
Spatially-Adaptive Pixelwise Networks for Fast Image Translation

Melaku Getahun¹ Dawit Sefiw¹ Mekan Hojaye¹ Oluwafemi Adejumobi¹

Abstract

This project explores the use of a new generator architecture for an efficient high resolution image to image translation tasks. We replicated the results of ASAPNet(A Spatially Adaptive Pixelwise Network) paper by training the generator and discriminator networks from scratch. We evaluated the performance of the model on the cityscapes dataset using the inference time, Frechet Inception Distance and mean Intersection-over-union metrics and compared our results to those reported by the authors. We also investigate the possibility of using the same generator network for denoising and deblurring tasks.

Github repo: <https://github.com/MelakuNe/Machine-Learning-Project-2023>

Presentation file: <https://github.com/MelakuNe/Machine-Learning-Project-2023/blob/main/ASAPNet.pdf>

1. Introduction

Image translation from one domain to other has been a focus of research in machine learning and computer vision for years. Recent advancements in the area have resulted in the development of conditional Generative Adversarial Networks (GANs) that can learn a direct mapping between different image domains. However, these models are computationally expensive, especially when operating on high-resolution images. To address this issue, a novel architecture called ASAPNet (A Spatially Adaptive Pixelwise Network) (Shaham et al., 2021) that uses a new generator designed for fast image-to-image translation is proposed. The generator operates pixelwise, processing each pixel solely using a structurally simplified Multi-Layer Perceptron (MLP) that is adaptive to the input image. Generating high-quality images

at very short execution time is proven to be possible by this architecture. This model is evaluated in various image to image translation domains such as segmentation maps and indoor images and comparison is made with state of the art models. In this work, we would like to make sure that the achievement and result mentioned in the paper is fully reproducible and put direction of adapting this architecture for other image restoration tasks such as image deblurring. There exists state-of-the-art deblurring architectures like DeblurGAN that use generative architecture networks to produce a clear and high quality image that accurately represent the original scenes. However, inference time of these models increase with visual quality. As we can leverage computational gain from the lower resolution processing technique, proving the applicability of ASAPNet's generator network for such a task could be a huge advantage.

We evaluate the model on the cityscape dataset to generate realistic city scenery from semantic segmentation maps. In addition to that, the expandability of the proposed generator network is tested for image denoising and deblurring tasks. In general, while achieving an equally perceptible image quality, the newly proposed architecture proves to generate images in considerably lower inference time compared to existing state-of-the-art tools.

2. Related work

Image to Image Translation:

Conditional Generative Adversarial Networks (GAN) have brought tremendous new things in generative models including image to image translation. However, the translation of images from one domain using this model has several limitations including not being able to work with high-resolution image and more realistic image generation is also a challenge. (Wang et al., 2018) proposed a generator network composed of two sub-networks called global generator network and local enhance network. The former one is trained with low resolution images then the result is appended to the later network and it is trained together with high resolution images. On the other hand multi-scale discriminator with 3 sub-networks is proposed that operates at different image scales to improve the limitation of discriminator challenge of differentiating high resolution images and synthesized images. Mostly, after convolutional batch normalization

¹Skolkovo Institute of Science and Technology, Moscow, Russia. Correspondence to: Melaku Getahun <melaku.getahun@skoltech.ru>, Dawit Sefiw <dawit.sefiw@skoltech.ru>, Mekan Hojaye <mekan.hojaye@skoltech.ru>, Oluwafemi Adejumobi <oluwafemi.adejumobi@skoltech.ru>.

layer is used and it has lack of representing the semantic information well in fact it degrades this important feature of an image. (Park et al., 2019) proposed a conditional normalization layer consisting of Batch Normalization and Adaptive Instance Normalization. When performing image to image translation task the semantic label is projected to an embedding space and convolution operation is applied to produce γ and β which has a spatial information in it unlike the common batch normalization layer with γ and β just a vector.

Denoising:

Image restoration has been also one area of interest in deep convolutional network research era, and one of the main area is removing noise from an image given noisy data. Convolutional neural network have been used widely to study the noise natural distribution in a given data, however there is always a bottleneck that these model are trained with a specific noise level which makes it not robust enough to handle not predefined noise level data. (Zhang et al., 2017) proposed a denoising CNN model, where VGG architecture is adopted in such a way that it is suitable for denoising purpose. To improve the denoising performance and reduce the training time residual learning is incorporated with batch normalization. The model is tested with any different dataset and has achieved better result than previous work, in addition to that since the model is suitable for parallel computation it has shown good inference time for an image of 512x512 with unknown noise level it took 60ms.

3. Spatially Adaptive Pixelwise Networks

This section discusses the design of Spatially Adaptive Pixelwise Networks. ASAPNet is designed to be input-adaptive, spatially-varying, and high-frequency detail synthesizing model. The design involves using a Multi-Layer Perceptron (MLP) with ReLU activations that map:

$$f_p(x_p, p) = f(x_p, p; \varphi_p) = y_p$$

where f denotes the MLP architecture — shared by all pixelwise functions — and φ_p are the (spatially-varying) weights and biases for the MLPs.

Output is modelled as a spatially-varying pointwise non-linear transformation of the high-resolution input, where each pixel is independent of the others. To preserve spatial information, each pixelwise function takes as input the pixel’s coordinates, in addition to its color value. The pixelwise functions are parameterized with spatially-varying parameters that are conditioned on the input image. The parameter vectors are predicted by a convolutional network that operates on a much lower resolution image. The model also uses positional encodings to synthesize high-frequency details. Authors have found that the MLPs with 5 layers and

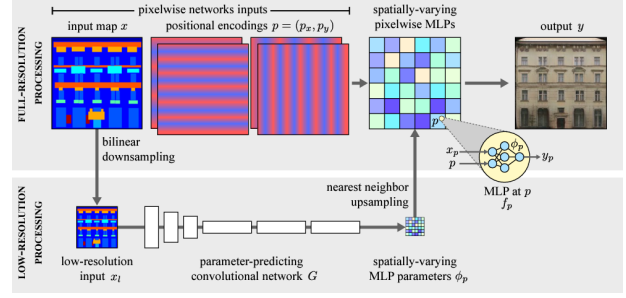


Figure 1. Architecture overview. The model first processes the input at very low-resolution, x_l , to produce a tensor of weights and biases φ_p . These are upsampled back to full-resolution, where they parameterize pixelwise, spatially-varying MLPs f_p that compute the final output y from the high-resolution input x . In practice, we found nearest neighbor upsampling to be sufficient, which means that the MLP parameters are in fact piecewise constant.

64 channels per layer provide fast execution and good visual quality. The training procedure follows best practices for the optimization procedure, including training the generator model adversarially with a multi-scale patch discriminator, using an adversarial hinge-loss, a perceptual loss, and a discriminator feature matching loss.

4. Experiments

We validate the performance of the proposed method for image segmentation task by making a comparison with the corresponding baseline model using inference time and performance metrics. FID(lower is better) which measures the similarity between the distributions of real and generated images, and mean intersection over union(higher is better) which is a measure of overlap between predicted and ground truth segmentations, are the metrics’ used for the segmentation task. We also experiment on the possible use of the model’s generator network for denoising and deblurring task.

Baselines. In image segmentation, the performance of ASAPNet model is compared with pix2pixHD, a model which takes a residual learning approach that enables high-resolution image translation. For denoising, only the generator network is trained from scratch and DnCNN is used as a baseline. We also combined our generator network with DeblurGAN’s discriminator to compare performance metrics against baseline DeblurGAN.

Datasets. For the segmentation task, We used the Cityscapes dataset, which contains images of urban scenes and their semantic label maps, split into 3000/500 training/validation images. We construct 256×512 versions of this dataset. Denoising experiment is performed using the Berkeley Segmentation Dataset(BSD) splitted into 432/68 training/validation images. For the training data, a syn-

thetic noise with different standard deviations of noise level ranging from 10 to 55 is added and validation data is also influenced by a synthetic noise with standard deviation of noise level 25. To test the model for deblurring task, 1151 pair of blurred and sharp images from the GOPRO test dataset is used.

Speed: The inference time of all models is benchmarked on NVIDIA GeForce GTX 1080ti. The ASAPNet paper generalizes that their model’s inference time is $2 - 6 \times$ faster than pix2pixHD. For the experiment we did with cityscapes dataset of size 256×512 , we observed that ASAPNet takes 18ms to generate an image while it takes 51ms for pix2pixHD, which approximates to a $3 \times$ faster inference time.

Qualitative comparisons: Examples of translated images are shown in figure 2. It can be deduced that images generated using ASAPNet have a visual quality at least comparable to the baseline. Shown in figure 6 is outputs of ASAPNet’s generator and the DnCNN baseline for a denoising task. Although current result is not satisfactory, the model has shown a glimpse for its applicability in this task.

Quantitative comparisons: Even though a specific numerical result of FID and mean intersection over union for the specified dataset and size is not mentioned by the authors, results they provided in graphical form is demonstrated for comparison with ours in figure 4 and 5. The authors visualized that ASAPNet achieves a better FID score than pix2pixHD, while our experiment shows an FID score of 70.6 and 65 for both models, respectively. The experiments found that the average intersection over union was 30 for the ASAPNet model and 29 for the pix2pixHD model. However, the authors claimed pix2pixHD achieved better results. Although there is a minor discrepancy on performance metrics reported by the authors and our experiment, It can be generalized that the results are on same scale and ASAPNet shows a performance comparable to existing tools while outperforming in inference time. For denoising and deblurring tasks, our model has achieved a Peak to Signal Ratio of 16 and 19, respectively, while the same metrics is 28.55 and 24.6 for their corresponding baselines.



Figure 2. Qualitative comparison of both models

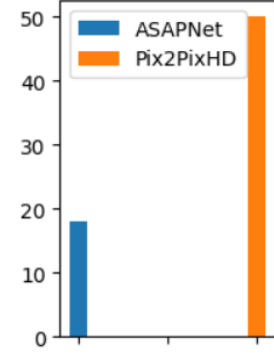


Figure 3. Inference Time

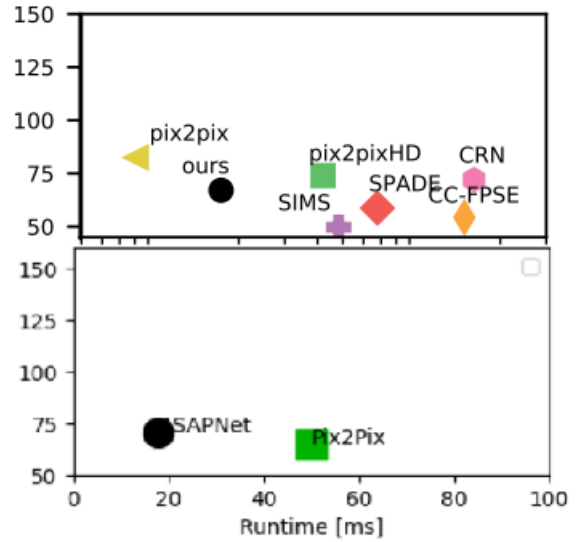


Figure 4. FID comparison of paper and replication work



Figure 6. Image outputs of denoising operation by the baseline DnCNN(left) and our generator network(right).

5. Conclusion

This project replicated a novel neural architecture for generative adversarial image-to-image translation problems that is faster than previous works but with same high level visual

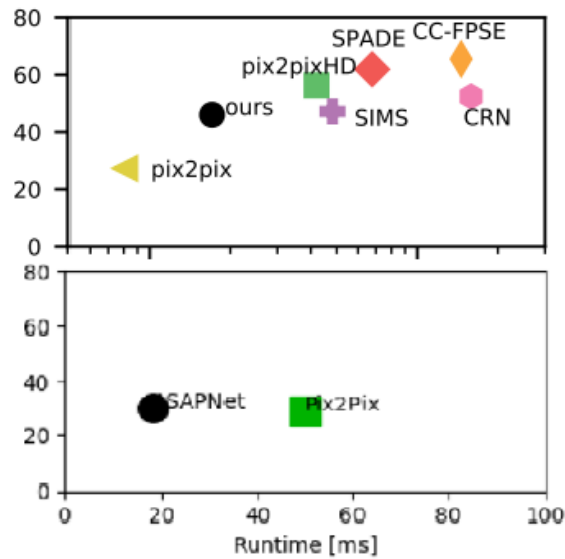


Figure 5. mean IoU of paper and replication work

quality. In addition to this, We showed our experiment on the subsequent use of the generator architecture for other image to image translation tasks such as denoising and deblurring. Although our experiment fell a bit short of proving the subsequent use of the generator network for such tasks, we strongly believe advanced development of the network can achieve an outstanding performance with an already granted advantage of inference time and computational cost.

References

- Park, T., Liu, M.-Y., Wang, T.-C., and Zhu, J.-Y. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2337–2346, 2019.
- Shaham, T. R., Gharbi, M., Zhang, R., Shechtman, E., and Michaeli, T. Spatially-adaptive pixelwise networks for fast image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14882–14891, 2021.
- Wang, T.-C., Liu, M.-Y., Zhu, J.-Y., Tao, A., Kautz, J., and Catanzaro, B. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8798–8807, 2018.
- Zhang, K., Zuo, W., Chen, Y., Meng, D., and Zhang, L. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155, 2017.

A. Team member’s contributions

Explicitly stated contributions of each team member to the final project.

Melaku Getahun (25% of work)

- Initially put forward the idea
- Training the replication part
- Experimenting with model parameters
- Preparing the GitHub Repo
- Preparing the Section 2 of this report
- Creating presentation

Dawit Sefiw (25% of work)

- Training the replication part
- Experimenting with model parameters
- Preparing the GitHub Repo
- Preparing the Section 3 of this report

Mekan Hojayev (25% of work)

- Training the research part
- Redacting presentation
- Preparing the GitHub Repo
- Preparing the Sections 4 of this report

Oluwafemi Adejumobi (25% of work)

- Training the research part
- Experimenting with model parameters
- Preparing the GitHub Repo
- Preparing the Section 5 of this report

B. Reproducibility checklist

Answer the questions of following reproducibility checklist.
If necessary, you may leave a comment.

1. A ready code was used in this project, e.g. for replication project the code from the corresponding paper was used.

☒ Yes.
☐ No.
☐ Not applicable.

General comment: If the answer is **yes**, students must explicitly clarify to which extent (e.g. which percentage of your code did you write on your own?) and which code was used.

Students' comment: None

2. A clear description of the mathematical setting, algorithm, and/or model is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

3. A link to a downloadable source code, with specification of all dependencies, including external libraries is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

4. A complete description of the data collection process, including sample size, is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

5. A link to a downloadable version of the dataset or simulation environment is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

6. An explanation of any data that were excluded, description of any pre-processing step are included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

7. An explanation of how samples were allocated for training, validation and testing is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

8. The range of hyper-parameters considered, method to select the best hyper-parameter configuration, and specification of all hyper-parameters used to generate results are included in the report.

☐ Yes.
☐ No.
☒ Not applicable.

Students' comment: because we were asked to use the paper parameters, so no selection of hyperparameters.

9. The exact number of evaluation runs is included.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

10. A description of how experiments have been conducted is included.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

11. A clear definition of the specific measure or statistics used to report results is included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

12. Clearly defined error bars are included in the report.

☒ Yes.
☐ No.
☐ Not applicable.

Students' comment: None

13. A description of the computing infrastructure used is included in the report.

☒ Yes.

☐ No.

☐ Not applicable.

Students' comment: None