



Learning to Adaptively Adjust Fovea Size for Human Eye Inspired Recurrent Neural Network

Alan Van Omen, Sehong Oh

April 26th, 2023

Outline

1. Introduction
2. Motivation & Related Works
3. Methodology
4. Experiments
5. Conclusion

Introduction

1. Human vision and computer vision

- Computer vision has a different way to see the world from humans
- Each system has its own pros and cons, and "capacity" and "vulnerability" are the distinct properties that differentiates them



```
82 62 63 64 65 66 67 67 69 70 71 72 72 73 73 73 72 72 71 70 69 67 66 66 66 65 63 62 61 60 6
81 62 63 64 66 66 67 68 68 69 70 71 72 72 73 72 72 71 70 69 68 66 66 65 63 62 61 60 6
51 62 63 64 66 66 68 68 69 70 71 72 73 73 72 72 71 71 69 68 67 66 66 65 65 64 63 62 61 6
81 63 64 64 66 67 68 68 69 70 71 71 73 73 74 73 73 71 70 69 68 66 66 65 64 63 62 61 61 6
81 63 64 65 67 68 69 70 71 71 72 55 69 72 72 71 71 70 69 68 67 66 65 64 63 62 60 60 6
63 64 65 66 67 68 69 70 71 72 12 4 5 11 68 72 71 71 69 69 68 67 66 65 64 62 62 60 59 5
63 65 66 66 68 68 69 70 71 71 72 18 4 4 7 6 6 71 70 69 68 67 66 65 64 63 61 60 59 59 5
63 65 67 67 68 69 70 71 72 64 4 27 24 54 33 20 28 64 68 67 66 65 64 63 62 61 59 58 5
64 65 66 66 68 69 70 71 11 24 24 12 17 24 28 60 37 43 30 62 66 68 67 66 65 64 63 61 60 59 58 5
65 66 67 67 68 69 71 10 6 6 5 34 36 12 47 34 17 29 54 43 63 67 66 65 64 63 62 60 59 58 5
64 65 66 66 68 69 10 6 6 5 5 7 16 19 4 17 44 27 24 60 67 66 66 65 65 64 63 61 60 59 58 5
63 64 65 65 67 30 6 6 5 5 6 8 9 20 27 51 78 41 44 66 65 65 65 65 64 63 62 60 59 58 5
63 64 65 65 34 5 5 5 5 5 5 4 19 6 7 34 61 20 59 65 65 64 64 64 63 62 61 60 59 57 9
63 64 64 65 14 5 6 5 5 4 5 4 18 7 5 4 19 10 11 65 64 64 64 63 61 60 62 61 60 59 58 5
63 64 64 65 33 7 4 5 6 6 7 10 6 5 5 4 21 24 18 64 64 64 63 62 64 65 62 62 60 59 58 5
64 64 64 64 65 50 4 4 4 5 11 16 6 6 4 6 35 16 26 66 64 64 63 61 72 67 63 62 61 59 58 5
64 64 64 64 65 46 4 4 4 5 6 9 8 5 20 10 43 56 29 57 64 64 63 61 70 67 62 64 65 59 59 5
64 64 64 65 66 27 5 4 4 5 6 6 6 18 66 20 57 60 46 57 75 70 62 81 70 67 62 61 60 59 58 5
49 58 62 65 57 5 5 6 6 6 6 18 59 22 60 58 44 22 63 71 72 60 69 68 61 60 58 59 59 5
62 62 57 59 28 5 5 5 5 5 5 70 30 43 61 62 64 5 28 64 60 62 56 63 65 66 67 61 53 5
52 52 32 33 3 6 5 5 5 5 6 11 39 21 31 51 50 45 46 18 32 36 33 23 44 70 71 51 42 27 3
50 50 51 39 5 5 5 5 6 6 6 6 42 69 34 42 39 43 37 26 29 26 29 26 35 33 18 1
52 53 51 22 5 5 5 5 6 6 5 44 56 17 51 54 53 54 56 51 22 54 54 55 55 54 53 53 52 5
54 54 53 5 5 5 5 6 6 6 13 52 42 21 51 54 51 49 49 50 22 41 45 42 42 41 40 41 44 43 4
52 52 54 5 5 5 5 6 6 6 6 28 56 22 54 53 51 51 51 51 44 23 51 51 49 50 49 48 46 4
54 54 52 53 38 7 5 6 6 6 6 40 54 62 51 53 56 55 52 51 51 52 52 50 49 48 46 45 46 4
51 52 51 53 27 14 5 4 5 4 7 47 51 11 49 47 49 49 52 52 54 29 33 14 48 46 47 47 46 46 4
48 50 51 53 25 14 17 8 4 4 17 46 40 18 43 47 46 49 52 54 53 53 54 15 50 49 46 47 47 47 4
49 49 49 49 22 12 20 24 6 14 35 51 39 48 48 50 51 51 49 51 51 52 50 41 58 48 47 47 45 45 4
51 49 50 50 22 13 19 36 13 12 42 50 49 73 50 50 50 49 48 49 48 49 45 51 46 44 44 44 42 45 4
47 49 49 47 20 16 26 21 15 36 48 42 61 47 48 51 47 50 51 51 49 47 47 52 47 47 44 43 45 4
```



Motivation

1. Computer vision models (CNNs) are vulnerable to adversarial noises^[1]

AllConv



SHIP

CAR(99.7%)

NiN



HORSE

FROG(99.9%)

VGG



DEER

AIRPLANE(85.3%)

2. Human eyes are more robust than computer vision

- Allowing models to fixate to different image regions can alleviate the effect of adversarial noise^[2]
- Non-uniform spatial sampling and varying receptive fields that mimic the retinal transformation in the primate retina can also improve the robustness against adversarial attacks^[3]

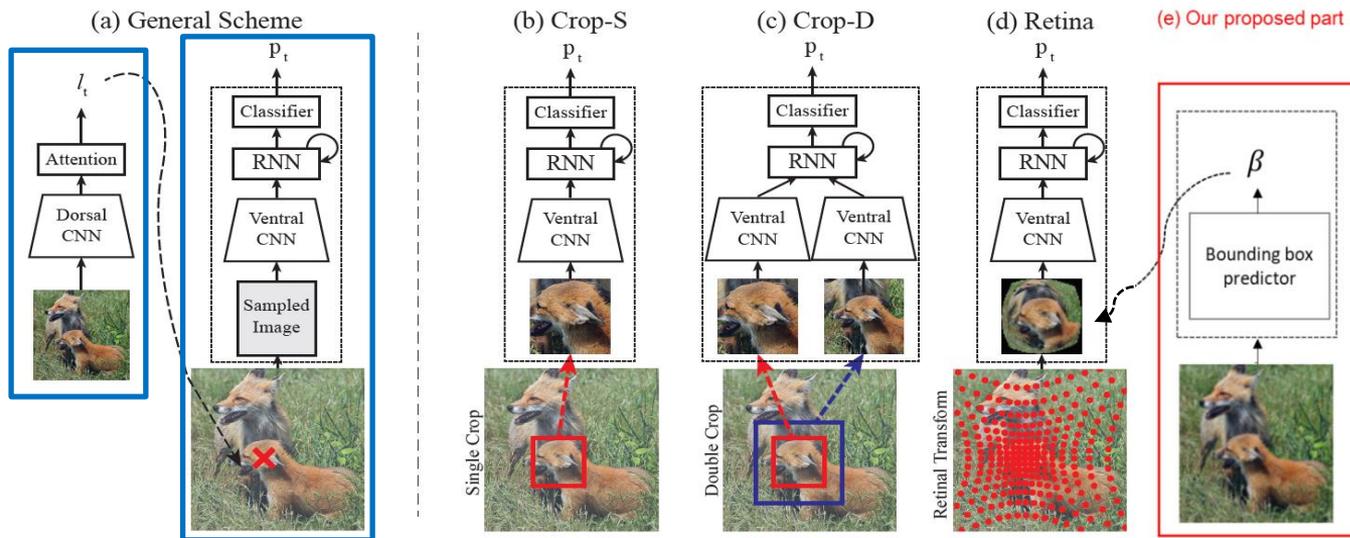
[1] Aleksander et al. Towards deep learning models resistant to adversarial attacks, 2017.

[2] Ricardo Gattass et al. Visual topography of v2 in the macaque, 1981.

[3] Pouya Bashivan et al. Neural population control via deep image synthesis, 2019.

Methodology

1. Dorsal Stream: or the "where" pathway is responsible for processing information related to spatial awareness and motion.
2. Ventral stream: or the "what" pathway, is responsible for object recognition, face recognition, and determining the color and shape of objects.



[4] Choi, Minkyu, et al. Human Eyes Inspired Recurrent Neural Networks are More Robust Against Adversarial Noises, 2022

Improvement

1. The hyper-parameter b in the model controls the degree of non-uniform sampling for the ventral stream (authors use constant $b=12$)

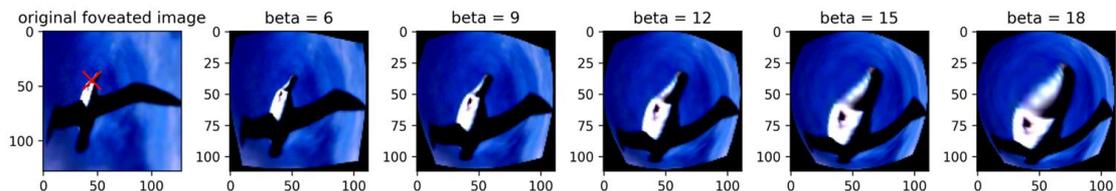
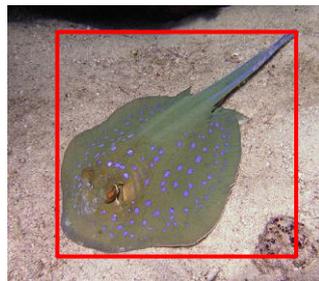


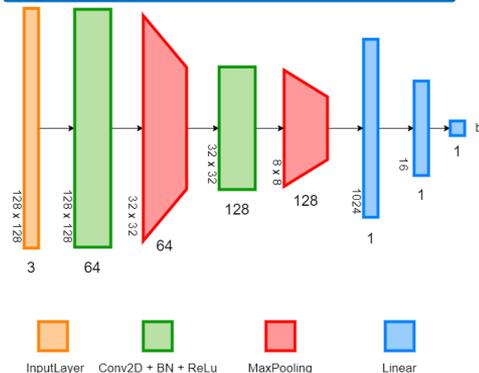
Figure 2. Examples of the effect of b on the retinal transformation, with fixation point centered at red x.

2. We implement 3 methods for adaptively changing the fovea size to improve robustness

Bounding Box Methods (YOLOv5)

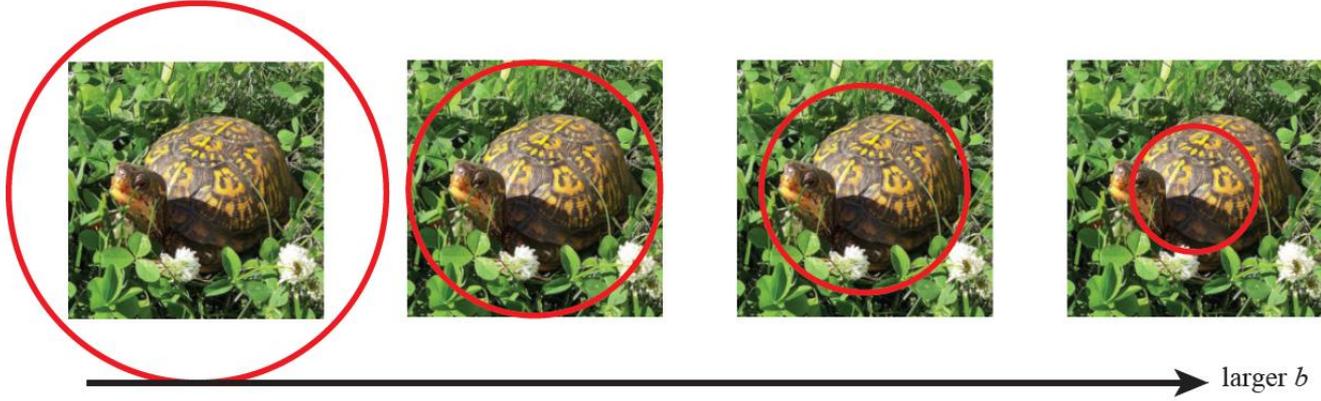


Convolution Neural Network



Effect of b on the fovea size

(a) Size of Fovea



(b) Retinal images



Improvement

1. The hyper-parameter b in the model controls the degree of non-uniform sampling for the ventral stream (authors use constant $b=12$)

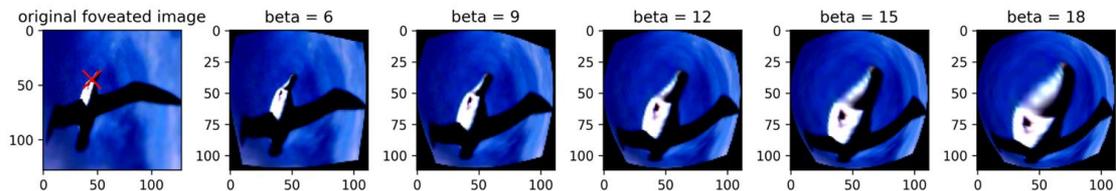
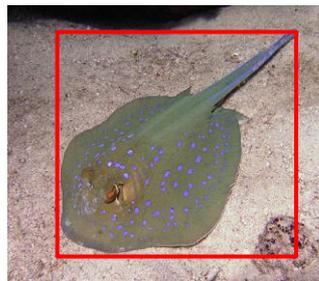


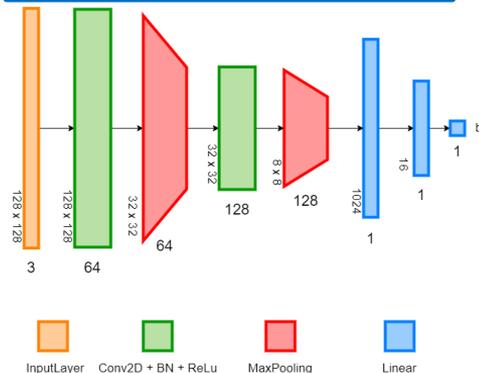
Figure 2. Examples of the effect of b on the retinal transformation, with fixation point centered at red x.

2. We implement 3 methods for adaptively changing the fovea size to improve robustness

Bounding Box Methods (YOLOv5)



Convolution Neural Network



Experiments

- Trained four models using various strategies for adapting the fovea size
 1. **b=12**: held beta constant like original paper
 2. **B-CNN**: learned beta concurrently in an end-to-end manner
 3. **C-BB**: estimated beta from the closest valid bounding box
 4. **L-BB**: estimated beta from the largest valid bounding box
- Learned adversarial noise with 100 steps of PGD over 1600 test images

| Parameter | Value |
|----------------------------------|-------|
| batch size | 4 |
| optimizer | Adam |
| initial learning rate | 1e-3 |
| epochs | 10 |
| retinal sampling grid size | 112 |
| time steps (1 fixation per step) | 4 |
| default b | 12 |

| Model | Top-1 Acc | ASR ($\epsilon = 5e - 3$) | |
|----------|--------------|-----------------------------|--------------|
| | | Targeted | Untargeted |
| $b = 12$ | 35.9% | 80.0% | 80.2% |
| b -CNN | 38.0% | 79.6% | 82.9% |
| C-BB | 39.8% | 78.9% | 80.6% |
| L-BB | 28.6% | 78.4% | 78.9% |

Conclusion

1. Successfully implemented 3 methods for adaptively setting the fovea size
2. Found that C-BB gave slightly higher robustness and better accuracy

Future work

1. Explore how to fine-tune bounding box model on target dataset.
2. Optimize the model efficiency and memory usage to allow for larger batch sizes and better stability.
3. Explore using some reward to promote robustness for learning beta in the end-to-end manner.

Reference

- [1] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [2] Ricardo Gattass, CG Gross, and JH Sandell. Visual topography of v2 in the macaque. *Journal of Comparative Neurology*, 201(4):519–539, 1981.
- [3] Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image synthesis. *Science*, 364(6439), 2019.
- [4] Choi, Minkyu, et al. "Human Eyes Inspired Recurrent Neural Networks are More Robust Against Adversarial Noises." *arXiv preprint arXiv:2206.07282* (2022).

Q & A

Thank you !