

Language-agnostic Semantic Consistent Text-to-Image Generation

SeongJun Jung¹, Woo Suk Choi¹, Seongho Choi¹ and Byoung-Tak Zhang^{1,2}

¹Seoul National University

²AI Institute (AIIS), Seoul National University

{seongjunjung, wschoi, shchoi, btzhang}@bi.snu.ac.kr

Abstract

Recent GAN-based text-to-image generation models have advanced that they can generate photo-realistic images matching semantically with descriptions. However, research on multilingual text-to-image generation has not been carried out yet much. There are two problems when constructing a multilingual text-to-image generation model: 1) language imbalance issue in text-to-image paired datasets and 2) generating images that have the same meaning but are semantically inconsistent with each other in texts expressed in different languages. To this end, we propose a Language-agnostic Semantic Consistent Generative Adversarial Network (LaSC-GAN) for text-to-image generation, which can generate semantically consistent images via language-agnostic text encoder and Siamese mechanism. Experiments on relatively low-resource language text-image datasets show that the model has comparable generation quality as images generated by high-resource language text, and generates semantically consistent images for texts with the same meaning even in different languages.

1 Introduction

In this paper, we consider multilingual text-to-image generation. There are two problems with multilingual text-to-image generation. The first problem is the language imbalance issue in text-to-image datasets. Most text-to-image generation datasets are in English, so it is difficult to construct text-to-image generation models for other languages. Furthermore, since existing multilingual datasets have a small amount of data, a discriminator overfitting may cause problems such as instability of learning in GAN. The second is that generative models have difficulty extracting semantic commonality between languages. This can produce different images for captions with the same semantics but different languages. In Yin et al. (2019), they treat the problem that captions with

same meanings in English create semantically different images. We extend this awareness between languages.

To solve those problems, we propose LaSC-GAN for text-to-image generation. LaSC-GAN consists of a language-agnostic text encoder and a hierarchical generator. Language-agnostic text encoder generates text embeddings to be used in the hierarchical generator for the first problem mentioned above. And we exploit the Siamese structure training to capture the semantic consistency between images generated in various languages.

Our main contributions are as follows: 1) By using a language-agnostic text encoder, images for low-resource language text can be generated only by learning the high-resource language. 2) Texts with the same semantics in different languages can generate semantically consistent images using the Siamese mechanism in hierarchical generator to extract semantic consistency between languages. We show the effect of each contribution in experiments using English MS-COCO (COCO-EN), Chinese MS-COCO (COCO-CN) and Korean MS-COCO (COCO-KO) datasets.

2 Related Works

2.1 Generative Adversarial Network (GAN) for Text-to-Image

Text-to-image generation using GAN has advanced a lot since GAN-INT-CLS (Reed et al., 2016). The discovery of hierarchical model architectures (Zhang et al., 2017, 2018; Xu et al., 2018) has produced realistic images that semantically match with texts. However, these models only considered image generation for a single language, and to the best of our knowledge the first paper dealing with multilingual text to image generation is Zhang et al. (2022). The model proposed in Zhang et al. (2022) requires learning for each language. However, our method can generate images from multi-

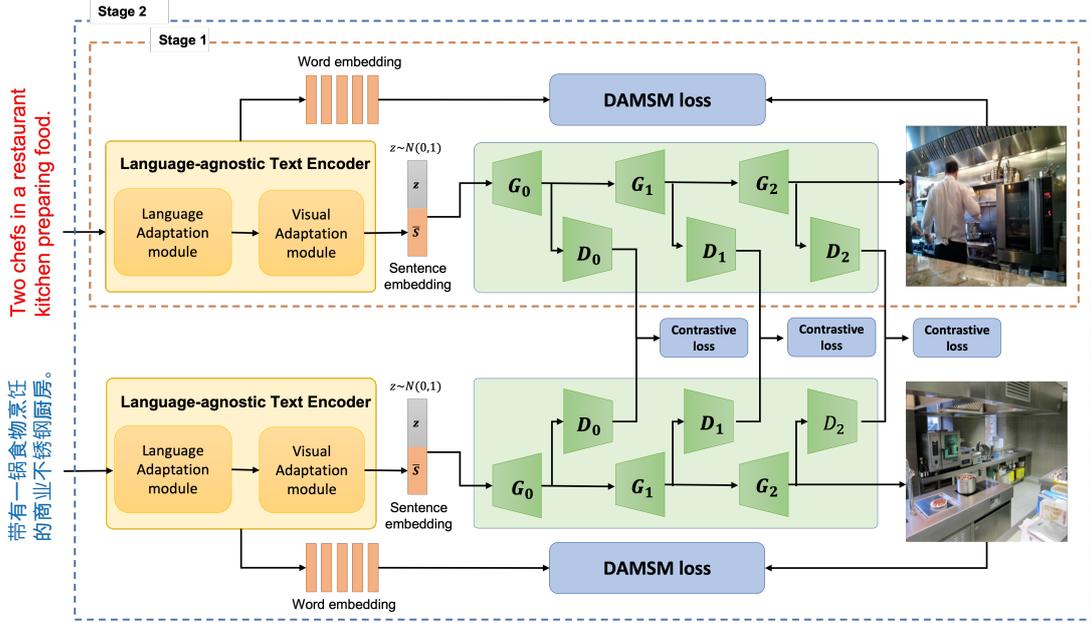


Figure 1: The architecture of LaSC-GAN. In stage1, the model is trained followed by (Xu et al., 2018) with COCO-EN. In stage 2, the text-to-image generation is trained with contrastive loss based on a Siamese structure with COCO-EN, CN, and KO.

lingual texts only by learning about high-resource language.

2.2 Multilingual Text Encoders

Multilingual text embedding models usually use the translation pairs datasets, and sometimes the translation pairs datasets and monolingual datasets are used together. Among these, language-agnostic BERT sentence embedding (LaBSE) (Feng et al., 2020) using MLM(Masked Language Model) pre-training was proposed.

3 Methods

We propose a LaSC-GAN for text-to-image generation. Our goal is to obtain as good visual quality of images created with low-resource language text as images generated with high-resource language text and to enable the model to reflect semantic consistency between languages in image generation. The LaSC-GAN consists of a language-agnostic text encoder and a hierarchical generator. The language-agnostic text encoder is used to obtain a text representation that will be fed as a condition to the generator. The hierarchical generator generates images for text conditions. Training strategy of the model consists of two stages as shown in Figure1. In stage 1, the model is trained followed by (Xu et al., 2018) using only a high-resource language dataset. In stage 2, the model is trained

with a Siamese structure with two model branches using data from different language pairs (EN-CN, EN-KO).

3.1 Model Architecture

Language-agnostic Text Encoder consists of a Language Adaptation module and a Visual Adaptation module. We use pre-trained LaBSE (Feng et al., 2020) for the Language Adaptation module and bi-directional LSTM for the Visual Adaptation module. We get language-agnostic token embeddings from each token embedding passed through the Language Adaptation module. Then, the obtained embedding is transferred to a visual representation space through the Visual Adaptation module and used as the text condition of the generator. Hidden states of each token in the bi-directional LSTM of the Visual Adaptation module are used as word embeddings, and the last hidden state is used as sentence embeddings. Our model can use 109 languages used in LaBSE training as inputs.

Hierarchical Generator uses the hierarchical generative adversarial network structure used in (Xu et al., 2018), which consists of 3 sub-generators (G_0, G_1, G_2). Each generator has an independent discriminator (D_0, D_1, D_2). The initial sub-generator generates a low-resolution image by putting the sentence representation (\bar{s}) input from the language-agnostic text encoder and a random

Metric	IS \uparrow			FID \downarrow			CLIP \uparrow		
	ST.1	ST.2	ST.1 \wedge ST.2	ST.1	ST.2	ST.1 \wedge ST.2	ST.1	ST.2	ST.1 \wedge ST.2
EN	14.89	-	-	97.41	-	-	0.227	-	-
KO	12.24	14.76	15.58	103.26	102.04	93.16	0.196	0.198	0.195
CN	14.98	16.14	16.55	97.26	93.64	93.40	0.213	0.214	0.212

Table 1: Quantitative results for each stage of LaSC-GAN. ST, EN, KO, and CN denote stage, English, Korean, and Chinese. ST.1 and ST.2 refer to models that have undergone only Stage 1 and 2 learning processes, respectively. And ST.1 \wedge ST.2 refer to a model using both learning processes together.

noise ($z \sim N(0, 1)$) from normal distribution. The following sub-generators generate a higher resolution image by using the previous generation result.

3.2 Training Strategy

In the first training stage, the model is trained followed by (Xu et al., 2018) using only a high-resource language dataset with DAMSM loss, and the parameters learned in the first stage are used in the second learning stage.

Then, in the second learning stage, we use the Siamese mechanism such as SD-GAN (Yin et al., 2019) to learn semantic commons between texts in different languages. In addition to the DAMSM loss, we compute contrastive loss as follows by using the visual features of the discriminator for the inputs to the two branches of the Siamese structure.

$$L = \frac{1}{2N} \sum_{n=1}^N y \cdot d^2 + (1 - y) \max(\epsilon - d, 0)^2 \quad (1)$$

where $d = \|v_1 - v_2\|_2$ is the distance between the visual feature vectors v_1 and v_2 from the two Siamese branches respectively, and y is a flag to mark whether the input descriptions are from the same image or not (i.e., 1 for the same and 0 for different). The hyper-parameter N is the length of the feature vector. The hyper-parameter ϵ is used to balance the distance value when $y = 0$.

4 Experiments

4.1 Datasets

We used MS-COCO (COCO-EN) (Lin et al., 2014) for stage 1. COCO-EN has 80K image train set and 40K image validation set. Each image has 5 English descriptions. We also used the multilingual versions of COCO-EN: COCO-CN and COCO-KO for stage 2. COCO-CN (Li et al., 2019) has 1 manually translated Chinese description for the 18K image train set and 1K image validation set.

We used the validation set index of COCO-CN for other languages as well. COCO-KO has Korean machine translation results for all descriptions of in COCO-EN. In stage 2, we use a subset of data from COCO-EN and COCO-KO that overlap with COCO-CN. In stage 2, EN-CN and EN-KO language pair datasets are used for training respectively. The models trained with EN-CN, EN-KO pair datasets are evaluated on the COCO-CN, COCO-KO validation set respectively.

4.2 Implementation Details

The hierarchical generator and discriminator followed (Xu et al., 2018), and the language-agnostic text encoder is comprised of LaBSE (Feng et al., 2020) and bi-directional LSTM. The Siamese mechanism learning method follows (Yin et al., 2019). We freeze the pre-trained parameters of LaBSE when learning the language-agnostic text encoder for stability of learning.

4.3 Metrics

We evaluated the visual quality of generated images using Inception Score (IS) and Fréchet Inception Distance (FID) used by Xu et al. (2018). In addition, we evaluated how much generated images are semantically similar to the conditioned texts through CLIP score used by Wu et al. (2021).

4.4 Zero-shot Language Text-to-image Generation

In this section, we shows the benefits of the language-agnostic text encoder. We trained only on the high-resource language dataset(COCO-EN) in stage 1. Thanks to the language-agnostic text encoder, our model can generate images from zero-shot languages. In Table 1, CN and KO are not used for learning in stage 1 but show metric scores that are not significantly different from EN used for learning. Figure 2 shows images generated in various languages using the stage1 model. The gener-

ated images from zero-shot language show similar visual quality to images generated with languages used for learning in Figure 2. In particular, our model can generate images from low-resource languages such as Thai(TH) and Nepali(NE).



Figure 2: Qualitative results of zero-shot language text to image generation of stage 1. In stage 1, the model was trained using only English texts. GT, EN, KO, CN, FR, TH, and NE denote ground-truth, English, Korean, Chinese, French, Thai, and Nepali respectively. The English description was translated into each language and used for generation.



Figure 3: Qualitative examples of the LaSC-GAN. The results of each stage with given a pair of language descriptions

4.5 Multilingual Semantic Consistent Text-to-image Generation

We conducted an experiment to show the effect of the Siamese mechanism training. In table 1, the model that performed stage 1 and stage 2 together showed better performance in IS and FID than the models performed separately. And we compute the FID of the language pairs (EN-KO, EN-CN) to shows that stage 2 helps the model to generate

FID	EN-KO	EN-CN
ST.1	57.74	51.04
ST.1 \wedge ST.2	57.32	49.56

Table 2: FID between languages in each stage.



Figure 4: Image generation results of the LaSC-GAN. The images were generated with sentences in which the nouns in the English description were replaced with Chinese and Korean nouns, respectively.

semantically consistent images if the semantics are the same in different languages. As shown in Table 2, it can be confirmed that the distance has gotten closer after stage 2. In addition, images generated from texts in different languages with the same meaning have similar images as shown in Figure 3. And Figure 4 shows the model can extract semantic commons between languages.

5 Conclusion

In this paper, we propose a LaSC-GAN for text-to-image generation. Through language-agnostic text encoder, the model can generate images with low-resource language texts in zero-shot setting. Furthermore, by Siamese mechanism, the model can extract high-level consistent semantics between languages when generating images. The experiments on COCO-EN, KO, and CN show that our proposed method can generate photo-realistic images from the relatively low-resource language text and extract semantic commons between languages for image generation.

Acknowledgements

This work was partly supported by the IITP (2015-0-00310-SW.Star-Lab/20%, 2018-0-00622-RMI/15%, 2019-0-01371-BabyMind/20%, 2021-0-02068-AIHub/15%, 2021-0-01343-GSAI (SNU)/15%) grants, and the CARAI

(UD190031RD/15%) grant funded by the DAPA and ADD.

References

- Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.
- Xirong Li, Chaoxi Xu, Xiaoxu Wang, Weiyu Lan, Zhengxiong Jia, Gang Yang, and Jieping Xu. 2019. Coco-cn for cross-lingual image tagging, captioning, and retrieval. *IEEE Transactions on Multimedia*, 21(9):2347–2360.
- 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*, pages 740–755. Springer.
- Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *International conference on machine learning*, pages 1060–1069. PMLR.
- Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. 2021. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*.
- 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*, pages 1316–1324.
- Guojun Yin, Bin Liu, Lu Sheng, Nenghai Yu, Xiaogang Wang, and Jing Shao. 2019. Semantics disentangling for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2327–2336.
- Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2017. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915.
- Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2018. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1947–1962.
- Han Zhang, Suyi Yang, and Hongqing Zhu. 2022. Cjetig: Zero-shot cross-lingual text-to-image generation by corpora-based joint encoding. *Knowledge-Based Systems*, 239:108006.