

Conditional Transfer with Dense Residual Attention: Synthesizing traffic signs from street-view imagery

*Clint Sebastian[†], *Ries Uittenbogaard[§], Julien Vijverberg[‡], Bas Boom[‡], and Peter H.N. de With[†]

[§] Department of Mechanical Engineering, Delft University of Technology, Delft, The Netherlands

[†]Department of Electrical Engineering, Eindhoven University of Technology, Eindhoven, The Netherlands

[‡] Cyclomedia Technology B.V, Zaltbommel, The Netherlands

Email: m.c.uittenbogaard@student.tudelft.nl, {c.sebastian, p.h.n.de.with}@tue.nl,

{jvijverberg, bboom}@cyclomedia.com

* denotes equal contribution

Abstract—Object detection and classification of traffic signs in street-view imagery is an essential element for asset management, map making and autonomous driving. However, some traffic signs occur rarely and consequently, they are difficult to recognize automatically. To improve the detection and classification rates, we propose to generate images of traffic signs, which are then used to train a detector/classifier. In this research, we present an end-to-end framework that generates a realistic image of a traffic sign from a given image of a traffic sign and a pictogram of the target class. We propose a residual attention mechanism with dense concatenation called Dense Residual Attention, that preserves the background information while transferring the object information. We also propose to utilize multi-scale discriminators, so that the smaller scales of the output guide the higher resolution output. We have performed detection and classification tests across a large number of traffic sign classes, by training the detector using the combination of real and generated data. The newly trained model reduces the number of false positives by 1.2 - 1.5% at 99% recall in the detection tests and an absolute improvement of 4.65% (top-1 accuracy) in the classification tests.

I. INTRODUCTION

Detection and classification of traffic signs in street-view imagery is vital for public object maintenance, map making and autonomous driving. This is particularly challenging when certain classes or categories of objects are scarce. Challenging cases occur in the automated detection and classification of traffic signs at a country-wide level on high-resolution street-view imagery. Manual efforts to find traffic signs in millions of high-resolution images is cumbersome and the detection/classification algorithm fails if it is not trained with the class-specific data or when traffic signs rarely occur. A possible approach to alleviate this problem is to generate realistic data using generative modeling and expand the training sets that have low amount of data or low recognition scores. However, generation of photo-realistic samples of traffic signs is difficult due to large variations in pose, lighting conditions and varying background.

Recent developments in deep learning can be applied to modeling of image-to-image translation problems. Generative Adversarial Network (GAN) is a class of deep learning algorithms that is used for generative modeling [1]. GANs formulate the generative modeling problem as zero-sum game

between two networks. A GAN consists of a generator network that produces samples from a given input sample or noise and a discriminator network that tries to distinguish if the generated sample is from the real or fake data distribution. Although Convolutional Neural Networks (CNNs) may be used to perform image-to-image translations, many of them apply a stochastic approximation to minimize an objective function and require paired data [2][3]. Alternatively, GANs try to achieve Nash equilibrium by generating a distribution that is close to the empirical one.

Conditional variants of GANs have recently produced impressive results in image-to-image translation applications [4][5][6][7][8]. A conditional GAN tries to map a Domain A to another Domain B , instead of using a noise vector as input. It learns to translate an underlying relationship between the Domains A and B without explicitly pairing inputs and outputs. To have a better consistency for the mapping, most Conditional GANs also reconstruct the output back to the input [4][7]. Apart from the mapping from Domain A to B , Domain B is also reconstructed back to A . The loss from the inverse mapping is also added to the objective function as a regularizer during training. Image analogy problems can be modeled using a Conditional GAN by providing an auxiliary piece of information [9]. The auxiliary information could be text, image or other data modalities. For example, an input image of a traffic sign can be paired with a pictogram to obtain an output image of a new traffic sign.

In this paper, we thus explore a conditional GAN for traffic sign generation, while using an auxiliary information of a pictogram. Specifically, we address the problem of retaining the original background while altering only the object information. In more detail, we propose a conditional GAN with Dense Residual Attention. To further improve the texturing and details of the generated traffic sign, we use multi-scale discriminators. We reduce the dependency of the mask generation process that is used in recent work [9]. Our method produces perceptually acceptable results without the implicit generation of a mask. However, we are able to obtain better results with a weak supervision on the traffic signs that have a complex pose. Finally, we improve classification and detection rates of rare traffic signs in high-resolution street-view imagery.

II. RELATED WORK

Unsupervised generative modeling using GANs has achieved state-of-the-art results in recent years. Many works have produced perceptually realistic images for problems such as image-to-image translation [4][5][6][7][8], image inpainting [10], super-resolution [3][11][12] and conditional analogy [9]. Deep Convolution GANs (DCGANs) introduced the convolutional architecture GANs that improve visual quality of generated images [13]. Recently, Wasserstein GANs (WGANs) introduced an objective function that improves model stability and provides a meaningful measure of convergence [14]. However, they require the discriminator (also known as critic) to lie within the space of 1-Lipschitz functions. To enforce the Lipschitz constraint, the weights are clipped to a compact space. To circumvent weight clipping, improved WGAN have proposed to add a gradient penalty term to the WGAN training objective. This gradient penalty term is the product of the penalty coefficient term and gradient norm of the critic’s output. This approach results in a better performance and stability [15].

Recently, GANs have gained popularity in image-to-image translation problems. GANs have been applied in a conditional setting to generate text, labels and images. Image-conditional models such as “pix2pix”, use a Conditional Adversarial Network to learn paired mappings between domains [6]. CycleGAN performs an unpaired domain transfer by adding a cycle consistency loss to the training objective [7]. Similarly, DiscoGAN also proposes an unpaired domain transfer where two reconstruction losses are added to the objective function [4]. Conditional Analogy GAN (CAGAN) propose to swap clothing articles on people [9]. Given a human model and a given fashion article, it produces a human model wearing the fashion article. The generator of the CAGAN architecture generates both the image and an implicit segmentation mask. The final output image is a convex combination of both the generated image and input image. However, in scenarios where the background is complex, the generation of an implicit mask becomes a challenging task. This is addressed in [9].

Attention mechanisms are necessary for proper guidance of feature propagation. Many efforts have been made towards incorporating an attention mechanism in deep neural networks [16][17][18]. Attention mechanisms have been widely adopted in Recurrent Neural Networks (RNNs) and Long Short Term Memory (LSTM) [19] for modeling sequential tasks. Residual Attention Networks propose an attention mechanism by stacking multiple attention modules for various classification tasks [20][21]. Each attention module consist of a mask and a trunk branch. The trunk branch performs feature processing using residual units and the mask branch uses a bottom-up top-down structure to guide the feature learning. Densely Connected Convolutional Network (DenseNet) proposes an architecture where every layer is connected to all higher layers [22]. This type of connectivity enables reusable features and provides high parameter efficiency, as low-dimensional feature maps are reused recursively. Concurrent

to our work, a multi-scale discriminator approach has been explored in [8] for image synthesis from semantic labels. In our method, we use a multi-scale discriminator to learn the finer details of the pictogram and the global features such as the pose and lighting condition of the traffic sign. The methodology is described in the following section.

III. METHOD

The task of transferring a pictogram to a given traffic sign can be formulated as follows. Given an image x_i of a traffic sign and pictogram p of the target class, the generator network G tries to produce a traffic sign image of the target class x_i^p . The discriminator network D distinguishes between x_i^p and x_j , where x_j is sampled from the real data distribution.

For training, we use the improved WGAN objective with cycle consistency loss at multiple scales [15][7]. The final training objective \mathcal{L}^s at a given scale s is expressed as:

$$\mathcal{L}^s = \min_G \max_{D^s} \mathcal{L}_{\text{WGAN-GP}}^s(G, D^s) + \mathcal{L}_{\text{cyc}}^s(G), \quad (1)$$

where $\mathcal{L}_{\text{WGAN-GP}}(G, D)$ is the WGAN loss function $\mathcal{L}_{\text{WGAN}}$ with gradient penalty and $\mathcal{L}_{\text{cyc}}(G)$ is the cycle loss. The training objective $\mathcal{L}_{\text{WGAN}}$ is expressed as the difference of the expected values of the fake and real outputs from the discriminator D . The WGAN adversarial loss with gradient penalty for our problem now becomes:

$$\begin{aligned} \mathcal{L}_{\text{WGAN-GP}}^s(G, D^s) = & \mathbb{E}_{x_i^p \sim \mathbb{P}_g} [D^s(x_i^p)] - \mathbb{E}_{x_j \sim \mathbb{P}_r} [D^s(x_j)] \\ & + \lambda \mathbb{E}_{\hat{x}_i \sim \mathbb{P}_{\hat{x}_i}} [(\|\nabla_{\hat{x}_i} D^s(\hat{x}_i)\|_2 - 1)^2]. \end{aligned} \quad (2)$$

Here, the output from the generator $G(x_i|p) \approx x_i^p$ and \hat{x}_i is a sampling from the distribution $\mathbb{P}_{\hat{x}_i}$. The sample \hat{x}_i is an interpolation between a pair of samples from the generated and real distributions \mathbb{P}_g and \mathbb{P}_r . The gradient penalty coefficient term λ is set to 10 for all the experiments in this research, which was adequate for our work and in accordance with [15]. To push the norm of the gradient towards unity, the gradient penalty is applied. For the cycle consistency loss, we use the same generator network to map multiple classes in a given category. Hence, cycle loss for our objective is defined as:

$$\mathcal{L}_{\text{cyc}}^s(G) = \mathbb{E}_{x_i^{p_a} \sim \mathbb{P}_g} \|x_i^{p_a} - G(G(x_i^{p_a}|p_b)|p_a)\|_2, \quad (3)$$

where p_a and p_b are pictograms of different traffic signs. The samples $x_i^{p_a}$ and $x_i^{p_b}$ are traffic signs conditioned on p_a and p_b . Note that $G(x_i^{p_a}|p_b) \approx x_i^{p_b}$ and $G(x_i^{p_b}|p_a) \approx x_i^{p_a}$. By transitive relation, $G(G(x_i^{p_a}|p_b)|p_a) \approx x_i^{p_a}$. Therefore we also minimize the L_2 distance between the input and reconstructed examples. An overview of the proposed network is presented in Figure 1. The generator has an encoder-decoder structure with a residual attention mechanism and dense connectivity at each scale. To discriminate between the real and generated samples at multiple scales, a discriminator is applied at each scale. The specifics of our contributions are addressed in the following subsections.

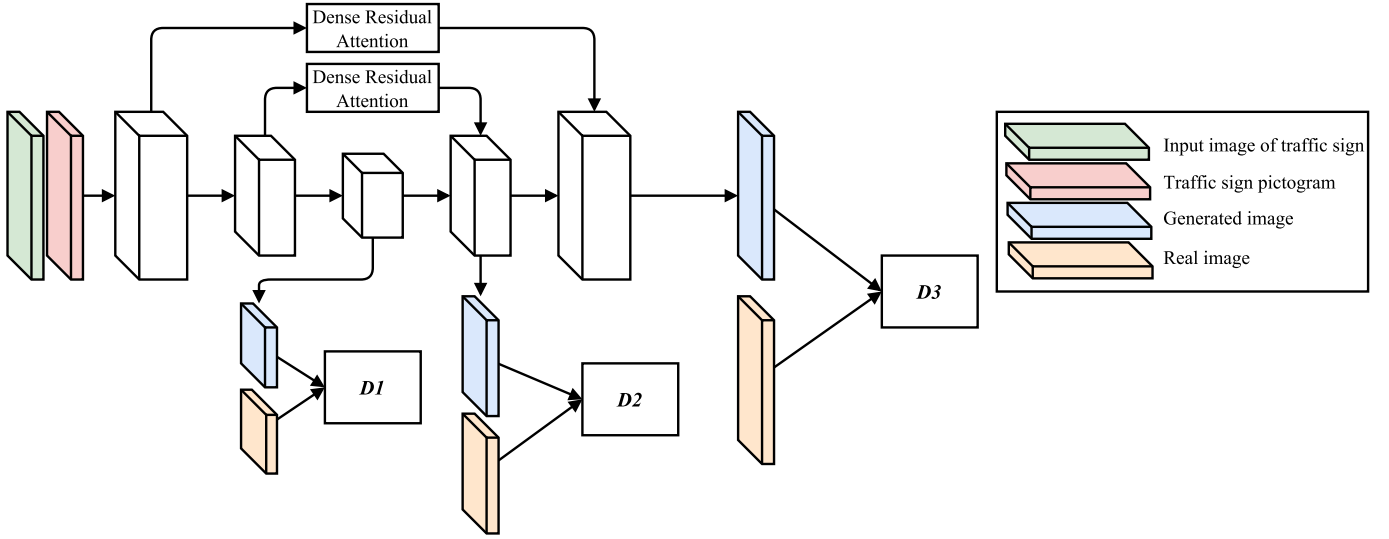


Fig. 1. Proposed network with multiple discriminators ($D1$, $D2$, $D3$). The network uses a Dense Residual Attention module and a discriminator at each scale. The Dense Residual Attention module receives the input from the feature maps F_e at the encoder side and outputs F_d^{out} . The output image generated from the auxiliary branch is supplied to the discriminator at a given scale. When the mask is applied, an element-wise product of the mask and the generated image is fed to the discriminator. The cycle loss is also obtained using the same generator network. Cycle loss not shown in figure. (Best viewed in color).

A. Dense connectivity

To enhance the feature propagation, we apply a dense connection between the encoder and decoder of the generator network. The dense connectivity is achieved by concatenating feature maps of the same size from the encoder to the decoder. Since the convolution over concatenated feature tensors is computationally expensive, we apply a 1×1 convolution across channels to reduce dimensionality [23]. The output feature tensor F_d^c at the decoder side is given as

$$F_d^c = [F_d, F_e]_{1 \times 1}, \quad (4)$$

where F_e and F_d are the feature maps at a given scale in the encoder and decoder. The expression $[\cdot, \cdot]_{1 \times 1}$ denotes the 1×1 convolution followed by the concatenation operation $[\cdot, \cdot]$. The 1×1 convolution reduces the dimensionality of F_d^c to the same size as F_e .

B. Attention mechanism through residual connections

The attention mechanism is necessary for learning the relevant features that must be passed to the subsequent layers. To retain the background of the input image, the features closer to the input should be preserved. The encoder side of the generator has features closer to the input and hence the information from the encoder is transferred to the decoder through a residual attention mechanism. Our method is similar to the approach proposed in [21]. However, we do not require a mask or trunk branch for feature processing, instead we couple the output feature maps from encoder and decoder. The proposed attention mechanism does not require any additional trainable parameters. The updated feature maps F_d^a after the attention mechanism is expressed as:

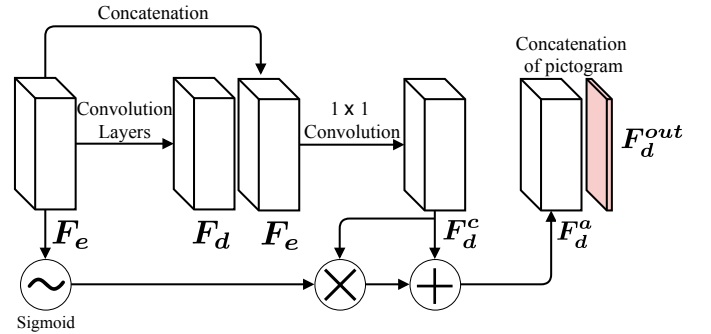


Fig. 2. Dense Residual Attention module followed by concatenation of the pictogram. Feature maps from the encoder F_e are concatenated with F_d , followed by 1×1 convolution to produce F_d^c . The Hadamard product of F_e and F_d^c followed by the addition of F_d^c results in F_d^a . The output obtained from the Dense Residual Attention module F_d^a with the pictogram p represents F_d^{out} .

$$F_d^a = F_d^c + \sigma(F_e) \odot F_d^c, \quad (5)$$

where σ denotes the sigmoid activation and \odot denotes the Hadamard product. At a given scale, we update the feature tensor in the decoder by the element-wise product of $\sigma(F_e)$ and F_d^c followed by the addition of F_d^c . We have also found that concatenating the pictogram p of the desired traffic sign at each scale improves the performance. The output tensor F_d^{out} at the decoder side is $[F_d^a, p]$. At larger scales, the concatenation of the traffic sign pictogram preserves the finer details of the target traffic sign. We refer to the combination of residual attention mechanism and dense connectivity as *Dense Residual Attention*. The Dense Residual Attention module is illustrated in Figure 2.

C. Multi-scale discriminators

To understand both the global and finer details, we train a discriminator at each scale. Concurrent to our work, multi-scale discriminators are used in [8]. However, our generator has auxiliary branches that generate an output image at each scale. We use multiple optimizers which are optimized jointly. We generate outputs after up-sampling the outputs from the previous layer by a factor of two. The generated image at a given scale is fed to the corresponding discriminator. To reduce computational cost, the generated image at a smaller scale use a discriminator of smaller depth. We use three discriminator networks ($D1, D2, D3$) which receive the input from the auxiliary branches of the generator. The discriminator at the smallest scale attends to global features such as lighting condition and pose of the traffic sign, whereas the concatenated pictogram p along with the discriminator at the largest scale captures the finer details of the traffic sign. Multi-scale discriminators simplify the transition going from the coarsest to the finest scale by retaining the global features.

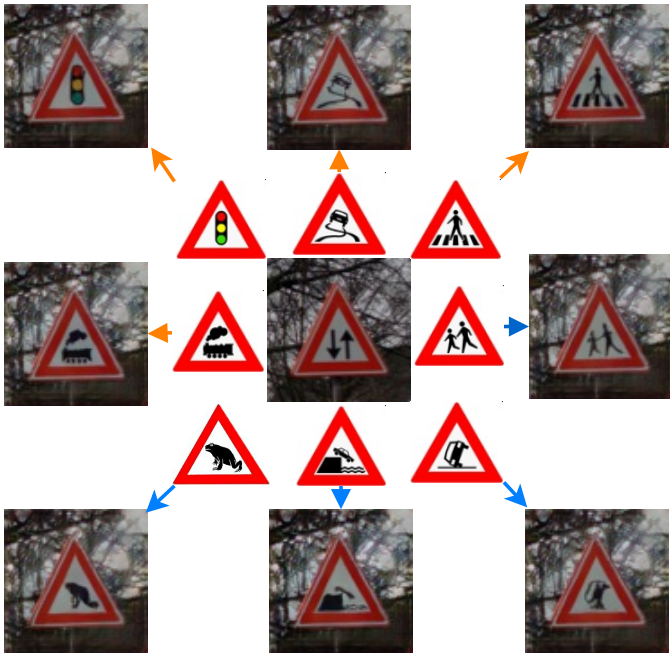


Fig. 3. Given an image of a traffic sign (center) and a pictogram, the trained model generates an image of a new traffic sign. Orange and blue arrows represent classes (not examples) inside and outside the training set.

D. Mask for weak supervision

In previous work [9], an implicit mask is generated to attend to the desired object. A convex combination of the input and the generated image produces the output image. This implicit mask generation is a tedious task when the desired object has varying light conditions and a cluttered background. With the methods described in the previous subsections, we have obtained perceptually appealing results when the pose of the traffic sign is not too skewed. However, we have found it beneficial to apply a mask to improve the performance. We use



Fig. 4. Examples of generated traffic signs using our method.

a rectangular bounding box on the desired object that intensity set to unity. The region of the mask around the bounding box has an intensity range within the unit interval, which changes during training starting from one. The final output that is supplied to the discriminator is the element-wise product of the mask and the generated image. The provided mask is only required during training and is not used during testing.

IV. EXPERIMENTS

A. Dataset

We have obtained the images of Dutch traffic signs within high-resolution street-view images from the company Cyclomedia and the pictograms from [24]. Out of a total of 313 classes, we select images of traffic signs from 55 classes that have a low amount of data and low recognition rates. We broadly partition the images into three categories, based on the appearance of the traffic signs as white triangles, white circles and blue rectangles. Each of the image and the pictogram has a resolution of 80×80 pixels.

B. Implementation details

We found that it is easier to transfer traffic signs from a pictogram within a category rather than transferring it across categories. Therefore, we use a model for each category (each category has several classes). For the generator network, a ResNet with encoder-decoder structure is applied as the backbone architecture. The residual connection is replaced by *Dense Residual Attention*, while multi-scale discriminators are applied. The generator up-samples the output at each scale using bilinear up-sampling, followed by convolution with residual units. We use a discriminator at each scale,

TABLE I

Classification performance of three categories of traffic signs that consist of 55 classes. Each class is approximately expanded by 300 examples.

Category information		Amount of training data		Classification performance (Top-1 score)		
Category	Number of classes	Real training data	Generated training data	Real data	Real + Generated data	Difference
White triangles	26	6498	7828	50.0%	55.3 %	+5.3%
White circles	23	19100	7200	68.6%	70.1%	+1.5%
Blue rectangles	6	6679	2074	65.1%	66.5%	+1.4%

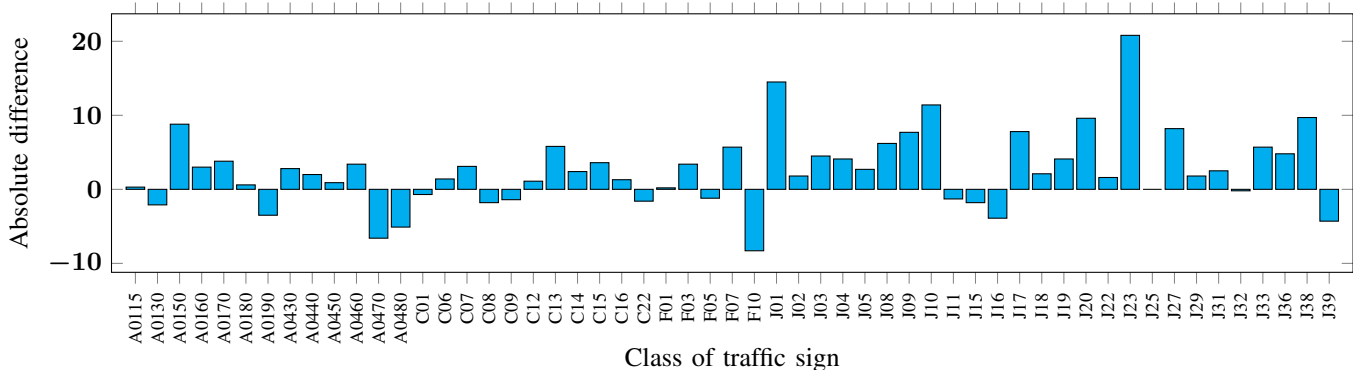


Fig. 5. Overview of the classification performance of traffic signs across 55 classes. The y axis represents the absolute difference of top-1 accuracy with respect to the baseline (trained with real data). The x axis represents the class of traffic sign which is described in [24].

resulting in a total of three discriminators. The discriminators have networks of varying complexity depending on the image scale. We use three residual units for the image generated at the highest resolution (80×80 pixels), two residual units at 40×40 pixels and a single residual unit at 20×20 pixels. The last layer of each of the discriminator network is a fully-connected layer. Adam optimizer [25] is applied to each discriminator. For the mask, we conduct experiments with both rectangular and circular masks across all classes. However, we have found that the circular mask around the desired object offers the best performance, irrespective of the shape of the traffic sign. Depending on the class of traffic sign, we have trained the model between 60K and 100K iterations. Generation of 1000 images took approximately 57 seconds on a NVIDIA Tesla P100 GPU (12 GB).

C. Ablation study

With a standard ResNet structure for the generator, we observe that the generated traffic sign has a poor geometry and texture. It also produces backgrounds that are unrealistic and have low perceptual quality. Besides this, we note that noise present at the lower scales of the network is up-sampled, resulting in large noisy patches. This is shown in Figure 6. We conduct an ablation study to understand the contribution of each component of the proposed network. The results using our method are illustrated in Figures 3 and 4.

1) *Dense Residual Attention*: The addition of residual attention mechanisms suppress the noise and produce a clearer background. We have observed a further improvement in the detail of the visual quality after concatenating the encoder feature maps. The dense concatenation partially improves the

geometry of the traffic sign. However, the dense concatenation without residual attention does not retain the background information from the input image.

2) *Multi-scale discriminators*: The multi-scale discriminator captures both the global and local details of the desired image from the input image and pictogram. Without the multi-scale discriminator, the edges and other geometrical features of the traffic sign are not sharp. The outputs at the smaller scale produce consistent geometries which guide the higher-resolution outputs. At the smallest scale, the details of traffic sign texture are absent, whereas the lighting condition and pose are learned. The concatenation of the pictogram at every scale captures the texture of the traffic sign.

3) *Mask*: With the addition of a mask, we can observe an improved performance in replacing the texture of the traffic sign. We also observe a better performance in situations where the traffic sign has a complex pose (skewed angle) or texture.

D. Detection and classification results using generated data

We have generated samples using our method and added the new data to the existing training set. For training, we use an ensemble of HOG-SVM and CNN detectors [26][27]. For detection, the baseline recall is 99% on the test set and we did not notice any significant improvement in the recall with the addition of generated sets. However, the number of false positives decreased between 1.2% - 1.5% in the detection tests. This is beneficial in reducing manual efforts to remove false-positives from the automatically extracted detections. We have also conducted classification tests using the generated data. The details of the classification results are shown in Table 1. Out of the 313 classes, 55 classes that have a low amount of

data, are expanded. Among the 55 classes, 41 classes produced a higher classification rate when trained with the combination of real and generated samples. The lowest score obtained for a class is 8.3% (average decrement by 2.92%) below the baseline, whereas the highest gain has 20.8% absolute improvement (average increment by 4.65%) over the baseline. We do not conduct detection and classification performance analysis of classes outside the training set as there is no real test set.

E. Failure cases and comparison with other methods

Figure 6 demonstrates examples of failure cases with our method. In the first and second row, we observe lack of details when the traffic sign in the pictogram are not thick or when traffic signs have a skewed angle. At the bottom row (Figure 6), we observe noisy outputs that progressively become larger, as the network scale increases in the generator. However, the residual attention mechanism exhibits a certain amount of robustness to this type of noise. We observed fewer cases with such noise compared to Conditional GANs.

We conducted experiments with other methods as well and results are presented in Figure 7. WGAN-GP generates images from noise, which results in smeared backgrounds and traffic signs. Boundary Equilibrium GAN (BEGAN) often results in mode collapse (image generated from noise) that produce a low variety of samples. Conditional Analogy GAN (CAGAN) uses an implicit mask generation fails with street-view images due to complex backgrounds which result in poor outputs. Conditional GANs (cGANs) produces outputs, which are often smudged and have the incorrect geometry of the traffic sign.

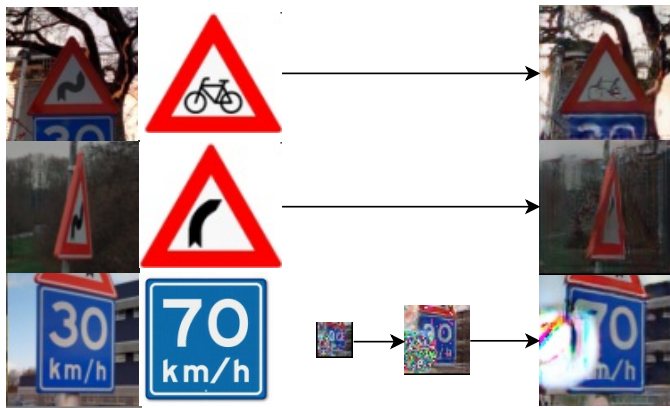


Fig. 6. Examples of failure using our method. Top row: Generated output lacks details in the traffic sign and background. Middle row: Skewed angled traffic signs have difficulty producing pictogram textures. Bottom row: Noise present in the lower layers progressively grow into a large noise.

V. CONCLUSIONS

We have presented a conditional GAN with Dense Residual Attention, which generates a new traffic sign conditioned on a given pictogram. The network utilizes multiple Dense Residual Attention modules that are composed of a residual attention mechanism and a dense concatenation. The Dense Residual



Fig. 7. Row 1 (top): WGAN-GP (generated from noise), Row 2: Boundary Equilibrium GAN (generated from noise), Row 3: Conditional Analogy GAN, Row 4: Conditional GAN, Row 5 (bottom): Conditional GAN with Dense Residual Attention and multi-scale discriminators (our method).

Attention module improves the visual quality of the traffic signs and suppresses the cases with progressively growing noise. We propose the use of multi-scale discriminators, which result in images of traffic signs that are both globally and locally coherent. The discriminator at smaller scale captures global features and steers the high-resolution output to produce images with more accurate geometries. The discriminator at the larger scale along with the concatenated pictogram assists in producing images of traffic signs with finer details. Comparison other methods reveals that the proposed method produces visually appealing results with finer details in the traffic signs and has fewer geometrical errors. We have further conducted detection and classification tests across a large number of traffic sign classes, by training our detector with the combination of real and generated data. The trained model reduces the number of false positives by about 1.2 - 1.5% at a recall of 99% in the detection tests while improving the top-1 accuracy on the average by 4.65% in classification tests.

REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [2] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *European Conference on Computer Vision*. Springer, 2016, pp. 649–666.
- [3] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711.

- [4] T. Kim, M. Cha, H. Kim, J. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," *arXiv preprint arXiv:1703.05192*, 2017.
- [5] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [6] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *arXiv preprint arXiv:1611.07004*, 2016.
- [7] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," *arXiv preprint arXiv:1703.10593*, 2017.
- [8] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional gans," *arXiv preprint arXiv:1711.11585*, 2017.
- [9] N. Jetchev and U. Bergmann, "The conditional analogy GAN: Swapping fashion articles on people images," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [10] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544.
- [11] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," *arXiv preprint arXiv:1609.04802*, 2016.
- [12] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [13] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [14] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.
- [15] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of wasserstein gans," *arXiv preprint arXiv:1704.00028*, 2017.
- [16] H. Larochelle and G. E. Hinton, "Learning to combine foveal glimpses with a third-order boltzmann machine," in *Advances in Neural Information Processing Systems 23*, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc., 2010, pp. 1243–1251.
- [17] J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, and B.-T. Zhang, "Multimodal residual learning for visual qa," in *Advances in Neural Information Processing Systems*, 2016, pp. 361–369.
- [18] H. Zheng, J. Fu, T. Mei, and J. Luo, "Learning multi-attention convolutional neural network for fine-grained image recognition," in *Int. Conf. on Computer Vision*, 2017.
- [19] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [21] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [22] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [23] Q. C. M. Lin and S. Yan, "Network in network," in *ICLR 2014*, 2014.
- [24] "Comprehensive overview of Dutch road signs," <http://www.verkeersbordenoverzicht.nl/>.
- [25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [26] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1. IEEE, 2005, pp. 886–893.
- [27] I. M. Creusen, R. G. Wijnhoven, E. Herbschleb, and P. de With, "Color exploitation in hog-based traffic sign detection," in *Image Processing (ICIP), 2010 17th IEEE International Conference on*. IEEE, 2010, pp. 2669–2672.