

Understanding the Effects of Biological Predispositions in Human Vision in Thwarting Adversarial Attacks

Varshini Reddy
MIT 6.819 - Advances in Computer Vision
Harvard University

varshinibogolu@fas.harvard.edu

Abstract

The structure of deep nets aspire loosely to emulate the human brain, beginning from the introduction of nodes to imitate neurons in the brains to adjusted strengths of connections between neurons to reflect their associations similar to that in the visual cortex. In short, there is a close association of cognitive and neuroscience with computer vision.

One of the grand challenges of computer vision lies in understanding how the brain recognizes objects in the visual world to improve machine vision. To this extent, this work aims to incorporate certain human vision traits such as acuity and initial blurred vision into known state of the art object recognition architectures, which is VGG-16 in our case. Additionally, given that humans are immune to adversarial attacks, we also try to understand whether human motivated computer vision architectures can help in mitigating the effects of adversarial attacks.

1. Introduction

Object detection and recognition has always been at the forefront of computer vision, even before era of neural networks. There have been many perspectives of how object recognition takes place. Humphreys and Bruce, in 1989, proposed a model of object recognition that fits a wider context of cognition. According to them, the recognition of objects occurs in a series of stages. First, a sensory image is generated, following a perceptual classification, where the information is compared with previously stored descriptions of objects. Another approach to this problem was proposed by Marr and Nishihara, where we store an object as a 3D model which can be used to make predictions, thus making them transformations invariant. This was followed by a proposal by Biedermann. Regardless of the approach, one thing remained constant among all proposals. The approach used to achieve object recognition was to

try emulate the human interpretation of vision detection and recognition.

Nature has caused humans to be born with multiple predispositions which might seem to be sensory function limitations initially, however they have been proven to be beneficial in the long run. An example of this is babies preferring sweet as compared to bitter food. This has been studied to help them eat energy packed and while discouraging the consumption of toxins. The motivation behind this work are 2 such biological predisposition. One babies starting at an acuity of 20/600 (which causes blurred and squinted vision) [5]. The second is the limited perception of colour in infants [1]. While the former helps focus on interpreting local patterns, the latter helps focus on detecting edges and other shapes rather than learning object recognition using only colours.

Regardless of the size of the image corpus' models, today, are trained on, such as ImageNet [2], they are remain susceptible to simple Adversarial attacks and failure in identifying Out-of-Distribution images. However, humans have the capacity to classifying objects or images they have never seen before successfully, though it may be into a broad category. The motivation behind this paper is to understand whether there are any traits of human vision which enable us to perform so well on unseen and adverse images.

2. Related Works

There have been immense efforts in understanding the workings of human brains and discovering the visual features and representations used by the brain to recognize objects. Further, efforts have been focused on applying this knowledge to improve computer vision. Recently, neural network models of visual object recognition, including biological and deep network models, have shown remarkable progress and have begun to rival human performance in some challenging tasks. These models are trained on image examples and learn to extract features and representations and to use them for categorization. It remains unclear, how-

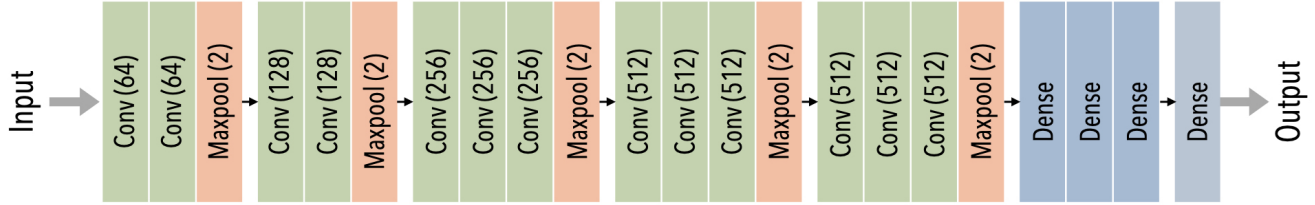


Figure 1. Baseline architecture VGG-16

ever, whether the representations and learning processes discovered by current models are similar to those used by the human visual system [9].

In a more recent work [3], the authors propose to improve machine vision by incorporating context representation of scenes as perceived by humans without objects (in most cases). With a systematic paradigm developed to capture human contextual expectations they were able to improve model performance as compared to generic image augmentation techniques. In yet another work [6], the authors talk of the differences in the perceptual distance of objects between human and computer vision. To this extent, they propose fixing the biases that cause machines to underestimate the distances between symmetric objects compared to human perception.

3. Proposed Method

In this work, we attempt to cover the effects of three human traits in general object detection and classification and in mitigating adversarial attacks.

3.1. Dataset and Pre-processing

For the purpose of this work, we use the a subset of the ImageNet dataset. In particular, we use the ImageNet 2012 Validation set for both model training and evaluation [7]. This dataset contains 50,000 images belonging to 1000 classes, all of which are balanced. However, owing to the limited computational capacity, we use only 15,000 images for both train and evaluation. Each image is resized to dimensions 224x224x3.

3.2. Proposed Model Architecture

For incorporating all these, we will be using the VGG-16 architecture [8] as the baseline. The weights of this model trained on the original ImageNet dataset is what is going to be used to reduce the required training epochs. The baseline architecture is shown in Fig. 1.

3.2.1 Reduced Visual Acuity

As discussed in the previous sections, infants are born with lower acuity than normal adults. As the child grows this acuity naturally improves. This initial disadvantage allows

us to learn local pattern relationships by forcing a smaller receptive field and blurred vision. To incorporate this into machine vision, we input the same image at various stages of the network as input. The image at each stage is pre-processed to have a varying level of blur, starting with highly blurred images. The blurring is done using a Gaussian kernel. The proposed architecture to achieve this is shown in Fig. 2.

3.2.2 Depth Perception

Another advantages humans have is the ability to perceive depth when they look at a scene. Though recent studies have shown that humans do not directly encode information in 3D, but rather as a sequence of 2D views, the essence of depth is still encoded in the information stored. This allows for an altered context and object representation as compared to how the learned representations are in computer models. In an effort to incorporate this we process the input image to compute its depth map. The input to the baseline architecture is a stack of the original input image with its corresponding depth map. This is depicted in 3.

3.2.3 Edge Detection

As seen in both, humans and machines, the kernels (Gabor filter in humans) learn small edges and other contrasting features initially followed by larger image features. To reduce the training requirement and find if reducing noise which makes the model sensitive we can improve its adversarial robustness. This is also to take into account the presence of rods and cones in the eyes, which enables both during day and night. For this find the edges of the image and input it to the baseline architecture. The structure to achieve this is depicted in Fig. 4.

3.3. Training and Evaluation

For each of these models, we use 70% of the 15,00 images as train and the remaining 30% for evaluation.

For the architecture that takes as input the blurred image, we freeze the weights of the baseline architecture. We train the weights for the added layers and the full connected section for the first set of epochs and finally train all the layers for the second set of epochs. Similarly, for the depth and

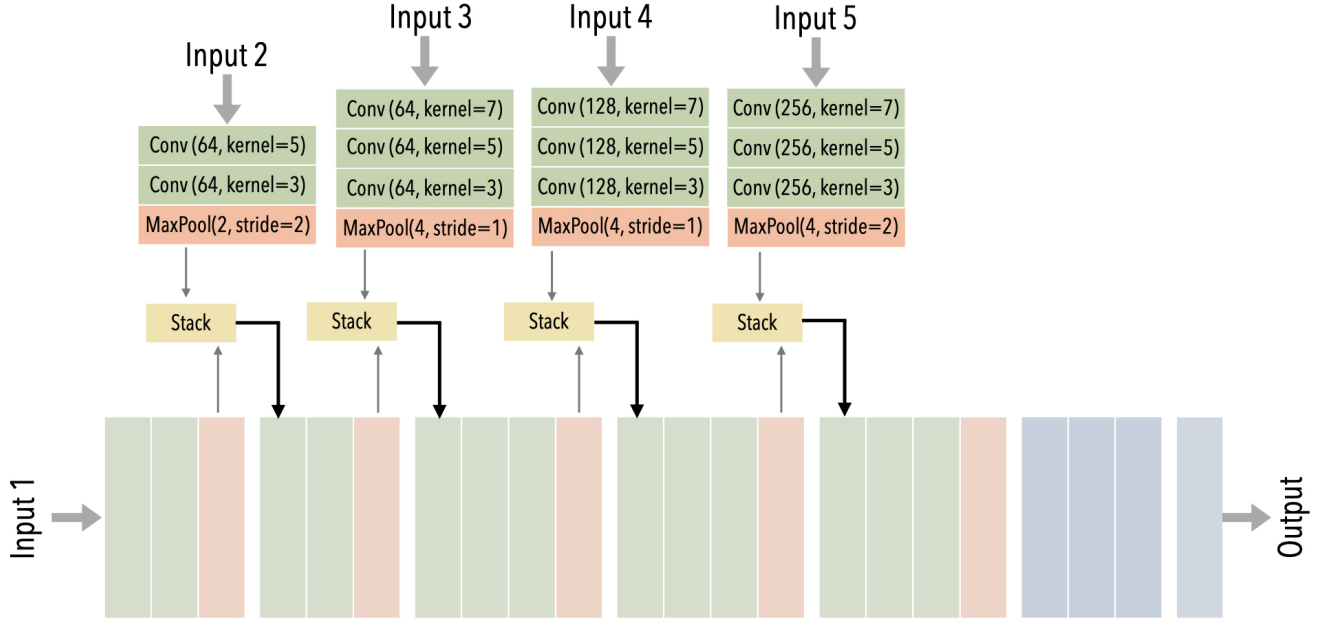


Figure 2. Altered baseline architecture to incorporate blurred images

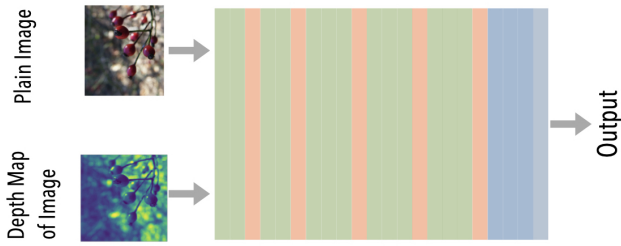


Figure 3. Altered baseline architecture to incorporate depth information

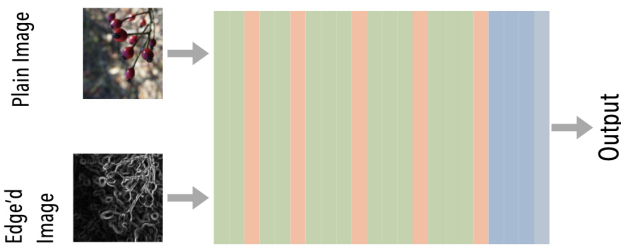


Figure 4. Altered baseline architecture to edge information

edge information incorporation experiments, we trained the model only for the weights that were not part of the original architecture along with the fully-connected layer. For the later part of the training, we briefly train all the layers.

To evaluate these models, we compute the performance as a function of accuracy (since it is a balanced dataset) on the original ImageNet test set. Additionally, to compute the adversarial robustness, we create adversarial samples from

the same test set and compare the model accuracy. For adversarial image generation we used an existing white-box attack PGD algorithm described in [4].

Thus, we mainly compute the difference between the accuracy of the baseline architecture, altered architecture where the images are not modified (all inputs are the plain image) and finally the altered architectures with the processed images.

4. Progress

- Baseline model with pre-defined weights evaluated on the test set and adversarial images.
- Defined the altered architecture for taking blurred images trained.
- Defined the algorithms for getting depth maps and edges of images.
- Defined the algorithms to get the adversarial images given a proxy model.

5. Plan till the Deadline

- Train the remaining architectures.
- Evaluate all models on the test images.
- Evaluate all models on the adversarial images.

References

- [1] Human infant color vision and color perception. *Infant Behavior and Development*, 2:241–273, 1979. [1](#)
- [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. [1](#)
- [3] Harish Katti, Marius V. Peelen, and S. P. Arun. Machine vision benefits from human contextual expectations. *Scientific Reports*, 9(1):2112, Feb 2019. [2](#)
- [4] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks, 2017. [3](#)
- [5] Daphne Maurer and Terri L. Lewis. Visual acuity: the role of visual input in inducing postnatal change. *Clinical Neuroscience Research*, 1(4):239–247, 2001. [1](#)
- [6] R. T. Pramod and S. P. Arun. Improving machine vision using human perceptual representations: The case of planar reflection symmetry for object classification. *IEEE transactions on pattern analysis and machine intelligence*, 44(1):228–241, Jan 2022. 32750809[pmid]. [2](#)
- [7] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [2](#)
- [8] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014. [2](#)
- [9] Shimon Ullman, Liav Assif, Ethan Fetaya, and Daniel Harari. Atoms of recognition in human and computer vision. *Proceedings of the National Academy of Sciences of the United States of America*, 113(10):2744–2749, Mar 2016. 26884200[pmid]. [2](#)