

1 We authors are grateful to the reviewers for their valuable comments. We will improve the final version by taken all the review
 2 comments and release the source code package to ensure the reproducibility. Below, we number and address comments of
 3 each reviewer in order.

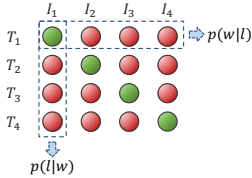


Figure 1: Illustration of **matched** & **mismatched** (T_n, I_n) pairs in a minibatch.

Table 1: Results of LeicaGAN trained with E_T^I and E_T^M of different loss weights.

Evaluation Metric	Inception Score		R-precision [%]	
	CUB*	Oxford-102*	CUB*	Oxford-102*
(1) LeicaGAN, $E_T^I(\alpha_1=0), E_T^M(\alpha_2=0)$	5.51±0.03	3.64±0.01	82.89	84.02
(2) LeicaGAN, $E_T^I(\alpha_1=1), E_T^M(\alpha_2=0)$	5.58±0.05	3.71±0.02	85.20	85.55
(3) LeicaGAN, $E_T^I(\alpha_1=5), E_T^M(\alpha_2=0)$	5.60±0.05	3.68±0.01	82.95	84.03
(4) LeicaGAN, $E_T^I(\alpha_1=1), E_T^M(\alpha_2=2)$	5.63±0.04	3.79±0.01	84.57	84.98
(5) LeicaGAN, $E_T^I(\alpha_1=1), E_T^M(\alpha_2=4)$	5.69±0.05	3.80±0.01	85.28	85.81
(6) LeicaGAN, $E_T^I(\alpha_1=1), E_T^M(\alpha_2=6)$	5.61±0.06	3.64±0.02	81.65	84.01

4 **R1.1** - Why is the probability of w matched by l calculated over the minibatch in Eq (3)?

5 \Rightarrow Assuming a minibatch of 4 pairs of (T_n, I_n), where $n \in \{1, 2, 3, 4\}$ as shown in Fig.1, we treat the **green** and **red**
 6 pairs as **matched** and **mismatched** pairs respectively. Therefore, for each text, we obtain 4 similarity scores indicating the
 7 closeness between the given text with the **matched** & **mismatched** images (vice versa). Afterwards, according to Eq.(4),
 8 $p(w|l)$ and $p(l|w)$ are estimated separately. Specifically, a softmax is applied among the matched and mismatched text-image
 9 (image-text) pairs in a batch along the row (column) as shown in Fig.1. They serve as different constraints for training the
 10 model at the same time.

11 **R1.2** - Why not use a symmetric similarity measure? An exploration of the effect of using a cosine similarity.

12 \Rightarrow The local-level similarity score $s_{w|l}$ or $s_{l|w}$ of the (T, I) pair is first calculated as: $\log(\sum_{i=1}^L \exp(\gamma_1 \cos(\tilde{l}_i, w_i)))^{1/\gamma_1}$ [37],
 13 which is a symmetric similarity result and we elaborate this equation here for a more clear explanation of the calculation.
 14 Careful revision of this part and further exploration about different similarity measures will be included in the final version.

15 **R1.3** - The imagination phase is just sampling a noise vector.

16 \Rightarrow The embedded textual features obtained in TVE convey different kinds of visually-relevant information. Therefore, in the
 17 imagination phase, an aggregation of these features and noise is applied to obtain an initial visual impression serving as
 18 the input for the subsequent creation phase. The noise is added to guarantee the diversity of the image generation. In the
 19 final version, we will explore more efficient and diverse modules to strengthen the impact of the imagination phase as we
 20 discussed in the section of Limitation.

21 **R1.4** - What is the effect of having the adversarial term in the encoder loss? An ablation study of this adversarial term.

22 \Rightarrow The loss weights α_1 and α_2 in training text-visual co-embedding models E_T^I and E_T^M , respectively, were carefully studied
 23 in Table 1. The detailed results will be reported in the final version. The comparison between (1) and (2) shows that the use
 24 of the adversarial loss obtains a performance gain. An intuitive illustration of the effect of the adversarial loss reducing the
 25 multi-modal domain gap has been shown in Fig.4 of [35].

26 **R2.1** - How can the model be able to control each word as presented in the experiments?

27 \Rightarrow The input sentence was first embedded into both (1) word-level feature matrices w to represent the semantic meaning of
 28 the words and (2) a sentence-level feature vector s to convey the semantic of the whole sentence. Afterwards, through the
 29 co-embedding of images and text, the correlations between them can be built. Therefore, during the decoding stage, with the
 30 guidance of the attention module, the generators focused on different words, resulting in different visual outputs accordingly.

31 **R2.2** - This paper only compares with AttnGAN and ablation studies of key components are missing.

32 \Rightarrow Before the NeurIPS submission deadline, StackGAN and AttnGAN were recognized as state-of-the-art models in T2I.
 33 Recently, Obj-GANs, SD-GAN, MirrorGAN, and DM-GAN are published in CVPR2019. However, it is not straightforward
 34 to compare LeicaGAN with all of them because they were tested on different datasets using different evaluation metrics and
 35 only two of them (MirrorGAN and Obj-GAN) released their source code packages in the last two months, after the NeurIPS
 36 submission deadline. To make a fair comparison, we will re-train and test these models under the same condition in the final
 37 version. Based on the results in Table 2 in the paper, we verified the effectiveness of the following key components: (1) the
 38 global-local attentive generator, (2) E_T^M , and (3) the collaborative attention module and its weights. Here, we report the
 39 ablation study for the weight choosing of the TVE models in Table 1. More details will be provided in the appendix and code.

40 **R2.3** - The paper names the text-image encoded feature as the prior knowledge, which is overclaimed.

41 \Rightarrow We name the text-image encoded features as the prior knowledge because the co-embedding mappings, which correlate
 42 textual and visual features with each other, are learned in the PKL phase. These cross-modal embeddings serve as the prior
 43 knowledge to guide the image generation of the subsequent phases. We can expect that this phase can be further improved by
 44 incorporating different types of prior. We agree the name of "prior knowledge" can be inappropriate and will change it to
 45 "prior" in the final version.

46 **R3.1** - Improving the limitations the authors already discussed in the paper.

47 \Rightarrow In the future, we plan to introduce more prior to the first PKL phase, such as fine-grained attributes, and more efficient and
 48 diverse modules will be explored to strengthen the impact of imagination. Additionally, more advanced network structures
 49 and training strategies will be employed in order to further enhance the generation ability of the proposed T2I model.