

# Unpaired Text-to-Image-to-Text Translation using Cycle Consistent Adversarial Networks

Jeremy Ma, Satya Krishna Gorti

## 1 Introduction

Text-to-Image synthesis is a challenging problem that has a lot of room for improvement considering the current state-of-the-art results. Synthesized images from existing methods give a rough sketch of the described image but fail to capture the true essence of what the text describes. The recent success of Generative Adversarial Networks (GANs) [2] indicate that they are a good candidate for the choice of architecture to approach this problem.

However, the very nature of this problem is such that a piece of text can map to multiple valid images. The lack of such a direct one-to-one mapping means that traditional conditional GANs [5] cannot be used directly. We draw our inspiration from the recent works of image-to-image translation [4][10] where a cycle consistent GANs have been trained and achieved very impressive results.

We believe that using a cycle GAN for text-to-image-to-text translation will generate better results than existing approaches and give more photo-realistic images. The added advantage of framing the problem in a cycle consistent manner would also mean that the architecture can not only be a text-to-image synthesizing network but also an image captioning network.

Therefore, we have two generators  $G$  and  $F$ . We train a mapping  $G : T_{emb} \mapsto Y$  and inverse mapping  $F : Y \mapsto T_{emb}$  in a cycle consistent manner, where  $T_{emb}$  is an embedding for the text that describes an image. The generators  $G$  and  $F$  have their corresponding discriminator  $D_g$  and  $D_f$ .

## 2 Related work

Generative Adversarial Networks (GANs) have achieved impressive results in problems such as image generation [site chintala]. Conditional GANs introduced in [5] build on top of GANs by learning to approximate the distribution of data by conditioning on an input.

In the recent past there have been attempts on text to image synthesis using conditional GANs such as [6][1][7][8]. We can see promising results in [6] by conditioning the GAN on text descriptions instead of class labels. [7] uses a similar approach but breaks the process of generation down into two stage process. Stage-1 produces a low-resolution image based on text description. Stage-2 takes Stage-1 result as input and generates high resolution photo-realistic images. [1] additionally conditions its generative process with both text and class information and has produces superior results compared to [scott]. The embeddings to represent text used in the aforementioned papers is Skip-Thought vectors [3].

Cycle consistent GANs have showed excellent results for multimodal learning problems, which lack a direct one-to-one correspondence with input and output and allows the network to learn many mappings at the same time as shown in [10][4][9]. But to the best of our knowledge, it still hasn't been used for text-to-image generation, which is a similar multimodal learning problem.

### 3 Project plan

### 4 Nice-to-haves

## References

- [1] Ayushman Dash, John Cristian Borges Gamboa, Sheraz Ahmed, Muhammad Zeshan Afzal, and Marcus Liwicki. Tac-gan-text conditioned auxiliary classifier generative adversarial network. *arXiv preprint arXiv:1703.06412*, 2017.
- [2] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [3] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *Advances in neural information processing systems*, pages 3294–3302, 2015.
- [4] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017.
- [5] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

- [6] Scott Reed, Zeynep Akata, Xincheng Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*, 2016.
- [7] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiao lei Huang, Xiaogang Wang, and Dimitris Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *IEEE Int. Conf. Comput. Vision (ICCV)*, pages 5907–5915, 2017.
- [8] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiao lei Huang, and Dimitris N. Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *CoRR*, abs/1710.10916, 2017.
- [9] Tinghui Zhou, Philipp Krahenbuhl, Mathieu Aubry, Qixing Huang, and Alexei A Efros. Learning dense correspondence via 3d-guided cycle consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 117–126, 2016.
- [10] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017.