

Multi-Modality Image Inpainting using Generative Adversarial Networks

Aref Abedjooy¹ and Mehran Ebrahimi²[0000-0002-3980-9582]

¹ Faculty of Science, Ontario Tech University, Oshawa, Ontario, Canada
Aref.AbedjooyDivshali@ontariotechu.net

² Faculty of Science, Ontario Tech University, Oshawa, Ontario, Canada
Mehran.Ebrahimi@ontariotechu.ca

Abstract. Deep learning techniques, especially Generative Adversarial Networks (GANs) have significantly improved image inpainting and image-to-image translation tasks over the past few years. To the best of our knowledge, the problem of combining the image inpainting task with the multi-modality image-to-image translation remains intact. In this paper, we propose a model to address this problem. The model will be evaluated on combined night-to-day image translation and inpainting, along with promising qualitative and quantitative results.

Keywords: Image-to-image translation · Image inpainting · Generative adversarial network · Deep learning

1 Introduction

Image-to-image translation and image inpainting are both challenging tasks. Translating an image from one form, or modality, to another may involve generating an entirely new and realistic version of the image. Image-to-image translation has different application domains including digital arts, medical imaging [6]. The image inpainting task entails filling in missing regions in an image so that the whole image appears realistic. Image inpainting is an important step in many photo editing tasks. For example, an image would have a missing area after the removal of an unwanted object.

An interesting challenge would be to translate and inpaint images at the same time. In both tasks, the model is required to generate realistic outputs. For example one may pose the problem of recovering an ideal day-time image of a scene, given a night-time image of the same scene where parts of the image is also missing. The combined image inpainting and translation tasks can be more challenging to create a realistic image.

Deep Learning techniques have been successful at image-to-image translation and image inpainting tasks and we plan to combine and apply the existing generative models to address the problem as described in the following Sections.

In the following Sections we cover the related work, followed by the methodology, experiments and results, ablation study, conclusions and discussions.

2 Related work

To the best of our knowledge, no studies have been conducted to combine image inpainting and image-to-image translation problems using Deep Learning (DL). Image generation

with generative adversarial networks (GANs) has gained remarkable progress recently. A wide variety of image synthesis tasks, such as image editing, image composition, etc., have been investigated extensively using GANs. GANs provide plausible results for image inpainting specifically [13][14] [8][12] [2] [9]. Solving various image-to-image translation tasks with conditional GANs was introduced by *pix2pix* [5]. Several studies have been conducted since then regarding the use of GANs to translate images [10] [15] [11] [1] [3].

3 Methodology

In this Section, we present the problem description along with a possible methodology to address the problem.

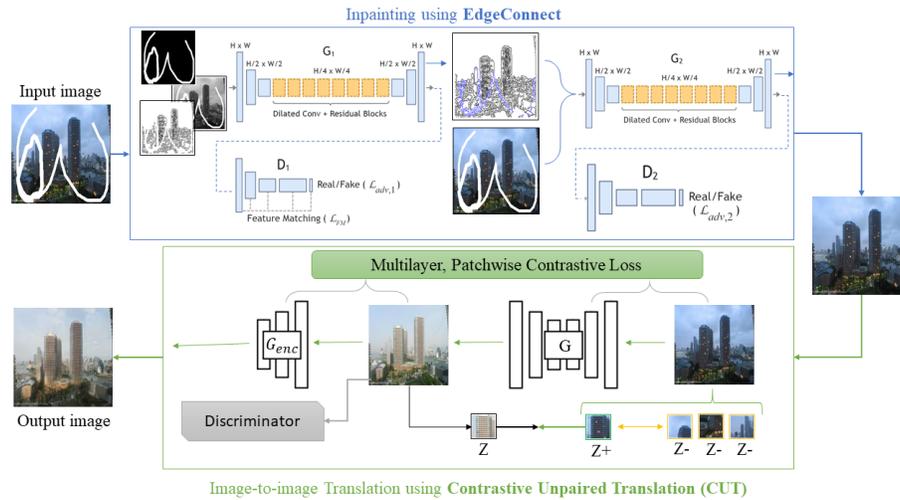


Fig. 1. Inpainting-First Model (M1), Inpainting module followed by an image-to-image translation module.

3.1 Problem description

The task of combining image inpainting and image-to-image translation can be formalized as follows. Here we introduce the problem in the context of day and nighttime images, that can be generalized to any arbitrary multi-modality image inpainting scenario. Given an incomplete input image in the nighttime mode N , with an arbitrary missing region, recover a corresponding complete daytime image D in which the missing region N_m has been completed with plausible content.

Human perception is perhaps the the most reliable tool for judging the perceived quality of a generated image. Ideally, the recovered image D should meet the following

criteria: It should contain meaningful coherent content with N so that human observer would accept it as a daytime image of N . It should have realistic-looking texture and color spectrum based on a daytime image. Finally, if the ground truth daytime image is available, the recovered image should ideally be similar to the ground truth based on suitable image similarity measures.

3.2 System overview

To address the described problem, we propose and investigate a two-stage model.

Inpainting first model The model involves first inpainting N using some inpainting method, e.g., the **EdgeConnect** [9] to obtain a complete nighttime image N_{EC} first. Once this is done, the new inpainted nighttime image without missing regions is translated into a realistic daytime image of the scene D using an image translation approach, e.g., the **CUT** [10] which is newer version of *CycleGAN* [15]. The inpainting-first model is illustrated in Figure 1. In the following Sections, we will also investigate the effect of switching the order of the inpainting and translation modules in the recovery process, i.e., considering a translate-first model.

Training There are two stages of training for the model.

In the **Inpainting-first** model, also referred to as M1, the *first training stage* involves training of the EdgeConnect model to inpaint a given N to generate N_{EC} . The EdgeConnect is trained in three stages: 1) training the edge completion model, 2) training the inpainting model and 3) training the joint model. To train the EdgeConnect model, incomplete grayscale image N_{gray} , edge map N_{edge} , and missing region mask N_m of training data are the inputs to the training at the edge model stage. The output would be the full edge map N_{edge} which contains hallucinated edges in N_m . During the training phase, the full edge map of nighttime ground truth $N_{gt_{edge}}$ computed using Canny edge detector is used to compare with N_{edge} to optimize the model. To train the inpainting model, predicted edge map N_{edge} and incomplete color image N are passed to generate N_{EC} . For the first stage of the model in this study the joint model of EdgeConnect was used [9].

In the *second stage* of the **Inpainting first** model, the inpainted nighttime image N_{EC} is translated into the corresponding daytime image using the CUT model to yield D_{CUT} . For training the CUT model in the *second stage*, N_{EC} of all training data is computed as the source domain and D_{gt} is considered as the target domain [10].

4 Experiments and Results

The proposed model is evaluated quantitatively and qualitatively in this Section. Dataset and setup are explained here as well.

4.1 Dataset

The *Transient Attributes dataset* [7] was used to create the **night2day** dataset. Each original sample is of size 256×512 containing two 256×256 images, i.e., one day image along with its corresponding night image. These images have been divided into 20, 110 night images and 20, 110 corresponding day images. The dataset was then split into training, validation, and test sets.

4.2 Experiments

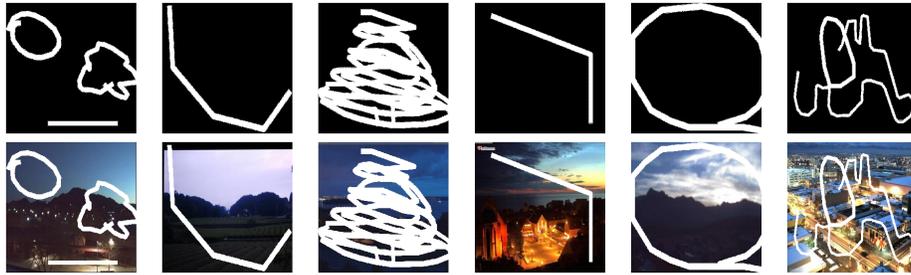


Fig. 2. Examples of masks from the QD-IMD data [4] applied to nighttime images.

All experiments were performed on Linux Ubuntu operating system and an NVIDIA GeForce GTX TITAN X GPU.

We used **Quick Draw Irregular Mask Dataset (QD-IMD)** [4] to apply random and irregular masks on night images. The QD-IMD research team believes that a combination of strokes drawn by the human hand is a good source of patterns for irregular masks [4]. This dataset contains 60,000 masks. The first row of Figure 2 shows some examples of masks taken from the QD-IMD.

To produce nighttime images with missing points, these masks are randomly applied to night images. Therefore, 20, 110 night-time images with irregular missing points for training, validation, and test phases are generated. More precisely, 2,011 or 10% for testing, 2,011 or 10% for validation, and 16,088 or 80% for training phases. Second row of Figure 2 illustrates examples of input nighttime images with missing points after applying the irregular masks. For each input image, the actual mask is provided to the model as well.

4.3 Model Evaluation

For testing the model, 2,011 nighttime images with missing points are used. Since the images are of different modalities, i.e, night and day and we have provided a two-step approach, different evaluation phases have been defined for a more meaningful quantitative evaluation of the results as explained below;

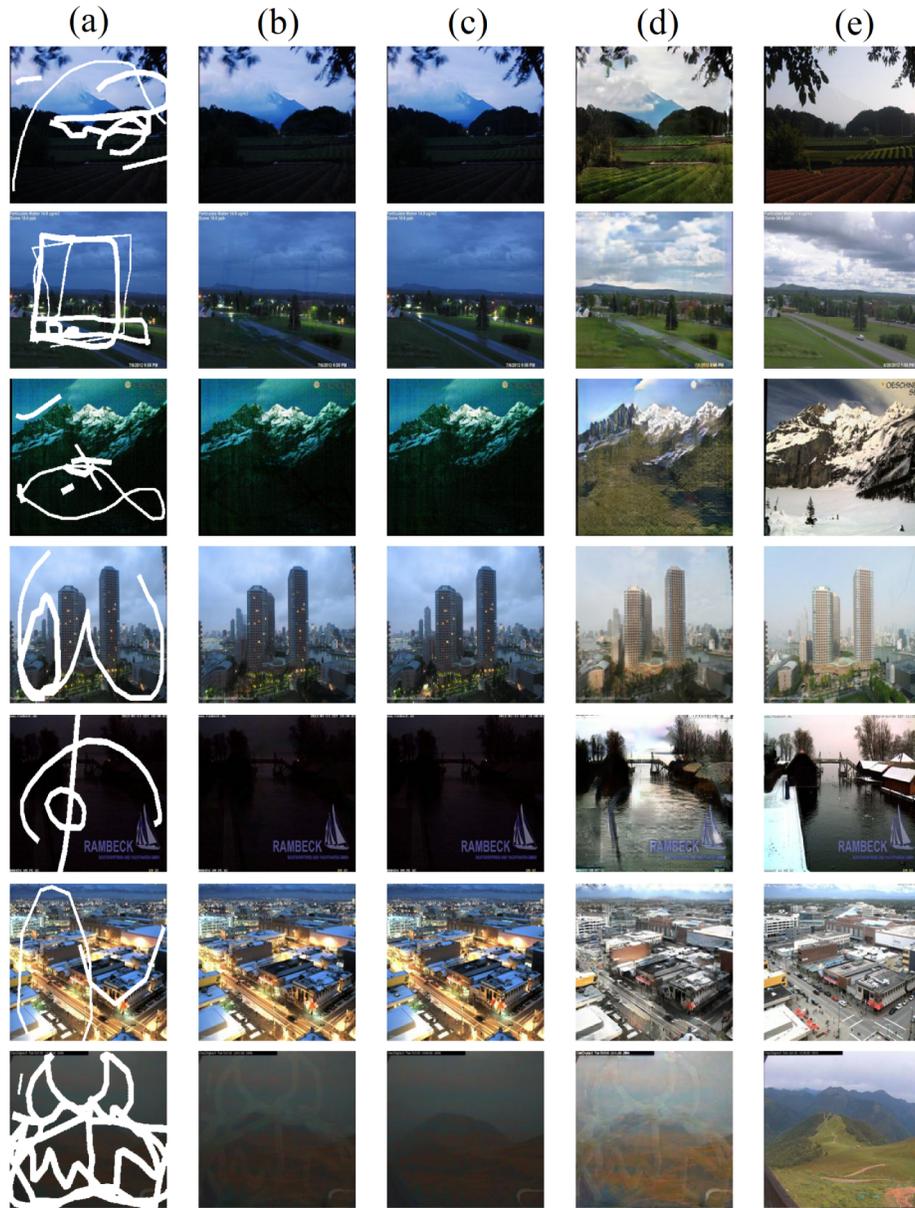


Fig. 3. Qualitative evaluation of model; from left to right; (a) Input image, (b) Inpainted night image, (c) Nighttime ground truth, (d) Generated output, (e) Daytime ground truth.

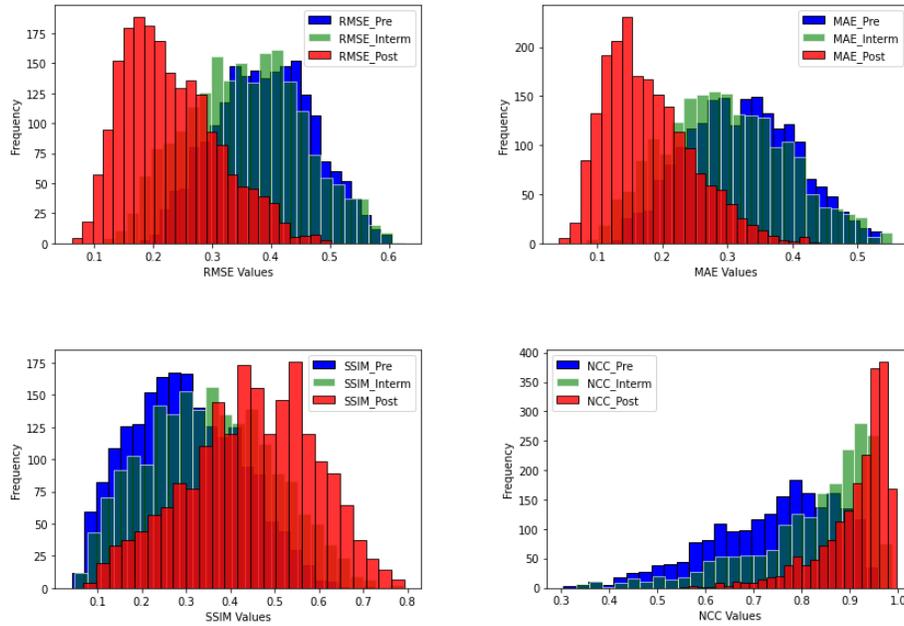


Fig. 4. Comparison of final results (Pre vs Intermediate vs Post) of the model using different measures.

- **Phase 1:** Comparing *the nighttime image with missing points (Input) vs daytime ground truth image,*
- **Phase 2:** Comparing *the inpainted night image after inpainting vs daytime ground truth image,*
- **Phase 3:** Comparing *the generated daytime image without missing points (output) vs daytime ground truth image.*

Root Mean Square Error (RMSE), Mean Absolute Error (MAE), Structural Similarity Index Measure (SSIM), Normalized Cross-Correlation (NCC), and Fréchet inception distance (FID) are calculated for each phase and histograms of these similarity measures over the test data are provided.

4.4 Results:

The total performance of the model can be evaluated by comparing Phase 1 with Phase 2 and Phase 3. We can refer to Phase 1 as “**Pre**” evaluation of similarities before applying the method. Results from Phase 2 are considered “**Intermediate**” after the first step, e.g, after inpainting. Finally phase 3 which compares the output with the ground truth, could be referred to as “**Post**”. For each measure, we compare the results of these phases. Figure 4 illustrates these comparisons.

Figure 4 shows that the *RMSE* (top-left) and the *MAE* (top-right) are shifted to the left from the “*Pre*” state to the “*Post*” state. Lower value of these measures for final

results in *Phase 3* indicates better performance. Similarly, the *Phase 3* or the “*Post*” state provides higher values for *SSIM* (bottom-left), and the *NCC* (bottom-right) in comparison with the “*Intermediate*” and the “*Pre*” states. Table 1 presents the *mean* and the *standard deviation* of these measures for **Pre**, **Intermediate**, and **Post**. The FID metric for each phase is computed as well. Based on all these measures, table 1 suggests that final results from the model are more consistent with ground truth.

Table 1. Comparing final results (**Phase 1-A vs Phase 2-A vs Phase 3**) of the model using different measures.

	Pre	Intermediate	Post
<i>RMSE</i> (↓)	0.39±0.08	0.36±0.10	0.23±0.08
<i>MAE</i> (↓)	0.16±0.06	0.14±0.07	0.06±0.04
<i>SSIM</i> (↑)	0.30±0.12	0.35±0.14	0.45 ±0.14
<i>NCC</i> (↑)	0.74±0.13	0.82±0.13	0.91±0.07
<i>FID</i> (↓)	271.01	80.28	56.77

In order to **evaluate the results of the model qualitatively**, Figure 3 provides some sample test images at different stages. From left to right, each row contains a) Input i.e. the nighttime image with missing points, b) The inpainted nighttime image without missing points (after inpainting), c) Nighttime ground truth image, d) Output i.e. the generated daytime image without missing points, and e) Daytime ground truth image. It can be observed that when enough details exist in the input image, the model produces visually plausible results.

4.5 Evaluation of a switched recovery order model

A question may be posed whether the model produces similar results if we change the order of the inpainting and translation modules, i.e., if the translation module was applied first followed by the inpainting. In order to answer this question, the so called **M2: Translation first model** was implemented and its results are presented in Figure 5.

The obtained results of **Phase 3** for the *inpainting first* and the *translation first* models allow us to compare their performances. In both models, phase 3 is comparing the generated daytime image (*Output*) vs daytime ground truth image. This comparison using the *RMSE* (top-left), the *MAE* (top-right), the *SSIM* (bottom-left), and the *NCC* (bottom-right) is illustrated in Figure 6. The *RMSE* histogram and *MAE* histogram for the *inpainting first model* are shifted to the *left*, indicating superior results. Similarly, the *SSIM*, and *NCC* histograms are shifted towards the *right*, which reaffirms better results based on these measures for the *inpainting first model*. Table 2 summarizes these comparisons.

A qualitative evaluation of the images generated using these two models, which are illustrated in Figure 3 and Figure 5 as well as quantitative evaluations above confirm that the **inpainting first model** provides a better solution.

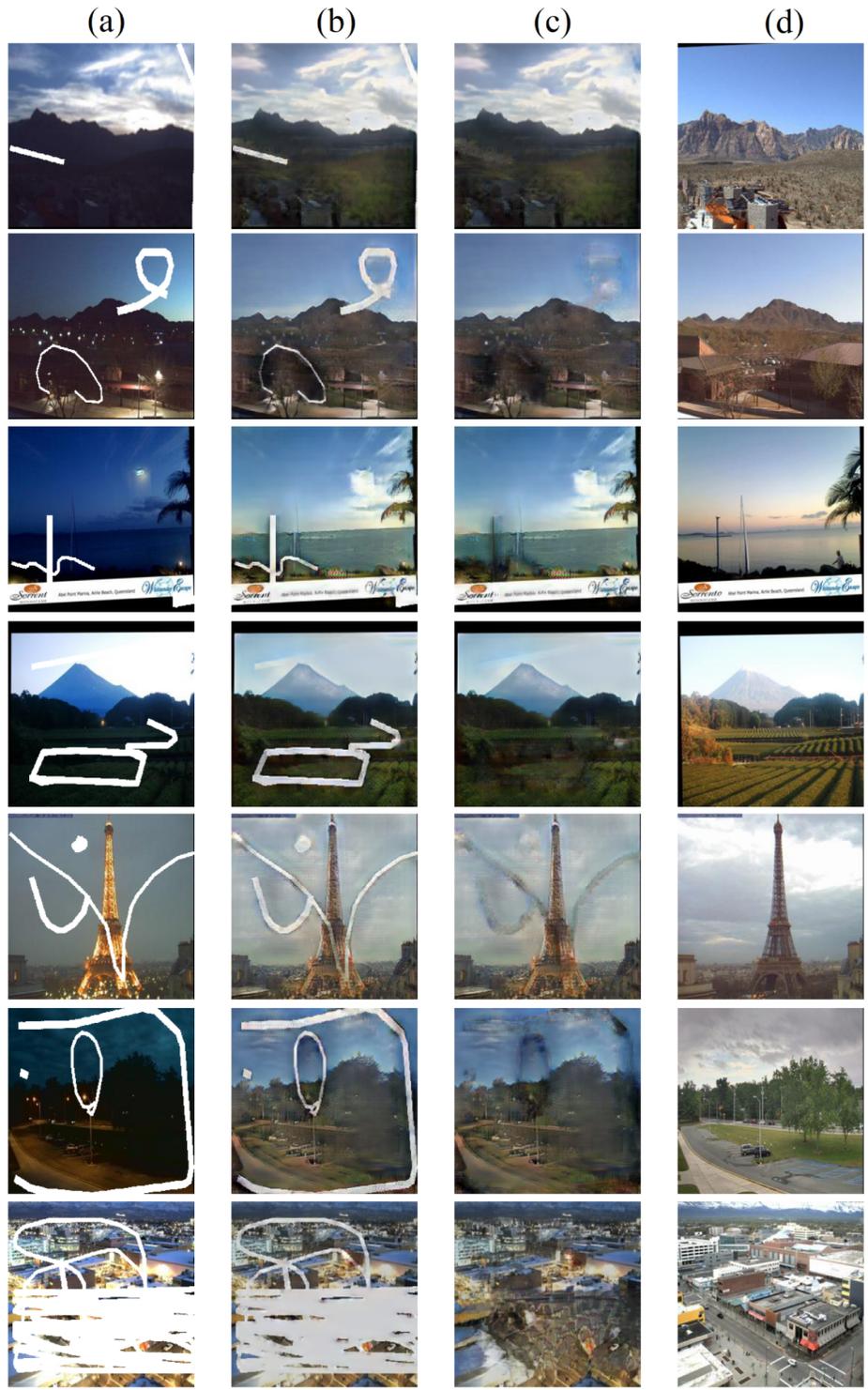


Fig. 5. Qualitative evaluation of the translation first model; from left to right; (a) Input image, (b) Generated day image with missing points, (c) Output image, (d) Day ground truth.

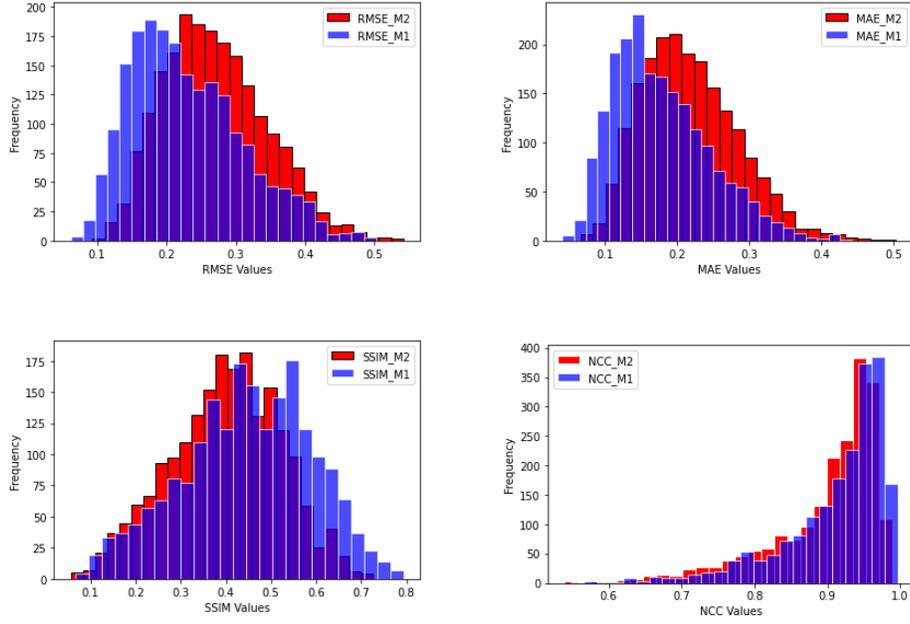


Fig. 6. Comparison of final results of the *M1: Inpainting first* with the *M2: Translation first* models using different measures.

Table 2. Comparing final results of the **M1: Inpainting first model** with **M2: Translation first model** using different measures.

	M2	M1
$RMSE(\downarrow)$	0.27 ± 0.07	0.23 ± 0.08
$MAE(\downarrow)$	0.08 ± 0.04	0.06 ± 0.04
$SSIM(\uparrow)$	0.40 ± 0.12	0.45 ± 0.14
$NCC(\uparrow)$	0.90 ± 0.07	0.91 ± 0.07
$FID(\downarrow)$	108.62	56.77

5 Ablation Study

Mask line thickness increases the number of missing pixels as well as the distance between those pixels and known neighboring pixels along the edge of the mask. Those neighboring pixels provide crucial information for hallucinating the missing points and translating a complete night image into a coherent day image.

In order to study the effects of the mask on the model’s performance, we choose an arbitrary night image from test data and apply a mask with a simple thin line from the *QD-IMD dataset* [4] dataset to it. At different steps, the line is dilated to increase its thickness. Therefore, various masks are generated, see Figure 7. We apply these masks to the sample image and test the *model* against it. For each iteration, Figure 7 illustrates the input images and results of the model.

Early steps of the model successfully recover the missing area and convert it into a plausible daytime image. After some steps, the recovered area is blurry and there is not much information to be used at the translation stage. Increasing the thickness of the line reduces the performance of both inpainting and translating as expected and the translation module is barely able to generate a daytime image.

6 Discussions and Future Work

To tackle the issue of merging the image-to-image translation task and the image-inpainting task, we proposed the *inpainting first* model by concatenating two existing modules and rigorously evaluated its performance after training. The model produces plausible images based on different image quality assessments both quantitatively and qualitatively. Furthermore, the ablation study indicated that the thickness input masks can drastically degrade the quality of the generated output.

An end-to-end generative network that performs both tasks simultaneously might provide more accurate results. In this manuscript, we focused on the night to day image inpainting and provided promising results.

A deep network that is able to generate a realistic image in a modality from an image with missing points in another given modality can be used in many image enhancement tasks including medical imaging applications.

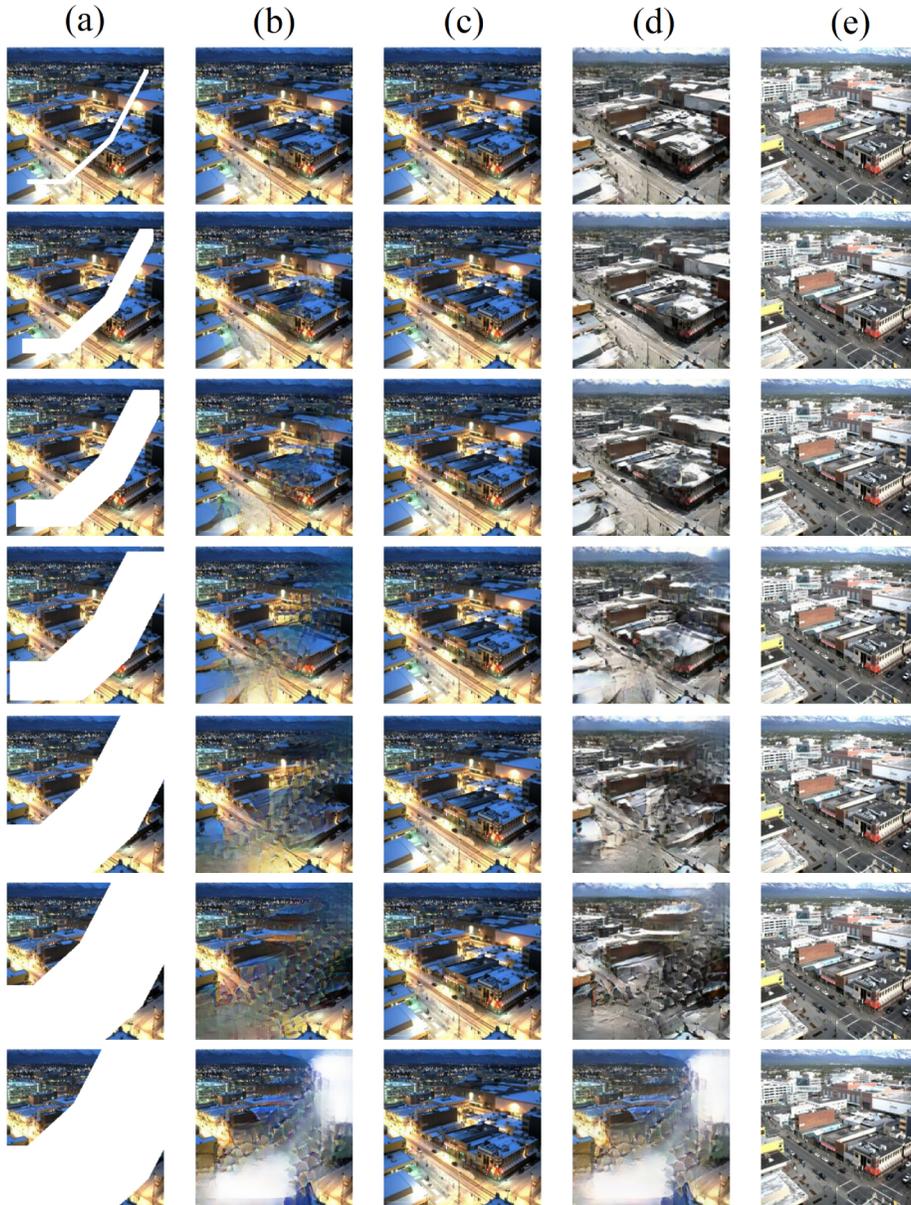


Fig. 7. Ablation study of the model; Analyzing mask dilation in input. From left to right; (a) Input image, (b) generated inpainted night image, (c) Night ground truth, (d) Output image, (e) Day ground truth.

Acknowledgements

This work was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC).

References

1. Baek, K., Choi, Y., Uh, Y., Yoo, J., Shim, H.: Rethinking the truly unsupervised image-to-image translation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14154–14163 (2021)
2. Chen, Y., Zhang, H., Liu, L., Chen, X., Zhang, Q., Yang, K., Xia, R., Xie, J.: Research on image inpainting algorithm of improved gan based on two-discriminations networks. *Applied Intelligence* **51**, 1–15 (06 2021). <https://doi.org/10.1007/s10489-020-01971-2>
3. Emami, H., Aliabadi, M.M., Dong, M., Chinnam, R.B.: Spa-gan: Spatial attention gan for image-to-image translation. *IEEE Transactions on Multimedia* **23**, 391–401 (2020)
4. Isakov, K.: Semi-parametric image inpainting (2018)
5. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks (2018)
6. Kong, L., Lian, C., Huang, D., Hu, Y., Zhou, Q., et al.: Breaking the dilemma of medical image-to-image translation. *Advances in Neural Information Processing Systems* **34** (2021)
7. Laffont, P.Y., Ren, Z., Tao, X., Qian, C., Hays, J.: Transient attributes for high-level understanding and editing of outdoor scenes. *ACM Transactions on Graphics (proceedings of SIGGRAPH)* **33**(4) (2014)
8. Li, J., Wang, N., Zhang, L., Du, B., Tao, D.: Recurrent feature reasoning for image inpainting (2020)
9. Nazeri, K., Ng, E., Joseph, T., Qureshi, F., Ebrahimi, M.: Edgeconnect: Structure guided image inpainting using edge prediction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (2019)
10. Park, T., Efros, A.A., Zhang, R., Zhu, J.Y.: Contrastive learning for unpaired image-to-image translation. In: European Conference on Computer Vision (2020)
11. Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., Cohen-Or, D.: Encoding in style: a stylegan encoder for image-to-image translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2287–2296 (2021)
12. Yi, Z., Tang, Q., Azizi, S., Jang, D., Xu, Z.: Contextual residual aggregation for ultra high-resolution image inpainting (2020)
13. Zhao, L., Mo, Q., Lin, S., Wang, Z., Zuo, Z., Chen, H., Xing, W., Lu, D.: Uctgan: Diverse image inpainting based on unsupervised cross-space translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
14. Zhou, T., Ding, C., Lin, S., Wang, X., Tao, D.: Learning oracle attention for high-fidelity face completion (2020)
15. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: IEEE International Conference on Computer Vision (ICCV) (2017)