

# Human sensitivity to perturbations constrained by a model of the natural image manifold

Ingo Freund

Centre for Vision Research and Department of Psychology,  
York University, Toronto, Ontario, Canada



Elee Stalker

Department of Psychology York University,  
Toronto, Ontario, Canada



Humans are remarkably well tuned to the statistical properties of natural images. However, quantitative characterization of processing within the domain of natural images has been difficult because most parametric manipulations of a natural image make that image appear less natural. We used generative adversarial networks (GANs) to constrain parametric manipulations to remain within an approximation of the manifold of natural images. In the first experiment, seven observers decided which one of two synthetic perturbed images matched a synthetic unperturbed comparison image. Observers were significantly more sensitive to perturbations that were constrained to an approximate manifold of natural images than they were to perturbations applied directly in pixel space. Trial-by-trial errors were consistent with the idea that these perturbations disrupt configural aspects of visual structure used in image segmentation. In a second experiment, five observers discriminated paths along the image manifold as recovered by the GAN. Observers were remarkably good at this task, confirming that observers are tuned to fairly detailed properties of an approximate manifold of natural images. We conclude that human tuning to natural images is more general than detecting deviations from natural appearance, and that humans have, to some extent, access to detailed interrelations between natural images.

easily complete missing image components with the most natural looking value (Bethge, Wichmann, & Wichmann, 2007) or detect visual deviations from naturalness (Gerhard, Wichmann, & Bethge, 2013; Wallis et al., 2017; Fründ & Elder, 2013). This precise tuning seems to be mostly restricted to foveal areas and humans are much less sensitive to deviations from naturalness in the periphery (Wallis & Bex, 2012). In fact, two physically different images that match in only a coarse set of image statistics in the periphery, typically appear to be the same (Freeman & Simoncelli, 2011).

Most of these studies are concerned with human sensitivity to deviations from naturalness. It is, however, less clear how to characterize human performance within the manifold of natural images. One challenge is that most experimental manipulations of natural images make the image itself appear less natural. For example, Bex (2010) manipulates images by local deformations. He finds that sensitivity to these deformations is tuned to the spatial frequency at which they occur. Although the deformations used by Bex (2010) only moderately alter the power spectrum of natural images, they considerably disrupt higher order statistical properties such as the phase spectrum of the images (Wichmann, Braun, & Gegenfurtner, 2006). Others, have manipulated images by introducing “dead leaves”—small homogeneous patches—in different locations (Wallis & Bex, 2012) or manipulating the correlation structure of the images (McDonald & Tadmor, 2006). For small patches presented in the visual periphery, observers often can not detect these manipulations. However, we argue that these manipulations still make the image appear less natural rather than studying visual processing within the domain of natural images. A possible solution to this problem is to use selective sampling (Sebastian, Abrams, & Geisler, 2017): Instead of attempting to manipulate the image directly, one chooses another natural image that *by*

## Introduction

The images that we experience in our everyday visual environment are highly complex and our visual system seems to be adapted to perform well with these natural images. However, they only comprise a small part of all possible digital images (Simoncelli & Olshausen, 2001; Geisler, 2008).

Humans are very quick at making simple decisions in natural images (Thorpe, Fize, & Marlot, 2001). We can

Citation: Freund, I., & Stalker, E. (2018). Human sensitivity to perturbations constrained by a model of the natural image manifold. *Journal of Vision*, 18(11):20, 1–15, <https://doi.org/10.1167/18.11.20>.

<https://doi.org/10.1167/18.11.20>

Received May 11, 2018; published October 31, 2018

ISSN 1534-7362 Copyright 2018 The Authors



*chance* shows the desired manipulation. Although this approach guarantees that the resulting “manipulations” remain natural, it is highly dependent on the indexing mechanism that is used to select the manipulated image. The selective sampling approach is dependent on the indexing mechanism used because, for any new feature, a new indexing mechanism would need to be implemented. More importantly, if there aren’t sufficiently many exemplars for a given feature, a new database would be needed. This dependence creates a limitation on generalizing the selective sampling approach to higher levels of visual processing.

Recent advances in machine learning might provide a means to constrain image manipulations to the domain of natural images. Here, we focus on a class of very powerful generative image models, known as generative adversarial nets (GANs; Goodfellow et al., 2014; Radford, Luke, & Chintala, 2016; Arjovsky, Chintala, & Bottou, 2017; Gulrajani, Ahmed, Arjovsky, Dumoulin, & Courville, 2017; Miyato, Kataoka, Koyama, & Yoshida, 2018; Hjelm et al., 2018). GANs learn a mapping, called the *generator*, from an isotropic Gaussian distribution to the space of images. One defining feature of GANs is the use of an auxiliary classification function, often called the *critic*, to judge how good the generator mapping is. Specifically, the critic attempts to predict if a given image has been generated by mapping isotropic Gaussian noise through the generator, or if the image is an instance from the training database. Generator and critic are trained in alternation, where the generator is trained to increase the errors of the critic and the critic is trained to decrease its own error (for example, see Goodfellow et al., 2014, for details). In general, generator and critic can be any possible transformation, but they are typically implemented as artificial neural networks with multiple hidden layers (Goodfellow et al., 2014; Radford et al., 2016). Although never studied quantitatively, images generated from GANs look quite similar to natural images and manipulations in a GAN’s latent space and seem to correspond in a meaningful way to perceptual experience. For example, Radford et al. (2016) start with a picture of a smiling woman, subtract the average latent representation of a neutral woman’s face and add a neutral man’s face to arrive at a picture of a smiling man. Similarly, Zhu, Krähenbühl, Shechtman, and Efros (2016) illustrate that projecting perceptually meaningful constraints back to a GAN’s latent space allows creation of random images with specified features (e.g., edges or colored patches) in the specified locations. Together, these experiences suggest that GANs recover a reasonably good approximation to the manifold of natural images.

In this study, we used a GAN to manipulate generated images and observers made perceptual judgments about these images. In Experiment 1, the

observers viewed a target image and decided which of two noisy comparison images corresponded to that target. Crucially, noise was either applied directly in pixel space or it was restricted to remain within an approximation to the manifold of natural images by applying it in the latent space of a GAN. We found that this task was more difficult if noise was applied in the latent space of a GAN, suggesting that noise in the GAN’s latent space actually changes image features that are relevant for image recognition, while noise that was directly applied in pixel space only resulted in degradation of the image without necessarily changing image content. In Experiment 2, observers were asked to detect changes of direction in videos that were constructed by moving along smooth paths through a GAN’s latent space and we found that observers performed significantly above chance, even for very small directional changes, suggesting that GANs not only recover a good approximation to the manifold of natural images, but that they also recover a perceptually meaningful parameterization of this manifold.

## Experiment 1: Sensitivity to perturbations within the approximate natural image manifold

### Method

#### Training generative adversarial nets

We trained a Wasserstein GAN (Arjovsky et al., 2017) on the 60,000  $32 \times 32$  images contained in the CIFAR10 dataset (Krizhevsky, 2009) using gradient penalty as proposed by Gulrajani et al. (2017). See Figure 1A for example training images. In short, a GAN consists of a generator network  $G$  that maps a latent vector  $\mathbf{z}$  to image space and a critic network  $D$  that takes an image as input and predicts whether that image is a real image from the training dataset or an image that was generated by mapping a latent vector through the generator network (see Figure 2 and Gulrajani et al., 2017, for details of the architecture of the two networks). The generator network and the critic network were trained in alternation using stochastic gradient descent. Specifically, training alternated between five updates of the critic network and one update of the generator network. Updates of the critic network were chosen to minimize the loss

$$\mathbb{E}_{\mathbf{z}}[D(G(\mathbf{z}))] - \mathbb{E}_{\mathbf{y}}[D(\mathbf{y})] + \lambda \|\nabla_{\mathbf{y}} D(\tilde{\mathbf{y}})\|, \quad (1)$$

and updates of the generator were chosen to maximize this loss. Here,  $\mathbb{E}_{\mathbf{z}}$  and  $\mathbb{E}_{\mathbf{y}}$  denote averages over a batch

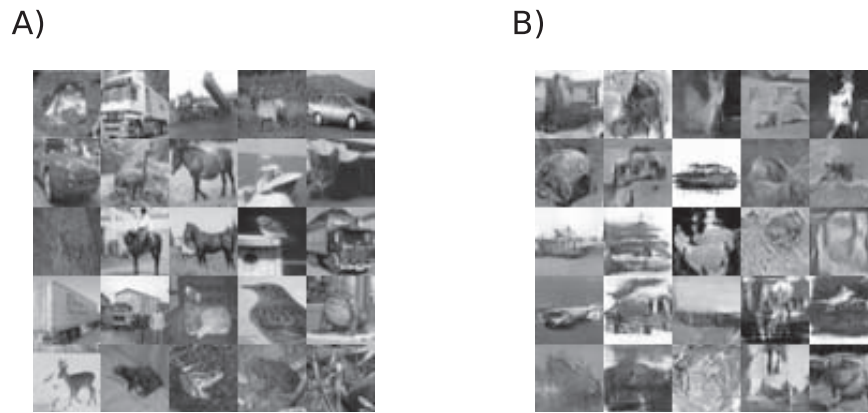


Figure 1. Example training images and samples. (A) Training images from the CIFAR10 database used to train the GAN model. (B) Samples generated from the trained GAN. Similar image samples were used in the experiment.

of 64 latent vectors  $z$  or training images  $y$ , respectively. Furthermore,  $\nabla_y$  denotes the gradient with respect to image pixels  $y$ , which was evaluated at random points along straight line interpolations between real and generated images (see Gulrajani et al., 2017, for details). We set  $\lambda = 10$  during training.

Networks with different numbers of hidden states (parameter  $N$  in Figure 2) were trained for 200,000 epochs using an ADAM optimizer (Kingma & Ba, 2015) with learning rate  $10^{-4}$  and  $\beta_0 = 0$ ,  $\beta_1 = 0.9$ . Specifically, we trained networks with  $N = 40, 50, 60, 64, 70, 80, 90$ , and 128 (see Figure 2). Wasserstein-2 error (Arjovsky et al., 2017) on a validation set (the CIFAR10 test dataset) was lowest, with  $N = 90$  in agreement with visual inspection of sample quality, so we chose a network with  $N = 90$  for all remaining analyses. Example images generated from this final network are shown in Figure 1B.

### Observers

Seven observers participated in the experiment. Two of them were authors, and the remaining five were students from various labs at the Centre for Vision Research at York University, Toronto, Ontario. One additional observer participated in the first session, but withdrew from the experiment afterwards and their data was excluded from the analysis. All observers reported normal or corrected-to-normal vision. Prior to participation, all observers provided written informed consent to participate and all procedures were approved by the York University Ethics Board.

### Procedure

Each individual observer judged image samples from the GAN in a spatial two-alternative forced choice match-to-sample task (see Figure 3). On every trial, the observer saw three images: a target image at the center,

one comparison stimulus on the left, and another comparison stimulus on the right. One of the comparison stimuli was a perturbed version of the target image generated from the GAN, while the other was an equally perturbed version of a separate image generated from the GAN. The observer was required to indicate which of the two comparison images matched the central target image by pressing a corresponding button on a computer keyboard (left arrow key if the left comparison stimulus matched, right arrow key if the right comparison stimulus matched). Stimuli were presented for up to 6 s or until the observer made a response, resulting in practically unlimited viewing time. Before each trial, a fixation point appeared on the screen for 500 ms. Each observer performed five sessions and each session consisted of 80 trials for each noise level and type, resulting in a total of 1,300 trials per observer (except O2, who performed 1,437 trials).

### Stimuli

All stimuli were samples from a GAN, converted to gray scale by averaging the red, green, and blue channels of the sample image. The target stimulus was always noise-free, while the two comparison stimuli were perturbed by one of three noise types (see Figure 3). Pixel noise was constructed by adding independent Gaussian noise to each pixel of the respective image. Fourier noise was constructed in the Fourier domain by replacing the image's phase component by random numbers. This resulted in an image with the same power spectrum as the source image but with completely random phases. A Fourier-perturbed image was constructed by adding a multiple of this power-spectrum-matched noise to the source image. Perturbations of the latent vectors of a GAN are closely correlated to pixel space perturbations ( $r = 0.82$ ,  $p < 10^{-60}$ ), but they are not exactly the same. Therefore, we constructed latent noise by manipulating the latent vector  $z$  from which an



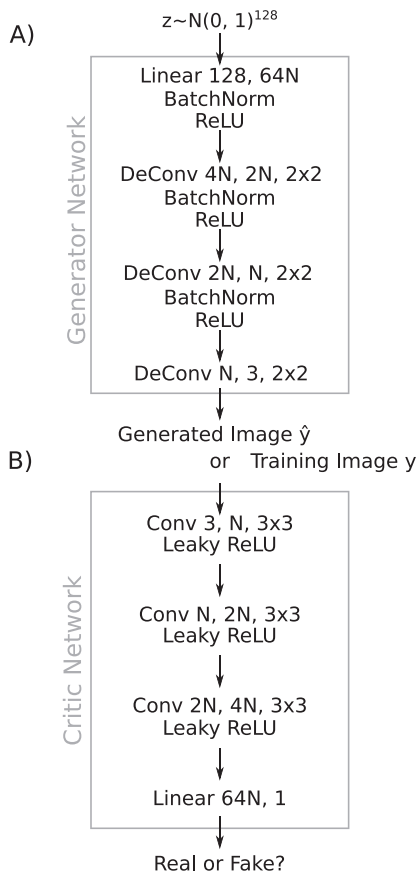


Figure 2. Architecture of the generative adversarial network. (A) Architecture of the generator network. Information flow is from top to bottom, network layers are: “Linear  $k, m$ ”; Affine transformation from  $k$  features to  $m$  features, “Conv  $k, m, n \times n$ ”; Convolutional layer from  $k$  channels to  $m$  channels using a kernel size of  $n \times n$ , “DeConv  $k, m, n \times n$ ”; like convolution but up-sampling before the convolution to increase spatial resolution and image size, “BatchNorm”; Batch normalization (Ioffe & Szegedy, 2015), “ReLU”; and rectified linear unit “ReLU( $x$ ) = max(0,  $x$ )” (Glorot, Bordes, & Bengio, 2011). The generator network maps a sample  $\mathbf{z}$  from an isotropic 128 dimensional Gaussian to a  $32 \times 32$  pixel color image. (B) Architecture of the critic network. Architecture component not used in A is: “Leaky ReLU” (He, Zhang, Ren, & Sun, 2015). The critic network receives as input either an image  $\hat{\mathbf{y}}$  generated by the generator network or a real training image  $\mathbf{y}$ , and it decides if the input image is real or not.

image was generated. To generate perturbed images with a predefined difference in pixel space, we started by adding independent Gaussian noise  $\zeta$  to  $\mathbf{z}$  and determining the corresponding image  $G(\mathbf{z} + \zeta)$ . We then used gradient descent on

$$(\|G(\mathbf{z}) - G(\mathbf{z} + \zeta)\| - \delta)^2 \quad (2)$$

to adjust  $\zeta$  such that the final difference between the target and the perturbed target had a predefined pixel space difference of  $\delta$ .

Stimuli were presented on a medium gray background ( $54.1 \text{ cd/m}^2$ ) on a Sony Triniton Multiscan G520 CRT monitor in a dimly illuminated room. The monitor was carefully linearized using a Minolta LS-100 photometer (Konica Minolta, Ramsey, NJ). Maximum stimulus luminance was  $106.9 \text{ cd/m}^2$ , minimum stimulus luminance was  $1.39 \text{ cd/m}^2$ . If the nominal stimulus luminance exceeded that range, it was clipped (for subsequent analyses, we also used the clipped stimuli). On every frame, the stimuli were re-rendered using random dithering to generate a quasi-continuous luminance resolution (Allard & Faubert, 2008). At a viewing distance of approximately 87 cm, each stimulus image subtended approximately  $0.65^\circ$  of visual angle and were separated by approximately  $0.13^\circ$  of visual angle. One pixel subtended approximately  $0.02^\circ$  of visual angle.

### Data analysis

For every observer, we estimated a psychometric function parameterized as

$$\begin{aligned} \text{prob}(\text{correct}|x) &= \psi(x) \\ &= \gamma + (1 - \gamma - \lambda)\sigma(a + bx), \end{aligned} \quad (3)$$

where  $\gamma = 0.5$  is the probability to guess the stimulus correctly by chance,  $\lambda$  is the lapse probability,  $\sigma$  is the logistic function and  $a$  and  $b$  govern the offset and the slope of the psychometric function. Here,  $x$  is the root-mean-square level of noise applied to perturb the respective images in dB relative to the screen’s background luminance. However, we note that different ways of scaling the noise (other than dB) did not impact our main results. We adopted a Bayesian perspective on estimation of the psychometric function (Fründ, Haenel, & Wichmann, 2011) and used weak priors  $\lambda \sim \text{Beta}(1.5, 20)$ ,  $a \sim N(0, 100)$ ,  $b \sim N(0, 100)$ , where  $a$  and  $b$  are expressed on the dB scale of the noise. Mean a posteriori estimates of the critical noise level  $x_c$  at which  $\psi(x) = 0.75$  and the slope of the psychometric function at  $x_c$  were obtained using numerical integration of the posterior (Schütt, Harmeling, Macke, & Wichmann, 2016).

To understand how the structure of the GAN’s latent space determined image matching performance, we reanalyzed data from the latent condition. For this reanalysis, we assumed that observers would pick the perturbed stimulus that is closer to the target with respect to some distance measure. More specifically, let  $\mathbf{t}$  denote the noise-free target stimulus and  $\tilde{\mathbf{t}}$  and  $\tilde{\mathbf{d}}$  denote the perturbed target and distractor stimuli, respectively. If an observer is picking the stimulus that is closer to the target, then  $c := d(\tilde{\mathbf{d}}, \mathbf{t}) - d(\tilde{\mathbf{t}}, \mathbf{t})$  should show a positive correlation with the observer’s trial by trial response accuracy. Here,  $d$  is a suitably defined distance measure. We used either the Euclidean

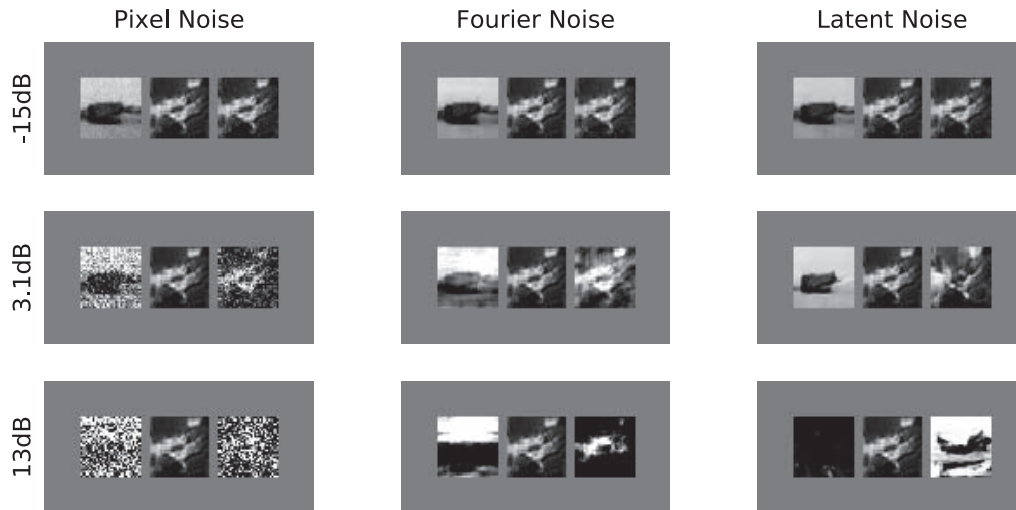


Figure 3. Design of the image-matching experiment. On every trial, the observer saw three images: a target image at the center, one comparison stimulus on the left, and another comparison stimulus on the right. Comparison stimuli were sampled randomly from the GAN and perturbed by different types of noise. Pixel noise (left column) was added as independent Gaussian noise to every pixel; Fourier noise (middle column) with the same power spectrum as the original image but with random phases was added to the original image; latent noise (right column) was independent Gaussian noise applied in the latent space of the GAN that was used to generate the images. The top row shows example experimental displays with low noise, the central row shows example experimental displays close to the critical noise level for latent stimuli, and the bottom row shows example experimental displays with high noise (amounts given on left). For illustration, the left stimulus is a random perturbed stimulus and the right stimulus is the perturbed target in every example display. During the experiment, the identities and locations of the comparison stimuli were randomized.

distance  $d_{\text{Euclidean}}(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_i (x_i - y_i)^2}$ , the radial distance  $d_{\text{radial}}(\mathbf{x}, \mathbf{y}) = \sqrt{(\|\mathbf{x}\| - \|\mathbf{y}\|)^2}$ , or the cosine distance  $d_{\text{cosine}}(\mathbf{x}, \mathbf{y}) = 1 - \langle \mathbf{x}/\|\mathbf{x}\|, \mathbf{y}/\|\mathbf{y}\| \rangle$ , where  $\|\mathbf{x}\|$  denotes the Euclidean norm of a vector  $\mathbf{x}$  and  $\langle \mathbf{x}, \mathbf{y} \rangle$  denotes the scalar product of vectors  $\mathbf{x}$  and  $\mathbf{y}$ . These distances were applied in either the GAN's latent space or directly in pixel space, after concatenating the respective stimulus' pixel intensities into one long vector. We then determined receiver operating curves (ROC) for predicting correct versus incorrect responses based on  $c$ . The area under the ROC is a measure for how well the respective distance measure predicts the observer's trial by trial responses (Green & Swets, 1966). To test if the area under the curve (AUC) was significantly different from chance, we performed a permutation test randomly reshuffling the correct/incorrect labels 1,000 times and taking the 95th percentile of the resulting distribution as the critical value. We also used permutation tests to determine if the AUC for two different distance measures was significantly different. For the pairwise comparisons, there are 128 possible reassignments of AUC values to the two conditions, and we computed all of them. The  $p$  values for these post hoc comparisons were corrected for multiple comparisons to control for inflation of the false discovery rate (Benjamini & Hochberg, 1995).

To gain insight into the image features that determined the observers' responses, we applied the ROC analysis to a number of image features as well.

Firstly, we calculated the luminance histogram (50 bins) for each image and calculated the distance difference  $c$  between luminance histograms of the respective images. Distance measures were computed over vectors of length 50, and each entry denotes the bin count. Secondly, to determine local dominant orientation at each pixel we first filtered the image with horizontal and vertical Scharr filters (Scharr, 2000) as implemented in scikit-image (van der Walt et al., 2014) giving local horizontal structure  $h$  and vertical structure  $v$ . The local orientation  $\phi$  was extracted from these two responses such that  $h = r \cos(\phi)$  and  $v = r \sin(\phi)$ , where  $r = \sqrt{h^2 + v^2}$ . We then determined the histogram (3 bins) of the local orientations across the image and calculated  $c$  as the distance difference between these orientation histograms. As a third feature, we calculated the edge densities of the two images by using the canny edge detector from scikit-image with a standard deviation of 2 pixels and calculating the fraction of pixels labeled as edges by this algorithm. As a fourth feature we determined the slope of the power spectrum in double logarithmic coordinates.

Finally, we used a standard method for image segmentation (Felsenwalb & Huttenlocher, 2004) to calculate segmentations of the images  $\mathbf{t}$ ,  $\mathbf{\tilde{t}}$ , and  $\mathbf{\tilde{d}}$ . We used the method implemented in scikit-image (van der Walt et al., 2014). Briefly, this algorithm iteratively merges neighboring pixels or pixel groups if the differences across their borders are small compared to

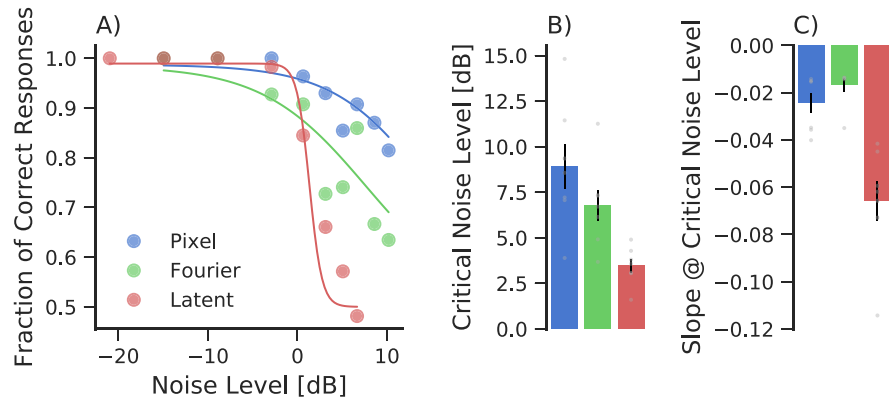


Figure 4. Higher sensitivity to noise perturbations applied within the manifold recovered by a GAN. (A) Psychometric function for one example observer. Each dot represents between 50 and 60 trials, solid lines are mean a posteriori estimates of the psychometric function (see Methods). All noise levels were quantified as root mean square difference to the target in pixel space and were normalized to the background luminance. (B) Average critical noise levels corresponding to 75% correct performance. Height of the bars denotes the mean across seven observers, error bars indicate standard error of measurement across observers. Light gray dots indicate results for individual observers. Colors are the same as in part A. (C) Average slope of the psychometric function at the critical noise level. Height of the bars denotes the mean across seven observers, error bars indicate standard error of measurement across observers. Light gray dots indicate results for individual observers. Colors are the same as in part A.

the differences within them. Each segmentation consists of a number of discrete labels assigned to the pixels of the original image. If two pixels belong to the same segmented region, the two labels associated with them should be the same. However, different segmentations may assign different labels to the same region. Thus, two segmentations  $s$  and  $\tilde{s}$  would be similar, if for many pairs  $(i, j)$  of pixels  $s_i = s_j$  implies  $\tilde{s}_i = \tilde{s}_j$ . When calculating the distance between two segmentations  $s$  and  $\tilde{s}$ , we therefore count the number of pixel pairs for which  $s_i = s_j$  and  $\tilde{s}_i = \tilde{s}_j$  and we normalize by the number of pixel pairs that are assigned to the same region by at least one of the two segmentations. This is captured by the distance measure

$$d_{\text{segm}}(s, \tilde{s}) = 1 - \frac{\sum_{ij} \mathbb{I}(s_i = s_j \text{ and } \tilde{s}_i = \tilde{s}_j)}{\sum_{ij} \mathbb{I}(s_i = s_j \text{ or } \tilde{s}_i = \tilde{s}_j)}, \quad (4)$$

where  $\mathbb{I}(A) = 1$  if the expression  $A$  is true and  $\mathbb{I}(A) = 0$  otherwise and the sums go over all pairs of pixels  $i, j$ . If the two segmentations define exactly the same regions (but possibly with different labels),  $d_{\text{segm}}$  will be 0, if the two segmentations are completely different, in the sense that one has only one region (the entire image) and the other assigns each pixel to its own region, then  $d_{\text{segm}}$  will be 1.

## Results

### Lower tolerance for noise applied within an approximation of the natural image manifold

In the first experiment, seven observers were required to judge which one of two noisy comparison images

matched a centrally presented target stimulus (see Figure 3). Figure 4A shows psychometric functions for one example observer as a function of noise level. In general, higher noise levels were associated with less correct responses. The observer's performance, indicated by level of response correctness, was least affected by noise that was applied independently to each pixel (critical noise level with 75% correct performance at  $14.8 \pm 1.43$  dB, posterior mean and standard deviation, blue dots and line in Figure 4A). Observer performance was more affected by noise that was applied in the pixel domain but matched the power spectrum of the original image (critical noise level at  $7.31 \pm 0.54$  dB, green dots and line in Figure 4A). Finally, the level of response correctness was most affected by noise that was applied in the latent space of the GAN and thus stayed within an approximation of the manifold of natural images (critical noise level at  $1.60 \pm 0.60$  dB, red dots and line in Figure 4A). Thus, this observer's level of response correctness was most affected by noise applied within a model of the manifold of natural images. Furthermore, the psychometric function for this observer was considerably steeper in the latent noise condition (slope at critical noise level was  $-0.11 \pm 0.09/\text{dB}$ ) than in the other two conditions (slopes at critical noise level were  $-0.015 \pm 0.006/\text{dB}$  for pixel noise and  $-0.013 \pm 0.0015/\text{dB}$  for Fourier noise), suggesting that the observer was relatively insensitive to low amplitude latent noise and then abruptly became unable to do the task, consistent with the idea that for latent noise above a certain level a categorical change happens, while noise in the pixel domain results in a more gradual decrease in image quality.

On average across all seven observers, the critical noise level was highest for independent pixel noise ( $8.91$

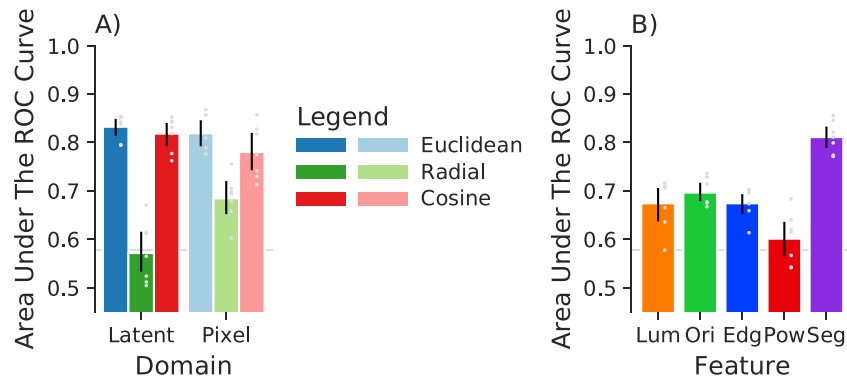


Figure 5. Trial-by-trial performance explanation. (A) Area under the ROC curve for different hypothetical distances applied in latent space or pixel space. Distance measures applied in the GAN’s latent space, corresponded to distances within a model of the manifold of natural images (darker colors). Color hue codes for the type of distance measure used. Error bars indicate 95% bootstrap confidence intervals. The light gray line marks the average critical value for chance performance. (B) Same as A but for image features. “Lum” corresponds to comparisons of luminance histograms, “Ori” corresponds to comparisons of orientation histograms, “Edg” corresponds to edge density, “Pow” corresponds to the slope of the image’s power spectrum, and “Seg” corresponds to comparisons of image segmentations (see Methods for detail).

$\pm 1.22$  dB,  $M \pm SEM$ ). It was comparable for Fourier noise,  $6.77 \pm 0.83$  dB and paired  $t$ -test pixel versus Fourier noise,  $t(6) = -1.59$ ,  $ns$ ), and it decreased significantly for latent noise,  $3.47 \pm 0.36$  dB, paired  $t$ -test pixel versus latent noise,  $t(6) = 3.59$ ,  $p = 0.011$ , and Fourier versus latent noise,  $t(6) = 3.89$ ,  $p = 0.0080$ , respectively (see Figure 4B). Thus, overall observers were most affected by noise that was approximately applied within the manifold of natural images by perturbing the GAN’s latent representation of the stimulus. We verified that this result also held for every individual observer. We further found that psychometric functions tended to fall off more steeply when noise was applied in the GAN’s latent space (average slope at critical noise level for latent noise was  $-0.066 \pm 0.0084$ /dB; see Figure 4C) than when noise was applied in pixel space (average slope at critical noise level for pixel noise was  $-0.024 \pm 0.0041$ /dB, for Fourier noise  $-0.017 \pm 0.0027$ /dB), replicating the observations from Figure 4A.

To summarize, we found that, in general, noise that was approximately applied within the manifold of natural images was more effective at interrupting observers’ performance than noise applied outside of the manifold of natural images in pixel space. We performed two additional sets of analysis to determine (a) if latent space or pixel space image differences were more predictive of observer’s trial-by-trial behavior, and (b) which image features were responsible for the decline in performance in the latent noise condition.

#### **Distance in latent space correlates with image matching performance**

We modeled observer behavior in the image-matching task, by assuming that on every trial, the observer

picks the comparison stimulus that appears closer to the target with respect to some appropriate distance measure. We looked at three different candidate distance measures and applied them in both the GAN’s latent space and directly in pixel space. To evaluate the relevance of each distance measure, we nonparametrically calculated area under the receiver operating curve (AUC) values for discrimination between correct and incorrect trials.

Figure 5A shows average AUC values for different hypothetical distance measures applied either in latent space or in pixel space. The simplest distance measure is the Euclidean distance, marked by the blue bars in Figure 5A (darker blue for Euclidean distance in latent space, lighter blue for Euclidean distance in pixel space). Note that Euclidean distance in pixel space is equivalent to the standard deviation of the difference image between stimuli that was used as a unified measure of perturbation strength in Figure 4 (without the transformation to dB). Euclidean distance in latent space was at least as predictive as Euclidean distance in pixel space (latent space,  $0.83 \pm 0.0092$ ; average AUC  $\pm SEM$ ; pixel space,  $0.82 \pm 0.014$ ; and permutation test  $p = 0.17$ ).

We performed the same analysis using the difference between the norms of either the latent vector or the pixel vector (green bars in Figure 5A). In pixel space, radial distance is equivalent to an observer who simply compares the RMS contrast of the images. Counter to the Euclidean distance, this observer would *first* take the standard deviation of the target and flanker images and then compare differences in standard deviations. In latent space, the norm of the latent vector seems to be related to contrast as well, but the relationship is more complex. Radial distance receives considerably lower AUC than Euclidean distance in both latent and pixel



space and was a much less reliable predictor of trial-by-trial performance. Radial distance in latent space was significantly less predictive than radial distance in pixel space (latent space,  $0.57 \pm 0.022$ ; pixel space,  $0.68 \pm 0.016$ ; and permutation test  $p < 0.05$  corrected) and for four out of seven observers, radial distance in latent space did not predict trial-by-trial choices significantly better than chance.

Finally, we analyzed how well cosine distance explained observers' responses (red bars in Figure 5A). Cosine distance is interesting for two reasons. First, cosine distance is equivalent to Euclidean distance except for the influence of the radial component. Second, cosine distance is closely related to correlation and an observer who relies on cosine distance essentially uses the target stimulus as a template and computes the correlation with each of the comparison stimuli to pick the stimulus that correlates best with that template.<sup>1</sup> Cosine distance applied in latent space was a better predictor than if it was applied in pixel space (latent space,  $0.82 \pm 0.012$ ; pixel space,  $0.78 \pm 0.019$ ; and permutation test  $p = 0.05$  corrected).

In the Appendix, we replicate these results using an analysis based on logistic regression.

### ***Distortions of mid-level features explain trial by trial performance***

In order to determine which image features were responsible for the decline in performance with perturbations in the latent space of GANs, we applied the same analysis to different image features (see Figure 5B).

We found that differences in the luminance distribution of the images are clearly predictive of trial-by-trial behavior (average AUC,  $0.67 \pm 0.018$ ; AUC was larger than 95th percentile of the null distribution in all seven observers). Yet, other features such as the difference in local orientation (average AUC,  $0.70 \pm 0.010$ ) or differences in the images' edge density (average AUC,  $0.67 \pm 0.011$ ) were equally good predictors of the observers' trial-by-trial behavior (permutation test not significant after correction for multiple comparisons).

One of the quantities that might have been relevant for our observers is the slope of the power spectrum (see for example Alam, Vilankar, Field, & Chandler, 2014). We evaluated to what extent this feature contributed to our observers' decisions and found that it is largely irrelevant at explaining the observers' trial-by-trial behavior: In three out of seven observers the AUC of this feature was not significantly different from chance performance and the average AUC for the slope of the images' power spectrum was significantly less than that for any other feature we studied.

In order to quantify how well differences in the mid-level structure of the perturbed images could explain

trial-by-trial responses, we calculated segmentations of all images using a standard segmentation algorithm (Felzenszwalb & Huttenlocher, 2004). As this feature is less common than the other image features, we show the distribution of this image feature in the Appendix. Although we do not believe that humans necessarily segment images using graph based optimization as the algorithm does, we believe that this approach provides at least a coarse approximation to the mid-level structure of the images. Differences in segmentation were considerably more predictive than differences in any other feature distribution (mean AUC for segmentation  $0.82 \pm 0.011$ , permutation test  $p < 0.05$  corrected). In fact, differences in segmentation were about as predictive of trial-by-trial behavior as cosine distance in latent space (permutation test  $p = 0.60$ ), suggesting that indeed distortions of the images' mid-level structure might be responsible for the decline in image-matching performance when noise was constrained to stay within the recovered manifold of natural images by applying it in the GAN's latent space.

Taken together, the results of image features and hypothetical distance measures suggest that the latent space of GANs captures processing beyond simple contrast differences. Specifically, trial-by-trial responses were most accurately explained by differences in image segmentations rather than low-level features such as luminance, orientation, or edge density.

## **Experiment 2: Sensitivity to directions in the recovered natural image manifold**

In Experiment 1, we found that human observers are particularly sensitive to image perturbations that stay within an approximation of the manifold of natural images. This was achieved by perturbing stimuli along a parameterization of the manifold of natural images as recovered by a GAN. We wondered if observers were also sensitive to other aspects of this parameterization, such as for example direction. We therefore asked observers to discriminate between videos that were created by walking along either straight paths in latent space (i.e., that contained no change in direction) or paths that had a sudden turn (i.e., that contained a change in direction).

### **Method**

#### ***Observers***

Five observers participated in the second experiment and the control condition for the second experiment.



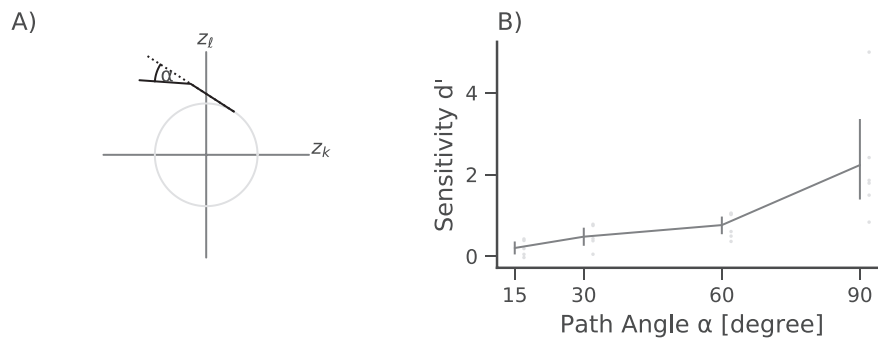


Figure 6. Observers can detect turns in paths along the manifold of natural images as recovered by a GAN. (A) Illustration of stimulus construction. Two latent axes  $z_l$  and  $z_k$  are shown for illustration, the actual latent space had a dimension of 128. The light gray circle marks the sphere of radius 10 in latent space. Shown are a straight path (dotted line, standard stimulus in the experiment) and a path that turns half way at an angle  $\alpha$  (solid line, target stimulus in the experiment). (B) Average sensitivity for detecting a turn in an image sequence along a path along the modeled manifold. Error bars mark 95% confidence intervals. Light gray dots mark individual observer results.

All five of them had participated in the first experiment as well. The procedures were approved by the Ethics Board of York University, and observers provided informed consent before participating. One observer (O3) accidentally did one session with incorrect angle logging. As the correct angles could not be recovered, we decided to exclude the corresponding trials from the analysis.

### Procedure

The experiment was a single interval design. On every trial, the observer either saw a video without a turn in latent space direction or a video that contained a turn in latent space direction and they had to decide if the presented video contained a turn or not. The probabilities for turn and no-turn videos were each 0.5. In order to avoid bias about the image features that would indicate a turn in latent space, the observer was instructed that there would be two classes of videos and that we were not able to describe the difference unambiguously in words. Instead, the observer first saw 100 trials with path angles of  $90^\circ$  (vs. straight) and received trial-by-trial feedback about their performance. After that, the observer performed two blocks of 100 trials for each path angle for the main experiment. For each block, the size of the possible path angle was kept constant. To allow the observer to calibrate their decision criterion to the size of turns presented in the respective block, we provided trial-by-trial feedback during the first 20 trials of each block and only analyzed the remaining 80 trials. Thus, in total we analyzed 160 trials per observer per path angle.

### Stimuli

Each video consisted of 60 frames. The first frame corresponded to a random point on the sphere with

radius 10. Successive frames were then constructed by taking steps of norm 0.5 in a random direction tangential to the sphere with radius 10 (see Figure 6A). Straight paths were constructed by simply taking 60 successive steps in the same direction. Paths with a turn were constructed by changing the direction of steps by an angle  $\alpha$  of  $15^\circ$ ,  $30^\circ$ ,  $60^\circ$ , or  $90^\circ$  after the first 30 frames. We will refer to this angle as the *path angle* in the following. At a frame rate of 60 Hz, each video had a duration of 1s and if the video contained a turn in latent space, that turn happened after 500 ms. Otherwise, the setup for Experiment 2 was the same as in Experiment 1.

For the control experiment, additional videos were constructed for paths through pixel space, through Fourier space, and through latent space. Videos for paths through pixel space were constructed in the same way as described above with the only difference that the respective vectors were sampled from an isotropic Gaussian. Videos for paths through Fourier space were similar to videos constructed in pixel space, but were additionally filtered to an approximate  $1/f$  power spectrum by multiplying their Fourier transform with  $1/(0.1 + f)$ . The exact size of the root mean square difference between successive video frames varied somewhat from trial to trial, but there were no statistically significant differences between frame by frame differences in pixel space, in Fourier space, or in latent space ( $p > 0.1$ ). Animated examples are available as supplementary material (Supplementary Files S1–S18).

### Data analysis

To evaluate observers' sensitivity to path angles, we calculated  $d'$  at each path angle. For single observers, confidence intervals for  $d'$  were determined by bootstrap with 1,000 samples.

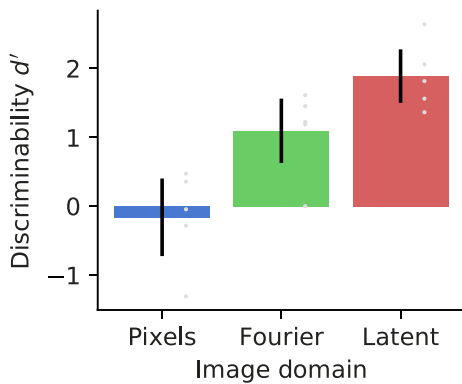


Figure 7. Turns in other image spaces are harder to detect. Average sensitivity for detecting a  $90^\circ$  turn in an image sequence along a path through pixel space, through Fourier space, and through latent space. Error bars mark 95% confidence intervals and gray dots mark performance of individual observers.

## Results

Figure 6B shows average discrimination performance as a function of path angle. Not surprisingly, observers were best at detecting turns of  $90^\circ$  in latent space (average  $d' = 2.23 \pm 0.59$ ). However, even for path angles of  $15^\circ$ , the smallest path angles tested, three out of five observers performed above chance (average  $d' = 0.20 \pm 0.090$ ; one-sided  $t$ -test against zero,  $t(4) = 2.25$ ;  $p = 0.0436$ ). This indicates that even small changes in latent space direction were detected by the observers.

Overall, observers were remarkably good at discriminating between different directions taken by movies in latent space. In fact half of the observers responded correctly in more than 80% of the trials (average  $78.1 \pm 2.56\%$ ,  $M \pm SEM$ ) at the largest turn.

Is sensitivity to these videos really specific for the latent space of a generative adversarial network? To control for the possibility that observers could simply detect any abrupt change in continuously morphed images, we performed a control experiment in which we used only  $90^\circ$  turns but varied the space in which the videos corresponded to straight paths. We found that observers were insensitive to turns in pixel space paths (average  $d' = -0.16 \pm 0.29$   $M \pm SEM$ ; Figure 7) and were moderately sensitive to turns in Fourier space paths (average  $d' = 1.09 \pm 0.24$ ). However, sensitivity to turns in paths through the GAN's latent space was highest (average  $d' = 1.89 \pm 0.20$ ). The pairwise comparisons between performance for pixel space versus Fourier space paths ( $p < 0.05$ , permutation test) and performance for Fourier space paths versus latent space paths ( $p < 0.05$ , permutation test) were both statistically significant with the observed configuration having the largest difference of all permutations. This

confirms that sensitivity to directions in latent space is a specific property of the representation recovered by generative adversarial networks.

## General discussion

We found that observers are sensitive to image manipulations that are restricted to remain within an approximation of the manifold of natural images. Errors made within the manifold of natural images recovered by a GAN seem to be related to changes in the configural structure of the images more than to changes in local image features. We further find that observers are sensitive not only to noise within the model of the manifold, but also to more subtle aspects of paths along the model of the manifold (see Experiment 2).

## Global versus local image models

Our results seem to contradict reports that find that observers are very sensitive to deviations from naturalness (Gerhard et al., 2013): How would observers be equally sensitive to changes within the manifold of natural images and deviations away from that manifold? We find that sensitivity to perturbations within an approximation of the manifold of natural images can be well predicted by the part and object structure of the image that is captured by standard segmentation algorithms (see Figure 5), while deviations from naturalness might be more related to the fine structure of local statistics. This is also visible in Figure 1: The two image classes in Parts A and B appear very similar at first glance and only upon closer inspection does it become clear that the examples in Part B are mostly meaningless. We believe that this point could be taken both as an advantage and a disadvantage. Clearly, the samples in Figure 1B do not match every aspect of the natural images in Figure 1A (although better matches can be achieved if the GAN is restricted to more narrow classes of objects; for example, see Radford et al., 2016; Gulrajani et al., 2017). However, the images capture a lot of the global and highly nonstationary properties of natural images, that texture models based on stationarity do not capture (Portilla & Simoncelli, 2000; Gatys, Ecker, & Bethge, 2015). We therefore believe that our approach is complementary to the local approach taken by studies that investigated texture processing (e.g., Freeman & Simoncelli, 2011; Gerhard et al., 2013; Wallis, Bethge, & Wichmann, 2016; Wallis et al., 2017).

## Small images

The images employed here are relatively small. Our images were only  $32 \times 32$  pixels in size. In contrast, Wallis et al. (2017) used image patches that were  $128 \times 128$  pixels to compare between texture images created from a deep neural network model and real photographs of textures. Other studies have used a range of image sizes (in increasing order of image size: Alam et al., 2014; Sebastian et al., 2017; Bex, 2010), but our images are closer to the range of image sizes used as patches of images (Gerhard et al., 2013) rather than entire images. However, training generative adversarial networks on larger images with similar image variability as the CIFAR10 network currently typically requires training class conditional networks as for example done by Miyato et al. (2018) when training on the entire ImageNet dataset (Russakovsky et al., 2015). Although this would have been possible in principle in this study as well, it would have implied that separate manifolds would be used for different classes and it would have made interpretation of our results considerably more complex. We therefore decided to restrict ourselves to smaller images.

## Dataset bias

Many publicly available data bases appear to be systematically biased (Wichmann, Drewes, Rosas, & Gegenfurtner, 2010): The images in these databases are presegmented in the sense that a photographer selected a viewpoint that they considered particularly “interesting,” or in that they selected which objects to put in focus. Although pictures from these databases may appear natural, conclusions drawn from these data bases may be misleading (Wichmann et al., 2010). The example pictures in Figure 1A clearly show such photographer bias. In every one of these examples, the perspective is clearly focused on one specific object, while typical natural scenes often contain multiple objects and a lot of not explicitly defined random texture (i.e., background). In general, the extent to which such bias would influence our results depends on the extent to which the biased pixel-by-pixel statistics in these images determine responses in our task. In the section Distortions of mid-level features explain trial-by-trial performance, we show that the main determinant of responses in our experiments was the segmentation structure of our images. Although these might also be influenced by photographer bias, we assume that the impact is more subtle. Furthermore, eye movements tend to center objects on or close to the fovea (Kayser, Nielsen, & Logothetis, 2006; ‘t Hart, Schmidt, Roth, & Einhauser, 2013), which might lead to similar global effects as a photographer focusing the

image on selected objects. However, the extent to which photographer bias impacts the image representation recovered by GANs is still an open question.

## Non-object images

Upon closer inspection, the examples in Figure 1B do not look exactly like real objects. Although each one of the examples can clearly be segmented into foreground and background, it is not always possible to actually name the foreground objects in Figure 1B; this seems to be easier for the training examples in Figure 1A, despite the relatively low resolution of the images. Thus, the samples from the GAN used in the present study would probably be easy to discriminate from real images if they were directly compared to real images (Gerhard et al., 2013; Wallis et al., 2017). It should thus be noted that the image representation learned by the GAN used in this study is only an approximation to the manifold of natural images (note however, that other studies training on the CIFAR10 dataset show samples of similar quality, e.g., Gulrajani et al., 2017; Roth, Lucchi, Nowozin, & Hofmann, 2017). Although image manipulations in this approximation appear to be convincing if the GAN has been trained on more restricted sets of training images (see for example Zhu et al., 2016), this cannot be guaranteed for all the stimuli used in this experiment. To this date, it is unclear how exactly the GAN samples used in this study match the perceived properties of the training data (even though the training data themselves might be biased; see section Dataset bias).

We believe however, that natural-appearing non-objects are still an interesting class of stimuli. For example, Huth, Nishimoto, Vu, and Gallant (2012) report that semantic content is important for shaping the response properties of large parts of anterior visual cortex, suggesting that many areas that are traditionally thought of as visual, are also semantic areas. Attempts to further test this claim have been restricted to correlational approaches (Khaligh-Razavi & Kriegeskorte, 2014), partly because it is difficult to generate non-object stimuli that otherwise fully match the properties of object stimuli (see also Fründ et al., 2008). Comparing responses to training images with objects (Figure 1A) and with images generated from a GAN but without the full semantic information (Figure 1B) could help resolve this point.

## Conclusion

In this study, we explored the potential of studying vision within an approximate manifold of natural images. To do so, we employed a generative adversarial network to constrain perturbations to remain within the



manifold of natural images and we find that observers are remarkably sensitive to image manipulations that are constrained in this way. We observe that perturbations within a model of the manifold of natural images tend to disrupt more global image structures such as figure–ground segmentation structure. This might prove useful in future studies that investigate such processes under more naturalistic conditions. The fact that GANs provide an approximate parametrization to the manifold of natural images encourages further use of these powerful image models to study vision under complex naturalistic stimulus conditions.

All code and data for this manuscript is available online at <http://doi.org/10.5281/zenodo.1308853>.

*Keywords:* natural images, image recognition, artificial neural networks, generative adversarial nets, noise perturbations

## Acknowledgments

This work was supported by a York University Bridging Grant and a York University Combined Minor Research Grant/Junior Faculty Fund.

Commercial relationships: none.

Corresponding author: Ingo Fruend.

E-mail: [ifruend@yorku.ca](mailto:ifruend@yorku.ca).

Address: Centre for Vision Research and Department of Psychology York University, Toronto, Ontario, Canada.

## Footnote

<sup>1</sup> Note that for cosine distance in latent space, linear template matching would be performed in the GAN's latent space, which corresponds to image space in a complex and highly nonlinear way.

## References

- Alam, M. M., Vilankar, K. P., Field, D. J., & Chandler, D. M. (2014). Local masking in natural images: A database and analysis. *Journal of Vision*, *14*(8):22, 1–38, <https://doi.org/10.1167/14.8.22>.
- Allard, R., & Faubert, J. (2008). The noisy-bit method for digital displays: Converting a 256 luminance resolution into a continuous resolution. *Behavior Research Methods*, *40*(3), 735–743.
- Arjovsky, M., Chintala, S., & Bottou, L. (2017). *Wasserstein GAN*. <https://arxiv.org/abs/1701.07875>.
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, *57*(1), 289–300.
- Bethge, M., Wiecki, T. V., & Wichmann, F. A. (2007). The independent components of natural images are perceptually dependent. In B. E. Rogowitz, T. N. Pappas, & S. J. Daly (Eds.), *Proceedings of SPIE, Human Vision and Electronic Imaging XII: Vol. 6492* (p. 64920A). Bellingham, WA: SPIE.
- Bex, P. J. (2010). (In) Sensitivity to spatial distortion in natural scenes. *Journal of Vision*, *10*(2):23, 1–15, <https://doi.org/10.1167/10.2.23>.
- Felsenzwalb, P. F., & Huttenlocher, D. P. (2004). Efficient graph-based image segmentation. *International Journal of Computer Vision*, *59*(2), 167–181.
- Freedman, D., & Diaconis, P. (1981). On the histogram as a density estimator: L2 theory. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, *57*(4), 453–476.
- Freeman, J., & Simoncelli, E. (2011). Metamers of the ventral stream. *Nature Neuroscience*, *14*, 1195–1201.
- Fründ, I., Busch, N. A., Schadow, J., Gruber, T., Korner, U., & Herrmann, C. S. (2008). Time pressure modulates electrophysiological correlates of early visual processing. *PLoS One*, *3*(2):e1675.
- Fründ, I., & Elder, J. (2013). Statistical coding of natural closed contours [Abstract]. *Journal of Vision*, *13*(9), 119–119.
- Fründ, I., Haenel, N. V., & Wichmann, F. A. (2011). Inference for psychometric functions in the presence of nonstationary behavior. *Journal of Vision*, *11*(6):16, 1–19. <https://doi.org/10.1167/11.6.16>.
- Gatys, L. A., Ecker, A. S., & Bethge, M. (2015). Texture synthesis using convolutional neural networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems 28* (pp. 262–270). Red Hook, NY: Curran Associates, Inc.
- Geisler, W. S. (2008). Visual perception and the statistical properties of natural scenes. *Annual Review of Psychology*, *59*, 167–192.
- Gerhard, H. E., Wichmann, F. A., & Bethge, M. (2013). How sensitive is the human visual system to the local statistics of natural images? *PLoS Computational Biology*, *9*(1), e1002873, <https://doi.org/10.1371/journal.pcbi.1002873>.
- Glorot, X., Bordes, A., & Bengio, Y. (2011). Deep

- sparse rectifier neural networks. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics. Journal of Machine Learning and Research, Vol. 15* (pp. 315–323). Fort Lauderdale, FL: PMLR.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., . . . Bengio, Y. (2014). Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems 27* (pp. 315–323). Red Hook, NY: Curran Associates.
- Green, D. M., & Swets, J. A. (1966). *Signal detection theory and psychophysics*. New York, NY: Wiley.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. (2017). *Improved training of Wasserstein GANs*. <https://arXiv:1704.00028>.
- He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on image classification. *Proceedings of the 2015 IEEE International Conference on Computer Vision* (pp. 1026–1034). Washington, DC: IEEE Computer Society.
- Hjelm, R. D., Jacob, A. P., Che, T., Trischler, A., Cho, K., & Bengio, Y. (2018). Boundary-seeking generative adversarial networks. *International Conference on Learning Representations*. <https://arXiv:1702.08431v4>.
- Huth, A. G., Nishimoto, S., Vu, A. T., & Gallant, J. L. (2012). A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron, 76*, 1210–1224.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In F. Bach & D. Blei (Eds.), *Proceedings of the 32nd International Conference on Machine Learning, Journal of Machine Learning Research: Vol. 37* (pp. 448–456).
- Kayser, C., Nielsen, K. J., & Logothetis, N. K. (2006). Fixations in natural scenes: Interaction of image structure and image content. *Vision Research, 46*, 2535–2545.
- Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Computational Biology, 10*(11), e1003915.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations*. <https://arXiv:1412.6980v9>.
- Krizhevsky, A. (2009). *Learning multiple layers of features from tiny images* [Technical Report]. Toronto, Canada: University of Toronto.
- McDonald, J. S., & Tadmor, Y. (2006). The perceived contrast of texture patches embedded in natural images. *Vision Research, 46*, 3098–3104.
- Miyato, T., Kataoka, T., Koyama, M., & Yoshida, Y. (2018). Spectral normalization for generative adversarial networks. *International Conference on Learning Representations*. <https://arXiv:1802.05957v1>.
- Portilla, J., & Simoncelli, E. (2000). A parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision, 40*(1), 49–71.
- Radford, A., Luke, M., & Chintala, S. (2016). Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*. <https://arXiv:1511.06434v2>.
- Roth, K., Lucchi, A., Nowozin, S., & Hofmann, T. (2017). Stabilizing training of generative adversarial networks through regularization. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett, *Advances in Neural Information Processing Systems 30* (pp. 2018–2028). Red Hook, NY: Curran Associates.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., . . . Fei-Fei, L. (2015). ImageNet large scale visual recognition challenge. *International Journal of Computer Vision, 115*(3), 211–252.
- Scharr, H. (2000). *Optimale operatoren in der digitalen bildverarbeitung* (Unpublished doctoral dissertation). IWR, Fakultät für Physik und Astronomie, University of Heidelberg, Heidelberg, Germany.
- Schütt, H., Harmeling, S., Macke, J. H., & Wichmann, F. A. (2016). Painfree and accurate Bayesian estimation of psychometric functions for (potentially) overdispersed data. *Vision Research, 122*, 105–123.
- Scott, D. W. (1992). *Multivariate density estimation: Theory, practice, and visualization*. New York: Wiley.
- Sebastian, S., Abrams, J., & Geisler, W. S. (2017). Constrained sampling experiments reveal principles of detection in natural scenes. *Proceedings of the National Academy of Sciences, USA, 114*, E5731–E5740.
- Simoncelli, E. P., & Olshausen, B. A. (2001). Natural image statistics and neural representation. *Annual Review of Neuroscience, 24*, 1193–1216.
- 't Hart, B. M., Schmidt, H. C. E. F., Roth, C., & Einhauser, W. (2013). Fixations on objects in natural scenes: Dissociating importance from saliency. *Frontiers in Psychology, 4*:455, 1–9.

- Thorpe, S., Fize, D., & Marlot, C. (2001). Speed of processing in the human visual system. *Nature*, *381*, 520–522.
- van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N. . . . Yu, T. (2014). scikit-image: Image processing in Python. *PeerJ*, *2*, e453.
- Wallis, T. S. A., Bethge, M., & Wichmann, F. A. (2016). Testing models of peripheral encoding using metamorphism in an oddity paradigm. *Journal of Vision*, *16*(2): 4, 1–30. <https://doi.org/10.1157/16.2.4>.
- Wallis, T. S. A., & Bex, P. J. (2012). Image correlates of crowding in natural scenes. *Journal of Vision*, *12*(7): 6, 1–19. <https://doi.org/10.1157/12.7.6>.
- Wallis, T. S. A., Funke, C. M., Ecker, A. S., Gatys, L. A., Wichmann, F. A., & Bethge, M. (2017). A parametric texture model based on deep convolutional features closely matches texture appearance for humans. *Journal of Vision*, *17*(12):5, 1–29. <https://doi.org/10.1157/17.12.5>.
- Wichmann, F. A., Braun, D. I., & Gegenfurtner, K. R. (2006). Phase noise and the classification of natural images. *Vision Research*, *46*, 1520–1529.
- Wichmann, F. A., Drewes, J., Rosas, P., & Gegenfurtner, K. R. (2010). Animal detection in natural scenes: Critical features revisited. *Journal of Vision*, *10*(4):6, 1–27. <https://doi.org/10.1157/10.4.6>.
- Zhu, J.-Y., Krähenbühl, P., Shechtman, E., & Efros, A. A. (2016). Generative visual manipulation on the natural image manifold. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.) *Computer Vision—ECCV 2016. Lecture Notes in Computer Science, Vol. 9909*. New York, NY: Springer.

## Appendix

### Distribution of segmentation distances

Although the distributions for most image features analyzed in Experiment 1 are relatively intuitive, this might not be the case for the segmentation distances. Figure A1 shows the distribution of segmentation distances  $d_{segm}$  across all trials.

### Logistic regression analysis of image features

We also performed the analysis of image features using logistic regression. Specifically, we used logistic regression to predict the trial-by-trial responses from the respective distance measure. We then used

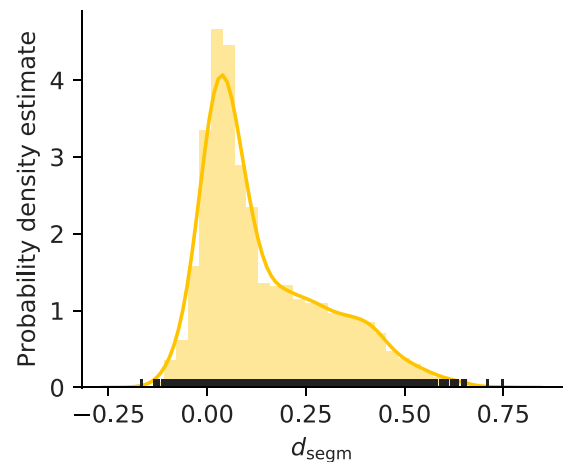


Figure A1. Distribution of segmentation distances across all trials. The yellow line is a kernel density estimate of the distribution, the bars show a histogram of the distribution. Kernel band width has been determined by Scott's rule (Scott, 1992), bin size has been determined by Freedman-Diaconis rule (Freedman & Diaconis, 1981). Dark gray ticks on the abscissa mark the locations of individual data points.

deviance to quantify how well the respective logistic regression model (and thus, the corresponding image feature) explained the observer's responses. Deviance is a generalization of the sum of squares error to the setting of logistic regression. Under the null hypothesis that the residuals are simply due to random fluctuations, deviance has a chi-square distribution with degrees of freedom equal to the difference between the number of data points and number of parameters in the model.

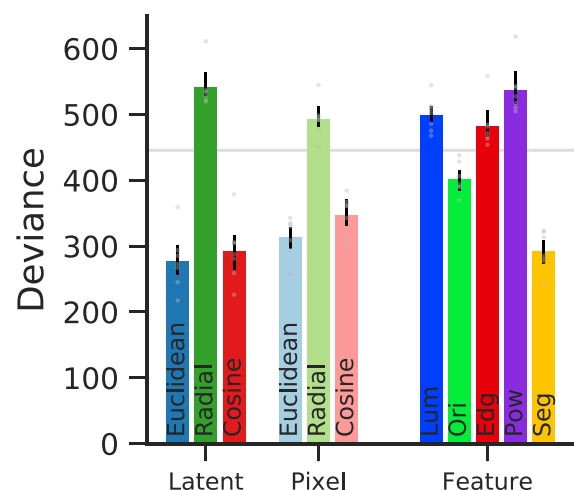


Figure A2. Logistic regression analysis of trial-by-trial errors. Colored bars indicate mean deviance across observers, error bars are 95% confidence intervals based on bootstrap across observers. Single observer deviances are marked by light gray dots. The horizontal gray line marks the 95th percentile of the distribution of deviances expected if residuals were due to chance.



Figure A2 shows deviances for the different image features. Overall, the results were similar to the analysis based on ROC curves: radial distances did not predict trial by trial errors very well, while Euclidean distance and cosine distance in latent space did. There was a significant difference between latent and pixel space for Cosine distance,  $t(6) = -3.86$ ,  $p <$

$0.05$  corrected, but not for Euclidean distance,  $t(6) = -2.46$ ,  $p = 0.048$  uncorrected.

Results for different image features were also very similar, with segmentation structure being most predictive. Low level features were generally not much better than chance performance.