





![](_page_0_Figure_5.jpeg)

Figure 2. CycleGAN and UNIT are able to perfectly reconstruct the input image despite the semantic inaccuracy of the translations. Adding low-amplitude Gaussian noise ( $\mu = 0, \sigma = 0.08$ ) to the translated image completely destroys the reconstruction. Our defense techniques (rows 3 and 4) rely on the translated image instead of embedded noise and are more robust to random perturbations.

## **Adversarial Self-Defense for Cycle-Consistent GANs** Dina Bashkirova<sup>1</sup>, Ben Usman<sup>1</sup> and Kate Saenko<sup>1 2</sup> <sup>1</sup>Boston University <sup>2</sup>MIT-IBM Watson AI Lab

Figure 5. Quantized reconstruction results of the original CycleGAN, CycleGAN with noise defense and CycleGAN with guess loss defense. The last column represents the translation from the corresponding ground truth semantic segmentation map to real frame for comparison.

### Results

![](_page_0_Figure_24.jpeg)

![](_page_0_Figure_25.jpeg)

![](_page_0_Figure_26.jpeg)

Figure 6. One-to-one translation results on SynAction dataset.

Method	acc. segm ↑	IoU segm↑	IoU p2p↑	RH↓	SN↓
CycleGAN	0.230	0.16	0.20	$27.43 \pm 6.14$	446.92
CycleGAN + noise*	0.240	0.17	0.23	$9.17\pm7.37$	94.15
CycleGAN + guess*	0.237	0.17	0.21	$11.38\pm7.03$	212.59
CycleGAN + guess + noise*	0.236	0.17	0.24	$\textbf{6.1} \pm \textbf{5.85}$	150.55
UNIT	0.080	0.04	0.06	$6.37 \pm 11.69$	361.52
MUNIT + cycle	0.131	0.08	0.17	$2.5\pm8.86$	244.95
pix2pix (supervised)	0.4	0.34		<u> </u>	_

Table 2. Results on the GTA V dataset. acc. segm and IoU segm represent mean class-wise segmentation accuracy and IoU, IoU p2p is the mean IoU of the pix2pix segmentation of the segmentation-to-frame mapping; RH and SN are the quantized reconstruction honesty and sensitivity to noise of the many-to-one mapping (B2A2B) respectively. \* -- our proposed defense methods. The reconstruction error distributions plots can be found in the supplementary material (Section 2).

Method	acc. segm↑	IoU segm↑	IoU p2p↑	$\mathbf{RH}\downarrow$	SN↓
CycleGAN	0.233	0.175	0.210	$21.77\pm5.16$	251.19
CycleGAN + noise*	0.242	0.187	0.218	$12.27\pm4.42$	222.18
CycleGAN + guess*	0.241	0.184	0.224	$7.47 \pm 2.38$	235.43
CycleGAN + guess + noise*	0.249	0.191	0.218	$\textbf{-0.45} \pm \textbf{2.26}$	238.25
UNIT	0.212	0.153	0.124	$19.63\pm6.07$	528.22
MUNIT + cycle	0.153	0.094	0.124	$21.43 \pm 7.85$	687.27
pix2pix (supervised)	0.301	0.234	9 <del></del>		9 <del></del>

Table 3. Results on the Google Maps dataset. The notation is same as in the Table 1.

### Acknowledgements

This project was supported in part by NSF and DARPA.

### References

- Richter, Stephan R., et al. "Playing for data: Ground truth from computer games." *European conference on computer vision*. Springer, Cham, 2016.
- Huang, Xun, et al. "Multimodal unsupervised image-to-image translation." Proceedings of the European Conference on Computer Vision (ECCV). 2018.
- Chu, Casey, Andrey Zhmoginov, and Mark Sandler. "CycleGAN, a master of steganography." *arXiv preprint arXiv:1712.02950* (2017).

# BOSTON UNIVERSITY

Figure 4. Illustration of sensitivity of cycle-consistent translation methods to high-frequency perturbations in one-to-many (left) and in many-to-one

• Zhu, Jun-Yan, et al. "Unpaired image-to-image translation using cycle-consistent adversarial networks." *Proceedings of the IEEE international conference on computer vision*. 2017. • Liu, Ming-Yu, Thomas Breuel, and Jan Kautz. "Unsupervised image-to-image translation networks." *Advances in neural information processing systems*. 2017.

• Sun, Ximeng, Huijuan Xu, and Kate Saenko. "A Two-Stream Variational Adversarial Network for Video Generation." arXiv preprint arXiv:1812.01037 (2018).