and then a too large loss for generators which cannot be optimized.

Thus we suggest a warmup phase before training  $ITA_T$ . For AttnGAN, we set a warmup phase to avoid this kind of crush. Due to DF-GAN using hinge losses, which cannot be larger than one, it can have no warmup phase. Refers to Table 1 for more parameter details.

We also suggest scaling the learning rate up when training with  $ITA$  with  $L_{db}$  or  $L_r$  due to their regularity.

### B.2 Applying $ITA_C$ and $L_r$

$ITA_C$  and  $L_r$  are based on  $\mathbb{C}_{ss|_r}$  and  $\mathbb{C}_{rr|_s}$  defined as:

$$\mathbb{C}_{ss|_r} = \mathbb{C}_{ss} - \mathbb{C}_{sr}\mathbb{C}_{rr}^{-1}\mathbb{C}_{rs}, \quad (36)$$

$$\mathbb{C}_{rr|_s} = \mathbb{C}_{rr} - \mathbb{C}_{rs}\mathbb{C}_{ss}^{-1}\mathbb{C}_{sr}. \quad (37)$$

We only used 30K random samples from CUB and COCO training sets, respectively, to obtain  $\mathbb{C}_{ss|_r}$  and  $\mathbb{C}_{rr|_s}$  for our experiments. The number of samples follows it of calculating FID [7]. It is rational to scale the number of samples up according to the size of the dataset. Nevertheless, we do not recommend using the whole training set for the calculation due to its memory consumption. Our calculated  $\mathbb{C}_{ss|_r}$  and  $\mathbb{C}_{rr|_s}$  will be released with our code.

## B.3 Applying SADA to Diffusion Models

For Diffusion models,  $ITA$  should be applied to conditional embeddings (including textual conditional embeddings). The  $GisC$  should be applied to features of generated images at each step.

## C MORE EXPERIMENTAL DETAILS

### C.1 Experimental Settings

We compare Mixup [16], Random Mask, Add Noise, with our proposed  $ITA_T$ ,  $ITA_C$  and  $L_r$  by applying them to AttnGAN [15], DF-GAN [10] with DAMSM encoders [15]. We evaluate each framework on two benchmark datasets, CUB [12] and COCO [9]. To delve deeper, we apply  $ITA_C$  with  $L_r$  to VQ-GAN+CLIP [14] with CLIP on COCO and the conditional DDPM [8] by specifically utilizing the MNIST dataset [4].

Specifically, experiments based on the conditional DDPM [8], specifically utilizes the MNIST dataset [4]. The methodology applied involved incorporating our proposed  $ITA_C$  on condition embeddings, with further integration of  $L_r$  on calculated feature shift of generated images from U-Net's bottleneck. We first train the bare diffusion model and then use its bottleneck's hidden feature as  $e_s|_r$  and the bottleneck's hidden feature of the next step as  $e_f|_s$ . Then other details will be as same as aforementioned.

**Table 1: Parameters for experiments.**

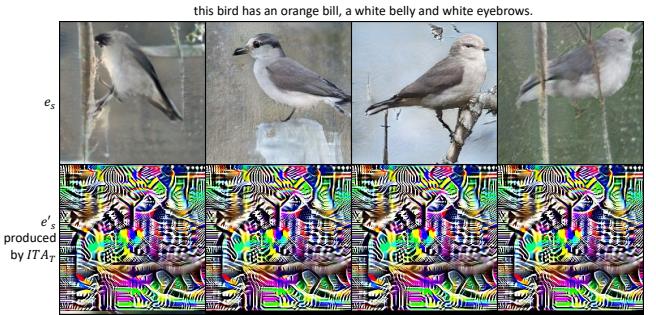
Dataset: CUB + $ITA_T$		
Backbone	Warm-up	$r$
AttnGAN	50	0
DF-GAN	100	0.2
Dataset: COCO + $ITA_T$		
Backbone	Warm-up	$r$
AttnGAN	0	0
DF-GAN	0	0.2
Dataset: CUB + $ITA_C$		
Backbone	$\beta$	$r$
AttnGAN	0.05	0
DF-GAN	0.05	0.2
Dataset: COCO + $ITA_C$		
Backbone	$\beta$	$r$
AttnGAN	0.01	0
DF-GAN	0.01	0.2
VQ-GAN + CLIP	0.05	0.2
Dataset: CUB + $ITA_C + L_r$		
Backbone	$\varphi$	
AttnGAN	0.01	
DF-GAN	0.01	
Dataset: COCO + $ITA_C + L_r$		
Backbone	$\varphi$	Learning Rate
AttnGAN	0.001	As original
DF-GAN	0.001	Doubled
VQ-GAN + CLIP	0.05	As original

Parameter settings follow the original models (including augmentations they used) of each framework for all experiments unless specified. For training settings, we train the model from scratch and also use their released model for tuning experiments. The discriminators are retrained during the tuning process for AttnGAN and DF-GAN since no released checkpoints are available; and we only tune the transformer part of experimental settings, which are specified in Table ???. Notice that we do not conduct  $ITA_T$  with  $L_r$  because  $e'_{s|r} - e_{s|r} \leq \epsilon \odot \beta \cdot d(\mathbb{C}_{ff|s})$  in Eq. (??) is required for  $L_r$ . See more training details and model parameter settings in Appendix C.

Especially,  $\mathbb{C}_{ss|r}$  for  $ITA_C$  and  $\mathbb{C}_{rr|s}$  for  $L_r$  are calculated on the training set using the encoders of the framework. We use 30K random samples from each dataset in our experiments. Limited sampling also leads to the possible implementation of  $ITA_C$  and  $L_r$  on super-large datasets.

## C.2 Implementing Mixup with AttnGAN and DF-GAN

We use the official code of Mixup [16] for our experiments. Augmented images and mixed textual embeddings are used for model training. All model settings follow the original settings of AttnGAN and DF-GAN.

**Figure 1: Collapse examples of DF-GAN with  $ITA_T + L_{db}$  using  $\alpha = 0.3$  generated on  $e_{s|r}$  and augmented  $e'_{s|r}$ .**

## C.3 Implementing with AttnGAN and DF-GAN

We randomly mask 15% tokens and use the original settings of AttnGAN and DF-GAN. AttnGAN with Random Mask collapsed multiple times during the training. We use the checkpoints that were saved before the collapse to resume the training.

## C.4 Parameter Settings

We train each backbone from the start on the CUB dataset and tune their released checkpoint on the COCO dataset. Due to no released checkpoints for discriminators of AttnGAN and DF-GAN, we retrain discriminators during the tuning phase. If there is no specification, we follow the original experimental settings of each backbone. Specified parameters used for producing final results in the paper are shown in Table 1. Notice that  $r$  for  $ITA_T$  can be set to zero due to the weak supervision of generative adversarial networks. Specifically, we double the learning rate for  $ITA_C + L_r$  tests due to their regularity.

**Table 2: Results of DF-GAN with  $ITA_T + L_{db}$  using different  $r$  values on the CUB dataset, training within 600 epochs.**

$r$	CS	FID
0	0.5791	13.96
0.1	0.5793	12.7
0.2	<b>0.5807</b>	<b>11.74</b>
0.3	$ITA_T$ collapses	

## D MORE RESULTS

### D.1 More Qualitative Results

*D.1.1  $ITA_T$  with different  $r$ .* As stated,  $r$  in Eq. (??) can control the augmentation strength. Larger  $r$  in  $ITA_T$  leads to more intensive augmentation. However, as shown in Table 2, an inappropriate large  $r$  can cause model collapse because  $S(\alpha, (e_{s|r}, G(e'_{s|r}))$  will lose its constraint, causing that  $e'_{s|r}$  is too different from  $e_{s|r}$  (i.e.,  $e'_{s|r}$  cannot maintain the main semantics of  $e_{s|r}$ ). Collapse examples are shown in Figure 1. It can be seen that using  $\alpha = 0.3$   $ITA_T$  cannot produce a semantic maintaining  $e'_{s|r}$  for  $G$ . Within the appropriate range, larger  $r$  offers better text-image consistency and image quality.

465

466

A small bird with a red head, breast, and belly and black wings.

523

467



524

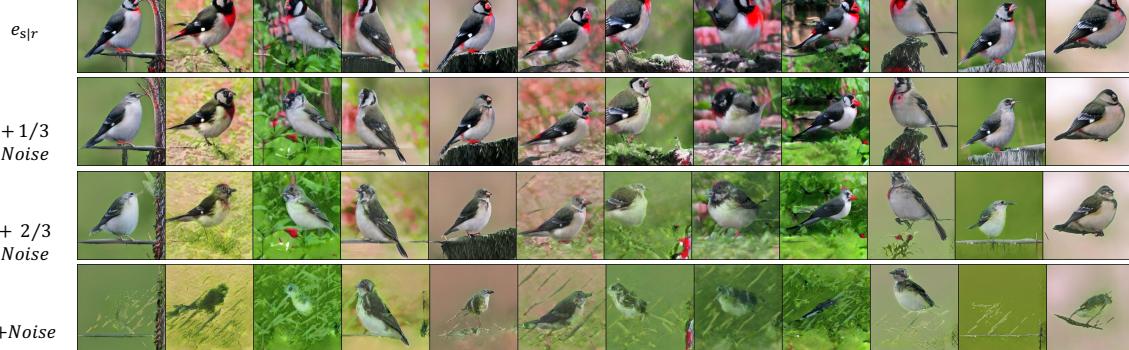
468

469

+ DiffAug

525

470



526

471

527

472

528

473

529

474

530

475

531

476

532

477

533

478

534

479

535

480

536

481

537

482

538

483

$\epsilon_{\text{slr}}$

+ 1/3 Noise

+ 2/3 Noise

+Noise

539

484

540

485

541

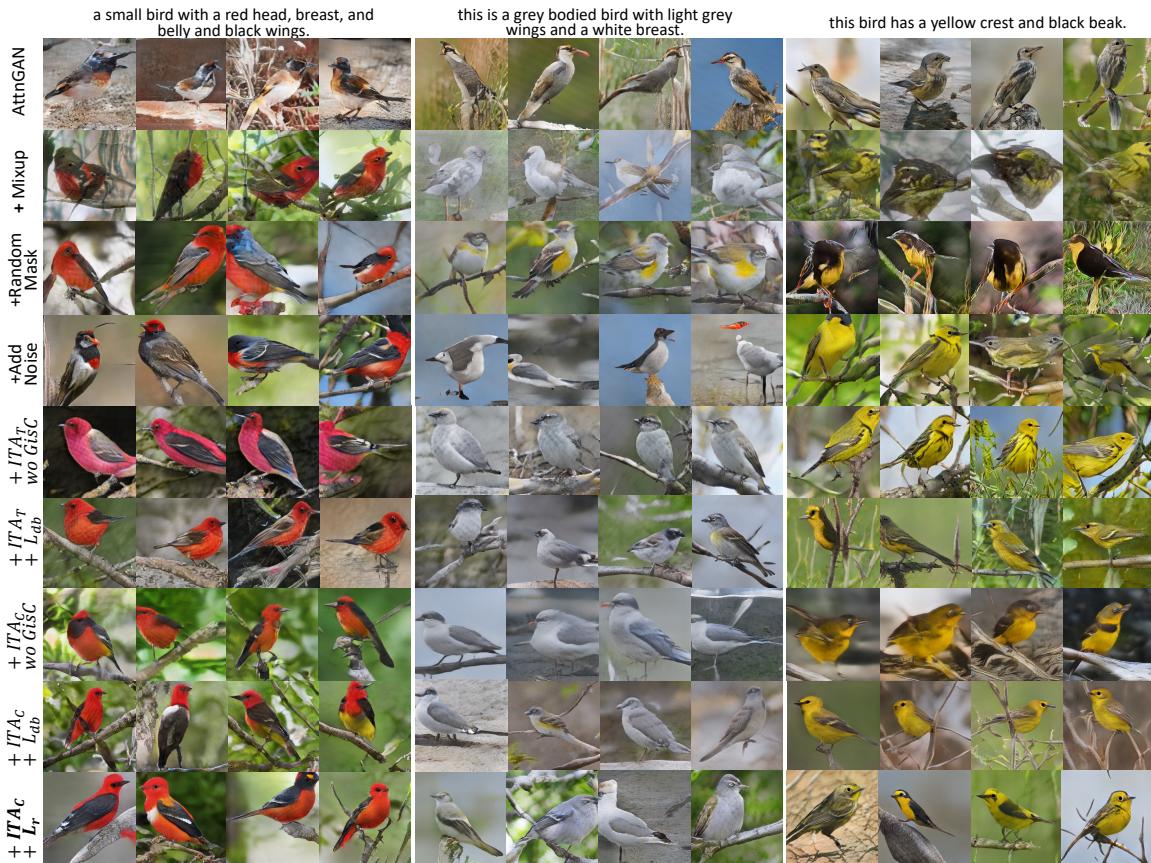
486

542

**Figure 2: Generated images of DF-GAN with DiffAug on CUB dataset with ascending scales of  $\epsilon$  are added.**

543

544



545

546

487

547

488

548

489

549

490

550

491

551

492

552

493

553

494

554

495

555

496

556

497

557

498

558

499

559

500

560

501

561

502

562

503

563

504

564

505

565

506

566

507

567

508

568

509

569

510

570

511

571

512

572

513

573

514

574

515

575

516

576

517

577

**Figure 3: Generated results of AttnGAN on CUB.**

578

579

580

581     D.1.2 *Semantic Collapse in DF-GAN with DiffAug.* We exhibit generated  
 582     images of DF-GAN with DiffAug on the CUB dataset. As the  
 583     ascending scales of  $\epsilon$  are added, the images illustrate the occurrence  
 584     of semantic collapse.

## 585     D.2 More Quantitative Results

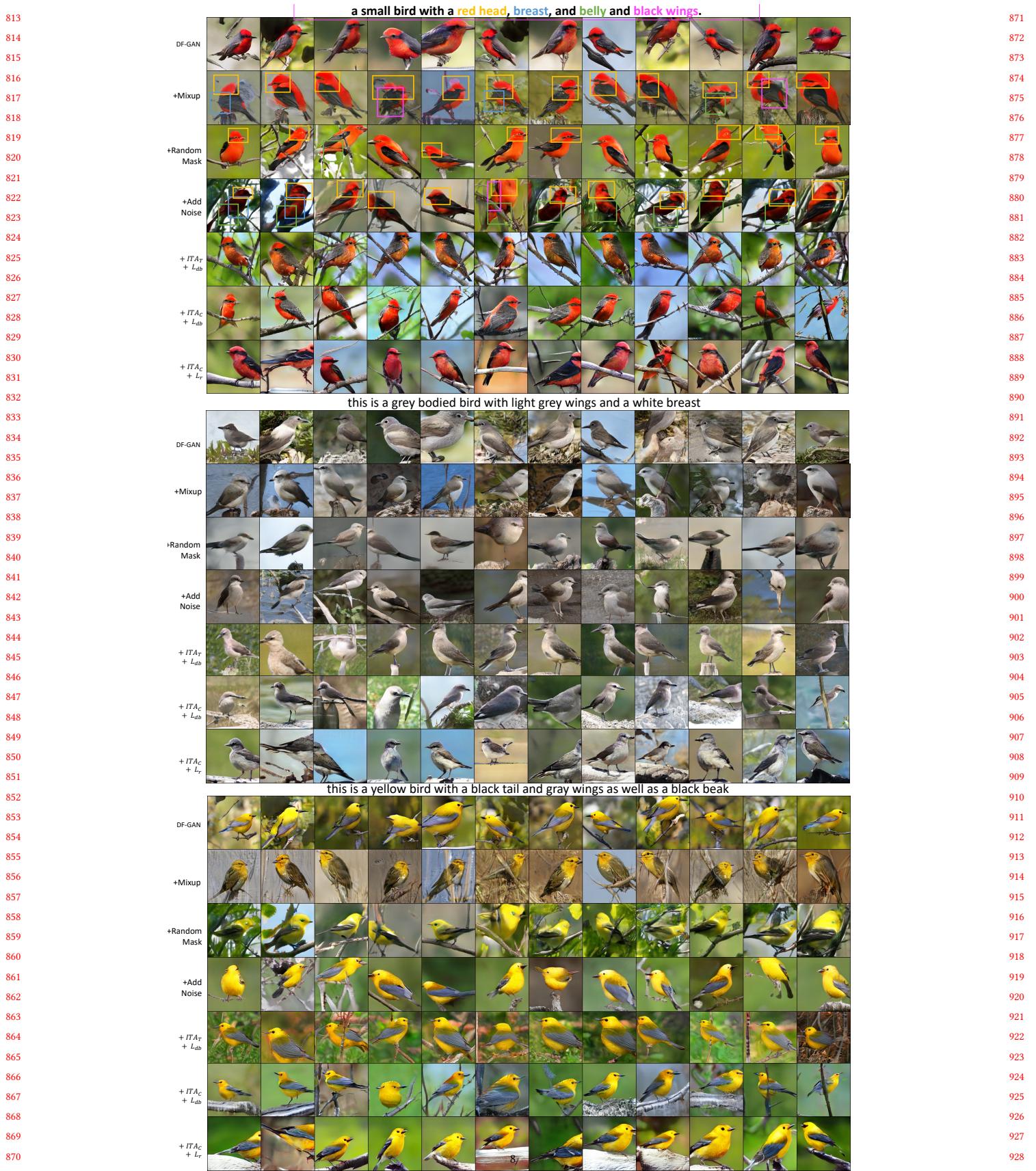
587     We show more generated examples of AttnGAN [15] as Figure 3  
 588     and Figure 4, DF-GAN [10] as Figure 5 and Figure. 6, and VQ-GAN  
 589     + CLIP [14] on COCO as Figure 7. The diversity improvement is  
 590     more obvious for AttnGAN than other backbones.

## REFERENCES

- [1] Yoshua Bengio, Grégoire Mesnil, Yann Dauphin, and Salah Rifai. 2013. Better mixing via deep representations. In *International Conference on Machine Learning*. PMLR, 552–560. 639
- [2] Shuxiao Chen, Edgar Dobriban, and Jane H Lee. 2020. A group-theoretic framework for data augmentation. *The Journal of Machine Learning Research* 21, 1 (2020), 9885–9955. 640
- [3] Jun Cheng, Fuxiang Wu, Yanling Tian, Lei Wang, and Dapeng Tao. 2020. RiFeGAN: Rich feature generation for text-to-image synthesis from prior knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10911–10920. 641
- [4] Li Deng. 2012. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine* 29, 6 (2012), 141–142. 642
- [5] Hao Dong, Jingqing Zhang, Douglas McIlwraith, and Yike Guo. 2017. I2t2i: Learning text to image synthesis with textual data augmentation. In *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2015–2019. 643
- [6] Rinot Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–13. 644
- [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in Neural Information Processing Systems* 30 (2017). 645
- [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851. 646
- [9] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*. Springer, 740–755. 647
- [10] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Changsheng Xu. 2022. Df-gan: A simple and effective baseline for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16515–16525. 648
- [11] Paul Upchurch, Jacob Gardner, Geoff Pleiss, Robert Pless, Noah Snavely, Kavita Bala, and Kilian Weinberger. 2017. Deep feature interpolation for image content changes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7064–7073. 649
- [12] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. 2011. The caltech-ucsd birds-200-2011 dataset. (2011). 650
- [13] Yulin Wang, Gao Huang, Shiji Song, Xuran Pan, Yitong Xia, and Cheng Wu. 2021. Regularizing deep networks with semantic data augmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021). 651
- [14] Zihao Wang, Wei Liu, Qian He, Xinglong Wu, and Zili Yi. 2022. CLIP-GEN: Language-Free Training of a Text-to-Image Generator with CLIP. *arXiv preprint arXiv:2203.00386* (2022). 652
- [15] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1316–1324. 653
- [16] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017). 654



**Figure 4: Generated results of AttnGAN on COCO.**



**Figure 5: Generated results of DF-GAN on CUB. Semantic mismatches and image quality degradation are highlighted for generated results of DF-GAN w/ Mixup, w/ Random Mask, and w/ Add Noise in the top group.**



Figure 6: Generated results of DF-GAN on COCO.

1045	a distinct looking planter with flowers sitting on a table		a big brown cow with horns standing in a big grassy field		a couch and chair are sitting in a room		a giraffe standing near a tree branch in the grass near a grove of trees		a train is parked on the train tracks at the station		1103
1046										1104	
1047										1105	
1048										1106	
1049										1107	
1050										1108	
1051	VQ-GAN With CLIP									1109	
1052										1110	
1053										1111	
1054										1112	
1055										1113	
1056										1114	
1057	+ $ITA_c$ + $L_r$									1115	
1058										1116	
1059										1117	
1060										1118	
1061										1119	
1062										1120	
1063										1121	
1064										1122	
1065										1123	
1066	VQ-GAN With CLIP									1124	
1067										1125	
1068										1126	
1069										1127	
1070										1128	
1071										1129	
1072	+ $ITA_c$ + $L_r$									1130	
1073										1131	
1074										1132	
1075										1133	
1076										1134	
1077										1135	
1078										1136	
1079										1137	
1080										1138	
1081										1139	
1082										1140	
1083										1141	
1084										1142	
1085										1143	
1086										1144	
1087										1145	
1088										1146	
1089										1147	
1090										1148	
1091										1149	
1092										1150	
1093										1151	
1094										1152	
1095										1153	
1096										1154	
1097										1155	
1098										1156	
1099										1157	
1100										1158	
1101										1159	
1102										1160	

Figure 7: Generated results of VQ-GAN + CLIP on COCO.