

---

# Foveated Downsampling Techniques

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Foveation is an important part of human vision, and a number of deep networks  
2 have also used foveation. However, there have been few systematic comparisons  
3 between foveating and non-foveating deep networks, and between different variable-  
4 resolution downsampling methods. Here we define several such methods, and  
5 compare their performance on ImageNet recognition with a Densenet-121 network.  
6 The best variable-resolution method slightly outperforms uniform downsampling.  
7 Thus in our experiments, foveation does not substantially help or hinder object  
8 recognition in deep networks.

## 9 1 Introduction

10 The retinas of humans, monkeys, and many other animals have a high-resolution fovea. In humans,  
11 this disproportionate representation of the central visual field carries through the whole visual cortex,  
12 and eye movements to foveate task-relevant features are an essential part of vision. Deep convolutional  
13 networks are inspired by the primate visual system, but they usually lack foveation, which may be a  
14 limitation in some contexts. In humans, foveation allows both the wide field of view needed for tasks  
15 like visual navigation, and the high resolution needed for tasks like reading, without impractical brain  
16 size or metabolic cost. Similar benefits may await deep networks. Some previous studies have used a  
17 rough approximation of natural foveation, made up of several distinct images at different resolutions.  
18 In contrast, resolution changes gradually in natural systems. This may have benefits, but it is not clear  
19 how to arrange such a representation for input to a convolutional network. A circular image with  
20 high magnification at the centre wastes pixels at the corners. A polar representation does not, but it  
21 sacrifices translational equivariance. In summary, while foveation could potentially have benefits for  
22 deep networks, it is not clear when, or how best to implement foveation.

23 To help fill this gap, we compare several foveated downsampling approaches to uniform downsampling  
24 in object recognition. In this context, the different foveated methods perform fairly similarly to each  
25 other, and the best performs slightly better than uniform downsampling (top-1 validation accuracy  
26 48.95% vs. 47.72%; Table 1). Therefore, foveation does not seem to be important for object  
27 recognition (which is unsurprising given the good performance of standard deep networks), but it  
28 does not greatly interfere either. This suggests that foveation could be incorporated into more general  
29 vision systems that perform multiple tasks, such as in robots that must recognize objects and also  
30 read text in the environment.

## 31 2 Methods

### 32 2.1 Network architecture and training

33 We trained deep networks on the ImageNet recognition task, with various kinds of downsampled  
34 images as input. In each case we used a DenseNet-121 [1] network with original hyperparameters

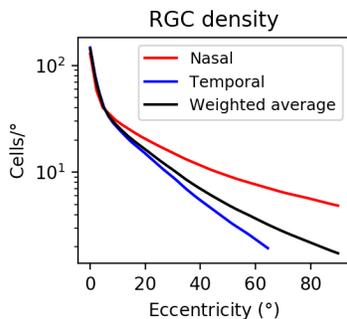


Figure 1: Estimate of retinal ganglion cell (RGC) density as a function of degrees from the fovea. We use estimates from [3], which provides data along the nasal-temporal axis. [4] shows that density is similar in temporal, dorsal, and ventral directions, but higher in the nasal direction. To calculate radially symmetric mean values, we sum nasal and temporal fits from [3] with weights 0.25 and 0.75 (to account for the fact that nasal density is atypical).

35 and training procedure, including random horizontal flips, batch size etc. We trained each network  
 36 for 90 epochs, using SGD (initial learning rate 0.1, reduced by 10x every 30 epochs).

## 37 2.2 Downsampling techniques

38 *Uniform downsampling:* As a baseline method, ImageNet images were uniformly downsampled to a  
 39 32x32 resolution.

40 *Multi-resolution downsampling:* We produced a simple foveated representation composed of four  
 41  $16 \times 16$  downsampled images with different magnifications. The first spanned the whole image, the  
 42 second spanned the central half of the width and height of the image, the third a quarter the width and  
 43 height, and the fourth an eighth. Several past papers have used a similar approach, e.g. [2].

44 *Polar retinal downsampling:* We sampled the image in polar coordinates, creating a rectangular  
 45 image ( $44 \times 23$  pixels) in which the long edge corresponded to the angle and the short edge the radial  
 46 distance from the fovea. The density of samples in the radial direction declined with greater distance  
 47 from the centre. We based the sampling density on retinal ganglion cell (RGC) density (see Figure 1).  
 48 We used gaussian filters with radially increasing widths to reduce artefacts. See example in Figure 2.

49 *Cartesian retinal downsampling:* We sampled the image with the same radially-varying density  
 50 as above, but created a circular image with strong barrel distortion (Figure 3), rather than a polar  
 51 representation. This resulted in a transformation that better retains the translational equivariance  
 52 property of convolutional networks, at the cost of wasting pixels in the corners.

## 53 2.3 Selection of image points to foveate

54 A saliency map was generated for each image with a DeepGaze II model [5]. This map estimated the  
 55 likelihood of a human orienting to each pixel. Human gaze often orients to areas of interest such as  
 56 faces and foreground objects, which often correspond to the target label. We selected the point of  
 57 highest saliency, subject to a constraint that avoided points near image edges (as selecting a point  
 58 near the edge would render much of the crop blank). Specifically, we only chose points around which  
 59 at least 80% of a  $256 \times 256$ -pixel crop would fall within the image boundaries (Figure 4). If the  
 60 resulting crop went outside the image boundaries, we extrapolated by copying edge pixels.

61 We sometimes chose multiple points in a single image. The highest-saliency points are typically close  
 62 together, and contain similar information. To avoid selecting multiple similar points, we modified the  
 63 saliency maps after each selection. Specifically, we subtracted a square-gaussian function from the  
 64 saliency map, with a peak equal to the saliency at the chosen point, and a width of 60 pixels.

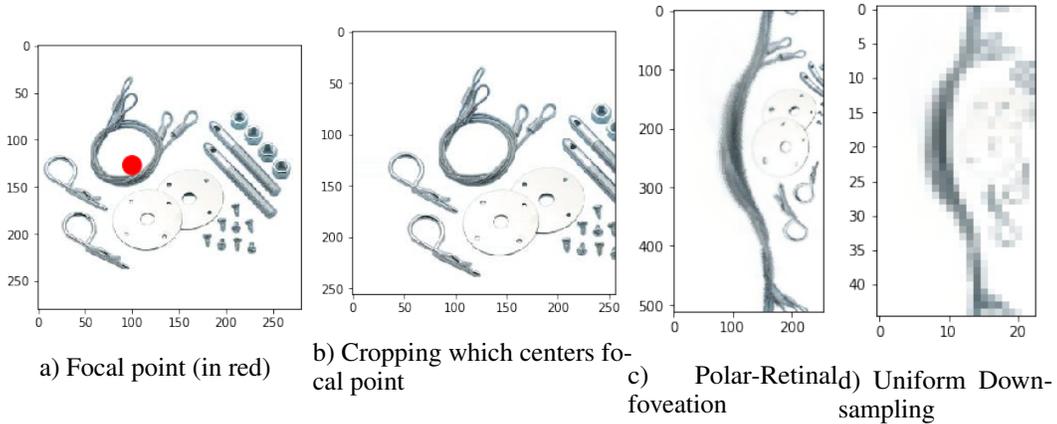


Figure 2: An example of polar-retinal downsampling. (a) The focal point (highest saliency) is determined (red dot). (b) The image is then cropped so the center is at the focal point. (c) The image is then 'foveated' resulting in pixels closer to the center becoming over-represented while pixels close to the edge are under-represented. In this case, the white area on the left of the foveated image is representing the white pixels inside the loop of steel wire of the source image. (d) The result is downsampled uniformly.

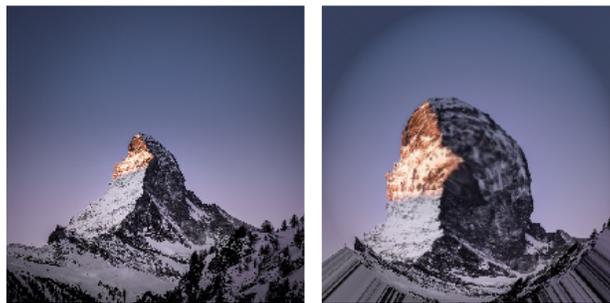


Figure 3: An image before and after cartesian-retinal downsampling. Much like polar foveation, the center of the image is over-represented in the downsample while the extremities are under-represented, proportional to RGC density data.

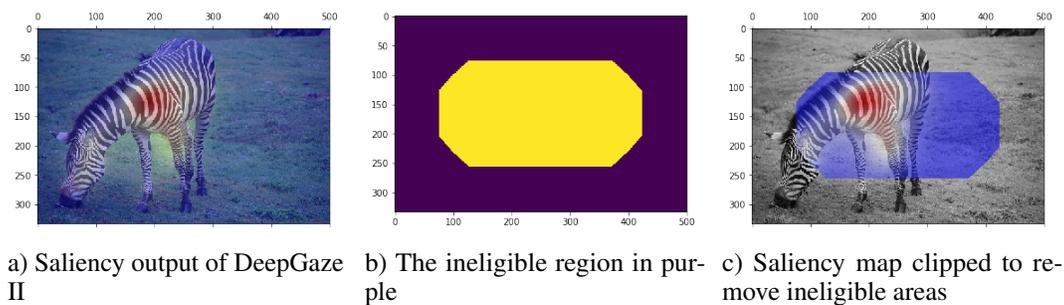


Figure 4: The process of finding a valid saliency map from which the point of highest saliency is chosen. (a) A DeepGaze II model determines a general saliency map. (b) An ineligible region is identified (in purple) where points would result in too much of the resultant crop (20% or more) falling outside the image. (c) The saliency map is clipped and normalized before points are chosen.

### 65 3 Results

66 Figure 5 shows training curves for each of the downsampling methods. During training, each crop  
67 surrounded one of the three most salient points (with sequential updating of the saliency map, as

Table 1: Performance on the validation set

Model	Top 1 Accuracy	Top 5 Accuracy
Most Salient: Uniform	37.66	63.03
Most Salient: Polar-Retinal	33.85	57.50
Top-3 Salient: Uniform	47.72	70.95
Top-3 Salient: Polar-Retinal	47.88	70.26
Top-3 Salient: Cartesian-Retinal	48.95	71.79
Top-3 Salient: Multi-Resolution	47.34	69.91

68 described in the Methods) at random. We also separately trained networks with the uniform and polar  
69 methods using the single most salient crop. Table 1 summarizes validation performance of the trained  
70 models. For Top-3 salient results, predictions were based on three foveations for each image (logits  
71 averaged across foveations).

## 72 4 Discussion

73 The cartesian method performed best in this study. Each of the foveated methods has a limitation that  
74 could potentially be improved in future work. The polar mapping sacrificed translational equivariance  
75 (e.g. the same edge detector could respond to a vertical edge at the bottom of the image and a  
76 horizontal edge at the side). This might be mitigated by adding rotational equivariance. The cartesian  
77 representation wasted pixels at the corners of the image, which limits computational efficiency. Our  
78 version of the multi-resolution representation arranged resolutions side-by-side, which introduced  
79 edge effects. The resolutions could also be treated as separate input channels. We did not do this  
80 because we wanted to hold constant the numbers of parameters and sizes of the representations across  
81 models. Given that foveated views seem not to impair object recognition performance, it would be  
82 interesting to explore potential benefits within more general vision systems.

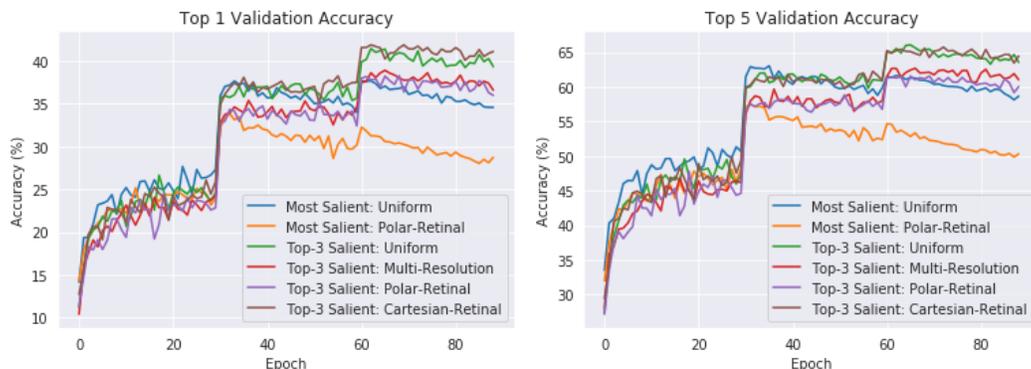


Figure 5: Top 1 (left) and Top 5 (right) validation accuracy during training

## 83 References

- 84 [1] Laurens van der Maaten Kilian Q. Weinberger Gao Huang, Zhuang Liu. Densely connected convolutional  
85 networks. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, 2017.
- 86 [2] Petrucio RT Medeiros, Rafael B Gomes, Esteban WG Clua, and Luiz Gonçalves. Dynamic multifoveated  
87 structure for real-time vision tasks in robotic systems. *J Real-Time Image Processing*, pages 1–17, 2019.
- 88 [3] Röhrenbeck J Boycott BB Wässle H, Grünert U. Cortical magnification factor and the ganglion cell density  
89 of the primate retina. *nature*, pages 643–646, 1989.
- 90 [4] Cowey A Perry VH, Oehler R. Retinal ganglion cells that project to the dorsal lateral geniculate nucleus in  
91 the macaque monkey. *Neuroscience*, pages 1101–1123, 1984.
- 92 [5] Matthias Bethge Matthias Kümmerer, Thomas S. A. Wallis. Deepgaze ii: Reading fixations from deep  
93 features trained on object recognition. 2016.